

Sapienza University of Rome

Master in Artificial Intelligence and Robotics

# Machine Learning

A.Y. 2025/2026

Prof. Luca Iocchi

## 4. Probability and Bayes Networks

Luca Iocchi

# Outline

- Uncertainty
- Probability
- Syntax and Semantics
- Inference
- Independence and Bayes' Rule

## Uncertainty

Consider action  $A_t = \text{leave for airport } t \text{ minutes before flight.}$

Will  $A_t$  get me there on time?

Problems:

- partial observability (road state, other drivers' plans, etc.)
- noisy sensors (traffic reports)
- uncertainty in action outcomes (flat tire, etc.)
- complexity of modelling and predicting traffic

# Uncertainty

Hence a purely logical approach either

- risks falsehood: “ $A_{25}$  will get me there on time”
- leads to conclusions that are too weak for decision making: “ $A_{25}$  will get me there on time if there’s no accident on the bridge and it doesn’t rain and my tires remain intact etc etc.”
- leads to non-optimal decisions ( $A_{1440}$  might reasonably be said to get me there on time, but I’d have to stay overnight in the airport ...)

# Probability

Representation of uncertainty with probabilities.

Given the available evidence,  $A_{25}$  will get me at the airport on time with probability 0.04

Given the available evidence,  $A_{60}$  will get me at the airport on time with probability 0.85

Given the available evidence,  $A_{1440}$  will get me at the airport on time with probability 0.999

# Probability

## Sample space

- $\Omega$  *sample space* (set of possibilities)
- $\omega \in \Omega$  is a *sample point/possible world/atomic event/outcome of a random process/...*

## Probability space (or probability model)

- Function  $P : \Omega \mapsto \mathbb{R}$ , such that
  - $0 \leq P(\omega) \leq 1$
  - $\sum_{\omega \in \Omega} P(\omega) = 1$

Example: rolling a die

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$P(\omega) = \{1/6, 1/6, 1/6, 1/6, 1/6, 1/6\}$$

## Event

An *event*  $A$  is any subset of  $\Omega$

Probability of an event  $A$  is a function assigning to  $A$  a value in  $[0, 1]$

$$P(A) = \sum_{\{\omega \in A\}} P(\omega)$$

Example 1:  $A_1 = \text{"die roll"} < 4$ ,  $A_1 = \{1, 2, 3\} \subset \Omega$

$$P(A_1) = P(1) + P(2) + P(3) = 1/6 + 1/6 + 1/6 = 1/2$$

Example 2:  $A_2 = \text{"die roll"} = 4$ ,  $A_2 = \{4\}$ ,  $P(A_2) = 1/6$

Example 3:  $A_3 = \text{"die roll"} > 6$ ,  $A_3 = \emptyset$ ,  $P(A_3) = 0$

Example 4:  $A_4 = \text{"die roll"} \leq 6$ ,  $A_4 = \Omega$ ,  $P(A_4) = 1$

## Random variables

A *random variable* (outcome of a random phenomenon) is a function from the sample space  $\Omega$  to some range (e.g., the reals or Booleans)  $X : \Omega \mapsto B$ .

Example:  $Odd : \Omega \mapsto Boolean$ .

$X$  is a variable and a function !

$X = x_i$  : the random variable  $X$  has the value  $x_i \in B$

$X = x_i$  is equivalent to  $\{\omega \in \Omega | X(\omega) = x_i\}$

Example:  $Odd = true \equiv \{1, 3, 5\}$

## Random variables

$P$  induces a *probability distribution* for a random variable  $X$ :

$$P(X = x_i) = \sum_{\{\omega \in \Omega | X(\omega) = x_i\}} P(\omega)$$

Example

$$P(Odd = true) = P(1) + P(3) + P(5) = 1/6 + 1/6 + 1/6 = 1/2$$

# Propositions

A proposition is the event (subset of  $\Omega$ ) where an assignment to a random variable holds.

