

# *Artificial Intelligence*

## Probabilistic reasoning: *representation*

Marco Piastra

# Probabilistic reasoning: *representation*

Possible worlds, events as subsets

Random variables: one value in each world

Probability

*Joint distribution*

*Marginalization*

*Conditionalization*

Independence, conditional independence

Graphical models

# Possibility (i.e. from logic to probability)

## ■ Objective knowledge and *plausible* knowledge

Each (rational) agent is supposed to hold some knowledge which is *objective* (to him/her)

Call  $\Gamma$  the theory representing such objective knowledge,

to the agent, the set of *possible worlds* is  $W \equiv \{ \langle U, v \rangle : \langle U, v \rangle \models \Gamma \}$

Example: the agent knows  $\Gamma \equiv \{ \varphi \vee \psi \}$

Therefore, only the worlds  $\{ \langle U, v \rangle : \langle U, v \rangle \models \{ \varphi \vee \psi \} \}$  are *possible* (to him/her)

Does this mean that the *event*  $\varphi \vee \psi$  did occur, already?

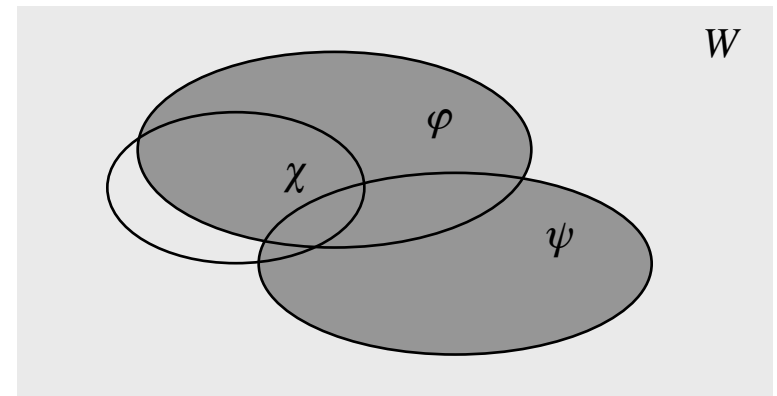
On the other hand, the agent might not know  $\chi$

Formally:

$$\varphi \vee \psi \not\models \chi$$

$$\varphi \vee \psi \not\models \neg\chi$$

To him/her, both  $\chi$  and  $\neg\chi$  are *plausible*  
(*plausible* = logically possible)



# Events as *subsets of possible worlds*

## ■ Subsets of *possible worlds*

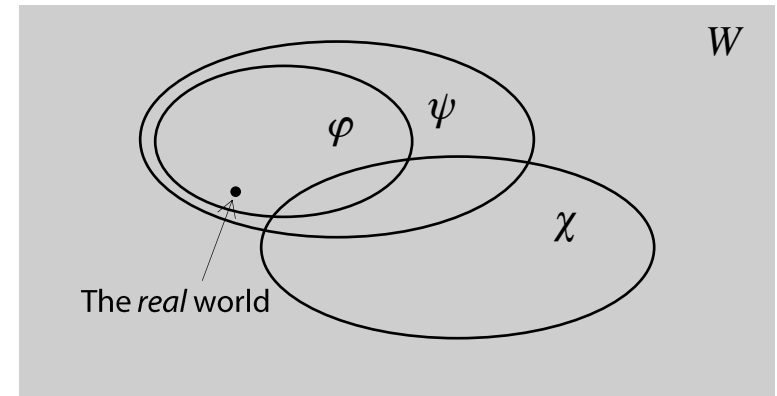
We already know that, given an interpreted logical language  $L$  (e.g. of the first order), every (closed) wff is associated to a *subset of possible worlds*

Formally, each (closed)  $\varphi$  is associated to the set  $\{\langle U, v \rangle : \langle U, v \rangle \models \varphi\}$   
(for simplicity, we keep  $U$  fixed here)

Intuition:

An **event** can be seen as a subset of *possible worlds*:  
an event is said to occur when the *real* world happens to belong to the corresponding subset of possible worlds

The agent is not supposed to know which world is the *real* one...



Note:

In 'classical' probability theory, the events need not be defined in a logical fashion (this fact has also some technical implications, to be clarified later on)

- Probability is a measure over the subsets of  $W$

In particular over  $W \equiv \{ \langle U, v \rangle : \langle U, v \rangle \models \Gamma \}$

i.e. the set of worlds that are deemed possible by an agent who knows  $\Gamma$

Technically,  $P(\cdot)$  is a *function* that assigns a measure (i.e. a real number) to each elements of a  $\sigma$ -algebra  $\Sigma$  of subsets of  $W$

## $\sigma$ -algebra (definition)

A collection of subsets  $\Sigma$  of a set  $W$  such that:

1)  $\Sigma$  is not empty

2) If  $\varphi \in \Sigma$  then  $\neg\varphi \in \Sigma$

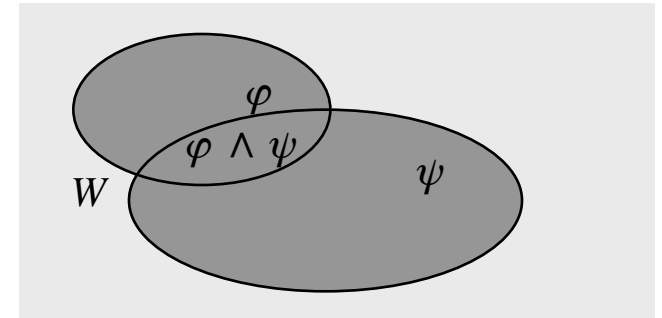
( $\neg\varphi$  is intended as the *complement* of  $\varphi$  in  $W$ )

3) For any *countable* collection of subsets  $\{\varphi_i\}$ ,  $\varphi_i \in \Sigma$ , we have  $\bigcup_i \varphi_i \in \Sigma$

Corollary:

The sets  $\emptyset$  e  $W$  belong to any  $\sigma$ -algebra generated on  $W$

Each element of a  $\sigma$ -algebra is an **event**



- Probability is a measure over the subsets of  $W$

In particular over  $W \equiv \{ \langle U, v \rangle : \langle U, v \rangle \models \Gamma \}$

i.e. the set of worlds that are deemed possible by an agent who knows  $\Gamma$

Technically,  $P(\cdot)$  is a *function* that assigns a measure (i.e. a real number) to each elements of a  $\sigma$ -algebra  $\Sigma$  of subsets of  $W$

$P(\cdot)$  is a *measure* defined over the  $\sigma$ -algebra  $\Sigma$

1) For each event  $\varphi \in \Sigma$ ,  $P(\varphi) \geq 0$

2)  $P(W) = 1$

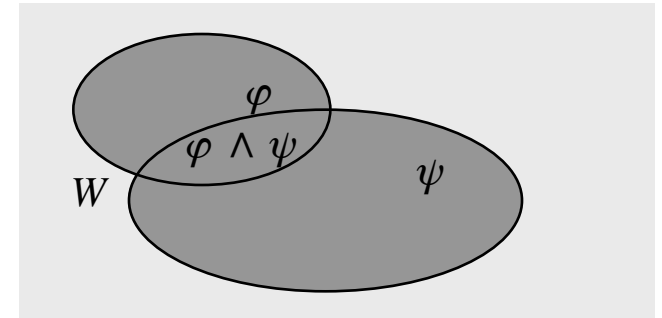
3) For every countable sequence  $\varphi_i$  of *disjoint* events

in  $\Sigma$  (*disjoint*  $\Leftrightarrow \varphi_i \cap \varphi_j \equiv \emptyset$  se  $i \neq j$ ):

$$P(\varphi_1 \vee \varphi_2 \vee \dots \vee \varphi_n) = \sum_i P(\varphi_i)$$

Corollary:

For any event  $\varphi \in \Sigma$ ,  $0 \leq P(\varphi) \leq 1$



# Partitions, random variables\*

## ■ Partition

A collection  $\varphi_i$  of *disjoint* events such that

$$\bigcup_i P(\varphi_i) = W$$

## ■ Random Variable

Let  $X$  be a variable having  $\{v_1, v_2, \dots, v_n\}$  as its domain.

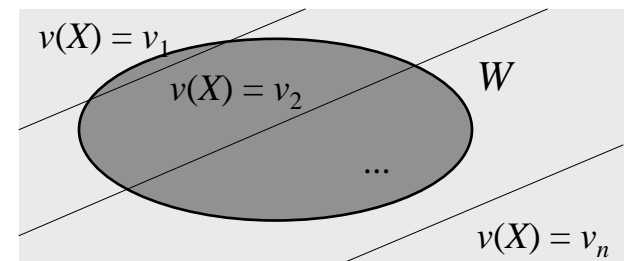
In each possible world,  $X$  has a specific value  $v_i$

The set of values  $v(X) = v_1, v(X) = v_2, \dots, v(X) = v_n$  define a *partition* of  $W$

- $X$  is a *random variable*
- Each constraint  $v(X) = v_i$  defines an event (i.e. a subset of  $W$ )
- *Given that*  $X=v_i$  e  $X=v_j$  are disjoint,  $P(X=v_i \vee X=v_j) = P(X=v_i) + P(X=v_j)$  whenever  $i \neq j$

Random variable having binary values are also said to be *bernoullian*

Random variables with vectorial values are also said to be *multinomial*



# Random variables, joint distribution\*

## ■ Many random variables

In practice, in a probabilistic representation, multiple random variables have to coexist