- event  $a \equiv A = \text{true} \equiv \{\omega \in \Omega | A(\omega) = \text{true}\}$

Propositions can be combined using standard logical operators, e.g.:

- event  $\neg a \equiv A = \text{false} \equiv \{\omega \in \Omega | A(\omega) = \text{false}\}$
- event  $a \wedge b =$  points  $\omega$  where  $A(\omega) = \text{true}$  and  $B(\omega) = \text{true}$
- event  $\neg a \vee b =$  points  $\omega$  where  $A(\omega) = \text{false}$  or  $B(\omega) = \text{true}$

$$P(\neg a \vee b) = \sum_{\{\omega \in \Omega | A(\omega) = \text{false} \vee B(\omega) = \text{true}\}} P(\omega)$$

## Syntax for propositions

- *Propositional* or *Boolean* random variables  
e.g., *PlayTennis*  
*PlayTennis = true* is a proposition, also written in lowercase *playtennis*
- *Discrete* random variables (*finite* or *infinite*)  
e.g., *Weather* is one of  $\langle \text{sunny}, \text{rain}, \text{cloudy}, \text{snow} \rangle$ .  
*Weather = rain* is a proposition  
Values must be exhaustive and mutually exclusive
- *Continuous* random variables (*bounded* or *unbounded*)  
e.g., *Temp = 21.6*, *Temp < 22.0*.
- Arbitrary Boolean combinations of basic propositions  
e.g., *playtennis*  $\wedge$  *Weather = rain*  $\wedge$  *Temp < 22.0*.

## Prior Probability

*Prior or unconditional probabilities* of propositions correspond to belief prior to arrival of any (new) evidence.

Examples:

$$P(\text{Odd} = \text{true}) = P(\text{odd}) = 0.5$$

$$P(\text{PlayTennis} = \text{true}) = P(\text{playtennis}) = 0.67$$

$$P(\text{Weather} = \text{sunny}) = 0.72$$

## Probability distribution

A *probability distribution* is a function assigning a probability value to all possible assignments of a random variable.

Note: sum of all values must be 1.

Examples:

$$P(\text{Odd}) = \langle 0.5, 0.5 \rangle$$

$$P(\text{PlayTennis}) = \langle 0.67, 0.33 \rangle$$

$$P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$$

Note: for real valued random variable  $X$ ,  $P(X)$  is a continuous function.

## Joint probability distribution

*Joint probability distribution* for a set of random variables gives the probability of every atomic joint event on those random variables (i.e., every sample point in the joint sample space).

Joint probability distribution of the random variables *Weather* and *PlayTennis*:

$P(\text{Weather}, \text{PlayTennis}) =$  a  $2 \times 4$  matrix of values:

<i>Weather</i> =	<i>sunny</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
<i>PlayTennis</i> = <i>true</i>	0.576	0.02	0.064	0.01
<i>PlayTennis</i> = <i>false</i>	0.144	0.08	0.016	0.09

*Every question about a domain can be answered by the joint distribution because every event is a set of sample points*

## Conditional/Posterior Probability

Belief after the arrival of some evidence.

I know the outcome of a random variable, how does this affect probability of other random variables?

Example:

I know that today *Weather* = *sunny*, how this information affects the random variable *PlayTennis*?

Notation:

$P(\text{PlayTennis} = \text{true} | \text{Weather} = \text{sunny})$ : conditional/posterior probability



## Conditional/Posterior Probability

In general, conditional/posterior probabilities are different from joint probabilities and from prior probabilities.

$$\begin{aligned} P(\text{PlayTennis} = \text{true} | \text{Weather} = \text{sunny}) &\neq \\ P(\text{PlayTennis} = \text{true}, \text{Weather} = \text{sunny}) &\neq \\ P(\text{PlayTennis} = \text{true}) \end{aligned}$$

Example (with  $pt \equiv \text{PlayTennis} = \text{true}$  and  $\text{sunny} \equiv \text{Weather} = \text{sunny}$ ):

$P(pt) = 0.67$ : prior

$P(pt, \text{sunny}) = 0.576$ : joint

$P(pt | \text{sunny}) = 0.8$ : posterior

## Conditional Probability Distributions

Conditional probability distributions: representation of all the values of conditional probabilities of random variables.