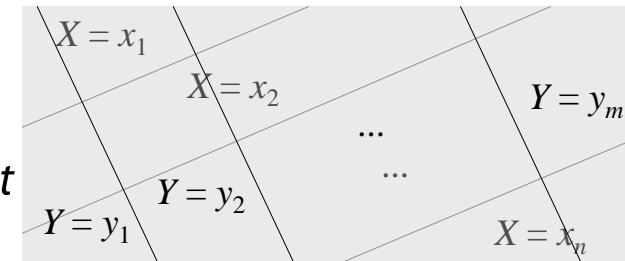
Example:

$X_i$  occurrence of a word  $I$  in the body of an email (0/1)

$Y$  classification of that email as spam (0/1)

Together, a collection of r.v.s define a partition of  $W$

Any combination of values of the above r.v.s defines an *event*



## ■ Joint probability distribution

for a given set of random variables, e.g.  $X, Y, Z$

It is a function  $P(X=x_i \wedge Y=y_j \wedge Z=z_k)$  that associates a real value to each individual combination of values  $\langle x_i, y_j, z_k \rangle$

Alternative notation:  $P(X=x_i, Y=y_j, Z=z_k)$  more frequently, just:  $P(X, Y, Z)$

Given that  $X, Y$  e  $Z$  define a partition of  $W$ : 
$$\sum_i \sum_j \sum_k P(X=x_i, Y=y_j, Z=z_k) = 1$$



# Marginalization

*Removing a random variable from a joint distribution*

Given a joint probability distribution

$$P(X=x_i, Y=y_j, Z=z_k)$$

The *marginal probability*  $P(X=x_i, Y=y_j)$  is obtained via summation:

$$P(X=x_i, Y=y_j) = \sum_k P(X=x_i, Y=y_j, Z=z_k)$$

A marginal probability, in general, is still a joint probability

# Conditional probability

## ■ Definition

$$P(A|B) = \frac{P(A,B)}{P(B)}$$

## ■ Meaning

It is a form of *inference*: switching from a set  $W$  to a set  $W'$  (i.e. a subset of the former)

*Therefore, from a probability measure to another one*

Consider an agent who thinks that  $W$  is the set of possible worlds

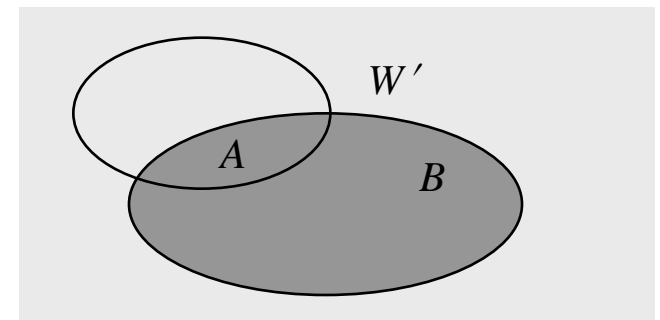
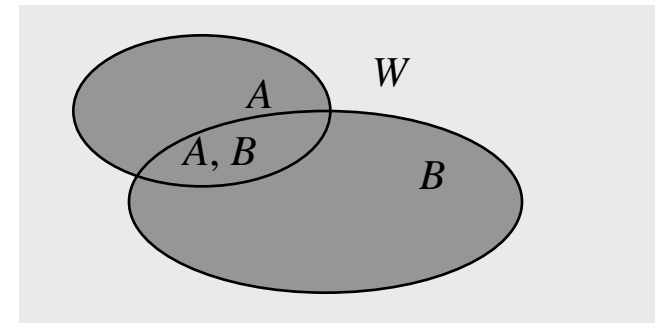
$P(A)$  is the probability (to him/her) that event  $A$  occurs

Suppose that the agent then learns that event  $B$  occurred

The event  $\neg B$  is now *impossible* (to him/her)

$W' \equiv B$  is the new set of possible worlds

$P(A | B)$  is the new probability of  $A$



# Bayes' Theorem (T. Bayes, 1764)



## ■ Definition

A relation between conditional and marginal probabilities

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}$$

$P(B | A)$  is also called *likelihood*  $L(A | B)$

$$P(A | B) \propto L(A | B) P(A)$$

The theorem follows from the definition of conditional probability (*chain rule*)