$P(\text{PlayTennis} | \text{Weather}) = 2 \times 4$  matrix

$P(\text{PlayTennis}   \text{Weather})$	<i>sunny</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
$\text{PlayTennis} = \text{true}$	0.8	0.2	0.8	0.1
$\text{PlayTennis} = \text{false}$	0.2	0.8	0.2	0.9

Different from joint probability matrix

## Conditional probability

Definition of conditional probability:

$$P(a|b) \equiv \frac{P(a \wedge b)}{P(b)} \text{ if } P(b) \neq 0$$

*Product rule*

$$P(a, b) \equiv P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$$

Example:

$$P(pt, sunny) = P(pt|sunny)P(sunny) = 0.8 \cdot 0.72 = 0.576$$

A general version holds for whole distributions, e.g.,

$$P(PlayTennis, Weather) = P(PlayTennis|Weather)P(Weather)$$

## Total probabilities

For a Boolean random variable  $B$

$$P(a) = P(a|b)P(b) + P(a|\neg b)P(\neg b)$$

In general, for a random variable  $Y$  accepting mutually exclusive values  $y_i$

$$P(X) = \sum_{y_i \in \mathcal{D}(Y)} P(X|Y = y_i)P(Y = y_i)$$

$\mathcal{D}(Y)$ : set of values for variable  $Y$

Example:

$$\begin{aligned} P(pt) &= P(pt|sunny)P(sunny) + P(pt|rain)P(rain) + \dots = \\ &= 0.8 \cdot 0.72 + 0.2 \cdot 0.1 + 0.8 \cdot 0.08 + 0.1 \cdot 0.1 = \\ &= \langle 0.8, 0.2, 0.8, 0.1 \rangle \cdot \langle 0.72, 0.1, 0.08, 0.1 \rangle = 0.67 \end{aligned}$$

## Chain Rule

- *Chain rule* is derived by successive application of product rule:

$$P(X_1, X_2) = P(X_1)P(X_2|X_1)$$

$$\begin{aligned} P(X_1, \dots, X_n) &= P(X_1, \dots, X_{n-1})P(X_n|X_1, \dots, X_{n-1}) \\ &= P(X_1, \dots, X_{n-2})P(X_{n-1}|X_1, \dots, X_{n-2})P(X_n|X_1, \dots, X_{n-1}) \\ &= \dots \\ &= \prod_{i=1}^n P(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$

## Inference by Enumeration

Start with the joint distribution  $P(\omega) = P(X_1, \dots, X_n)$

For any proposition  $\phi$ , sum the atomic events where it is true:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$$

Example:  $P(X_1, X_2, X_3)$  ( $X_i$  Boolean),  $\phi = x_1 \vee x_2$

$$P(\phi) = P(x_1, x_2, x_3) + P(x_1, x_2, \neg x_3) + P(\neg x_1, x_2, x_3) + \dots + P(\neg x_1, x_2, \neg x_3)$$

## Normalization

$$P(\phi | \psi) = \frac{P(\phi \wedge \psi)}{P(\psi)}$$

Denominator  $P(\psi)$  can be seen as a *normalization factor*  $\alpha$  constant wrt  $\phi$ .

Example:

$$\begin{aligned} P(\text{PlayTennis} | \text{sunny}) &= \alpha P(\text{PlayTennis}, \text{sunny}) = \\ &\alpha \langle 0.576, 0.144 \rangle = \langle 0.8, 0.2 \rangle, \text{ (with } \alpha = 0.72) \end{aligned}$$

$\alpha$  does not depend on *PlayTennis*

Note:  $\text{argmax}$  operator not affected by normalization factors

$$\text{argmax}_{x_i \in \mathcal{D}(X)} P(X = x_i | Y) = \text{argmax}_{x_i \in \mathcal{D}(X)} P(X = x_i, Y)$$

## Independence

$A$  and  $B$  are *independent* iff

$$P(A|B) = P(A) \quad \text{or} \quad P(B|A) = P(B) \quad \text{or} \quad P(A, B) = P(A)P(B)$$

Example:

$$P(\text{PlayTennis}, \text{Weather}, \text{Odd}) = P(\text{PlayTennis}, \text{Weather})P(\text{Odd})$$

$$P(\text{pt}, \text{sunny}, \text{odd}) = P(\text{pt}, \text{sunny})P(\text{odd}) = 0.576 \cdot 0.5 = 0.288$$

# Independence

$P(\text{PlayTennis}, \text{Weather}, \text{Odd})$  has 16 entries

$P(\text{PlayTennis}, \text{Weather})$  and  $P(\text{Odd})$  have  $8 + 2 = 10$  entries

Example:  $n$  independent biased coins, reduced size from  $2^n$  to  $n$

Absolute independence is powerful, but rare.

Complex systems have hundreds of variables, none of which are independent.

## Conditional independence

General formulation:

$X$  conditionally independent from  $Y$  given  $Z$  iff  $P(X|Y, Z) = P(X|Z)$

$$P(X, Y|Z) = P(X|Y, Z)P(Y|Z) = P(X|Z)P(Y|Z)$$

$$P(Y_1, \dots, Y_n|Z) = P(Y_1|Y_2, \dots, Y_n, Z)P(Y_2|Y_3, \dots, Y_n, Z) \cdots P(Y_n|Z)$$

$Y_i$  conditionally independent from  $Y_j$  given  $Z$

$$P(Y_1, \dots, Y_n|Z) = P(Y_1|Z)P(Y_2|Z) \cdots P(Y_n|Z)$$

## Conditional independence

Chain rule + Conditional independence

$$\begin{aligned} P(X, Y, Z) &= P(X|Y, Z)P(Y, Z) = P(X|Y, Z)P(Y|Z)P(Z) \\ &= P(X|Z)P(Y|Z)P(Z) \end{aligned}$$

$$P(Y_1, Y_2, \dots, Y_n, Z) = P(Y_1|Z)P(Y_2|Z) \dots P(Y_n|Z)P(Z)$$

In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in  $n$  to linear in  $n$ .

Example:  $Y_i, Z$  boolean variables: joint space:  $2^{n+1}$ , conditional independence space =  $4n + 2$ .

## Bayes' Rule

- Product rule  $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

$$\Rightarrow \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \alpha P(X|Y)P(Y)$$

Useful for assessing *diagnostic* probability from *causal* probability:

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

## Maximum likelihood problem

Given  $Y$ , what is the most likely value of  $X$  ?

$$\begin{aligned} \operatorname{argmax}_{x_i \in \mathcal{D}(X)} P(X = x_i | Y) &= \operatorname{argmax}_{x_i \in \mathcal{D}(X)} \frac{P(Y | X = x_i) P(X = x_i)}{P(Y)} \\ &= \operatorname{argmax}_{x_i \in \mathcal{D}(X)} P(Y | X = x_i) P(X = x_i) \end{aligned}$$

$P(Y)$  is constant wrt  $X$  and can be eliminated

Example:

$$\begin{aligned} \operatorname{argmax}_{v_i \in \{y, n\}} P(\text{PlayTennis} = v_i | \text{sunny}) &= \\ \operatorname{argmax}_{v_i \in \{y, n\}} P(\text{sunny} | \text{PlayTennis} = v_i) P(\text{PlayTennis} = v_i) \end{aligned}$$

## Bayes' Rule and conditional independence

Bayes rule

$$P(Z | Y_1, \dots, Y_n) = \alpha P(Y_1, \dots, Y_n | Z) P(Z)$$

$Y_1, \dots, Y_n$  conditionally independent each other given  $Z$

$$P(Z | Y_1, \dots, Y_n) = \alpha P(Y_1 | Z) \cdots P(Y_n | Z) P(Z)$$

Effects conditionally independent each other given a cause.

$$P(\text{Cause} | \text{Effect}_1, \dots, \text{Effect}_n) = \alpha P(\text{Cause}) \prod_i P(\text{Effect}_i | \text{Cause})$$

Total number of parameters is *linear* in  $n$

## Maximum likelihood problem

Given  $Y_1, \dots, Y_n$ , what is the most likely value of  $X$  ?

$$\begin{aligned} \operatorname{argmax}_{x_i \in \mathcal{D}(X)} P(X = x_i | Y_1, \dots, Y_n) = \\ \operatorname{argmax}_{x_i \in \mathcal{D}(X)} P(Y_1 | X = x_i) \dots P(Y_n | X = x_i) P(X = x_i) \end{aligned}$$

Example:

$$\begin{aligned} \operatorname{argmax}_{v_i \in \{y, n\}} P(\text{PlayTennis} = v_i | \langle \text{sunny}, \dots \rangle) = \\ \operatorname{argmax}_{v_i \in \{y, n\}} P(\text{sunny} | \text{PlayTennis} = v_i) \dots P(\text{PlayTennis} = v_i) \end{aligned}$$

## Bayesian networks

Graphical notation for conditional independence assertions and hence for compact specification of full joint distributions.

Syntax:

- a set of nodes, one per variable
- a directed, acyclic graph (link  $\approx$  “directly influences”)
- a conditional distribution for each node given its parents:  
 $P(X_i | \text{Parents}(X_i))$

In the simplest case, conditional distribution represented as a *conditional probability table* (CPT) giving the distribution over  $X_i$  for each combination of parent values.



## Burglar BN Example

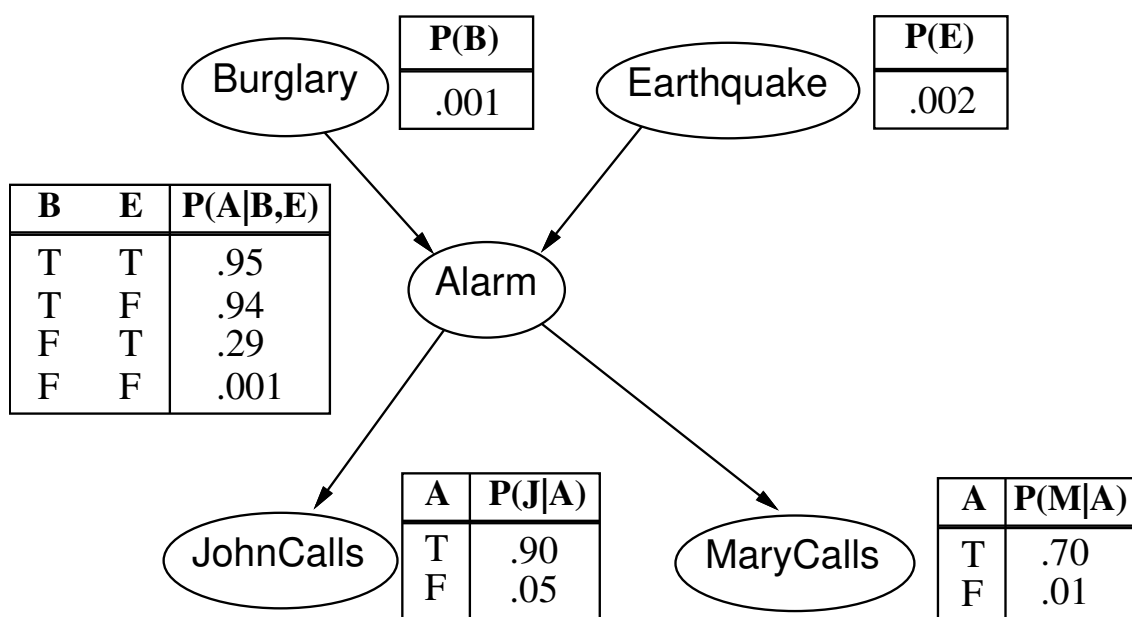
I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes the alarm is set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

Network topology reflects "causal" knowledge:

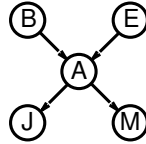
- A burglar can set the alarm
- An earthquake can set the alarm
- The alarm can cause Mary to call
- The alarm can cause John to call

## Burglar BN Example



## Compactness

A CPT for Boolean variable  $X_i$  with  $k$  Boolean parents has  $2^k$  rows for the combinations of parent values



Each row requires one number  $p$  for  $X_i = \text{true}$  (the number for  $X_i = \text{false}$  is just  $1 - p$ )

If each variable has no more than  $k$  parents, the complete network requires  $O(n \cdot 2^k)$  numbers

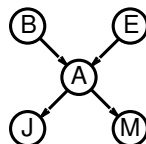
I.e., grows linearly with  $n$ , vs.  $O(2^n)$  for the full joint distribution

For burglary net,  $1 + 1 + 4 + 2 + 2 = 10$  numbers (vs.  $2^5 - 1 = 31$ )

## Computing joint probabilities

All joint probabilities computed with the chain rule:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i))$$



e.g.,  $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned} &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\ &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\ &\approx 0.00063 \end{aligned}$$