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

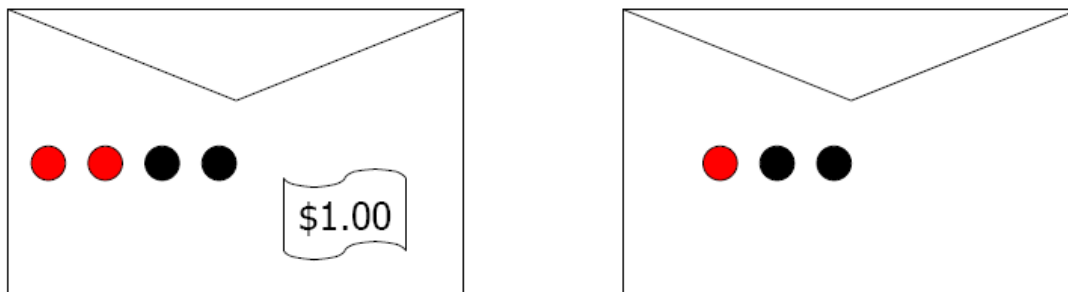
Given the definition of marginalization:

$$P(B) = \sum_A P(A, B) = \sum_A P(B | A) P(A)$$

also follows (Bayes' theorem alternative formulation):

$$P(A | B) = \frac{P(B | A) P(A)}{\sum_A P(B | A) P(A)}$$

# Example: information and bets



- Two envelopes, only one is extracted

One envelope contains two red tokens and two black tokens, it is worth \$1.00

One envelope contains one red token and two black tokens, it is valueless

The envelope has been extracted.

Before posing you bet, you are allowed to extract one token from it

a) The token is black. How much do you bet ?

b) The token is red. How much do you bet ?

Purpose: showing that Bayes' Theorem makes the representation easier

# Independence, conditional independence

## ■ Independence (also *marginal independence*)

Two events are independent iff their joint probability is equal to the product of the marginals

$$\langle A \perp B \rangle \Rightarrow P(A, B) = P(A) P(B)$$

## ■ Conditional independence

Two events are conditional independent, given a third event, iff their joint conditional probability is equal to the product of the conditional marginals

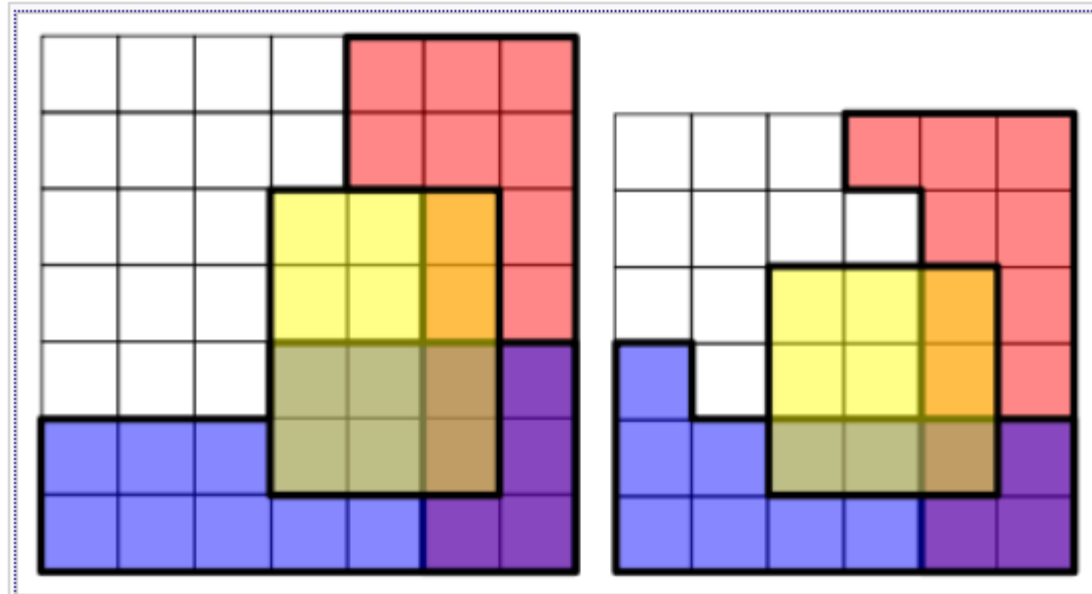
$$\langle A \perp B \mid C \rangle \Rightarrow P(A, B \mid C) = P(A \mid C) P(B \mid C)$$

$$\Rightarrow P(A \mid B, C) = \frac{P(A, B \mid C)}{P(B \mid C)} = \frac{P(A \mid C) P(B \mid C)}{P(B \mid C)} = P(A \mid C)$$

CAUTION: the two forms of independence are distinct!

$$\langle A \perp B \rangle \not\Rightarrow \langle A \perp B \mid C \rangle, \langle A \perp B \mid C \rangle \not\Rightarrow \langle A \perp B \rangle$$

# Independence, conditional independence



These are two examples illustrating **conditional independence**. Each cell represents a possible outcome. The events  $R$ ,  $B$  and  $Y$  are represented by the areas shaded red, blue and yellow respectively. And the probabilities of these events are shaded areas with respect to the total area. In both examples  $R$  and  $B$  are conditionally independent given  $Y$  because:

$$\Pr(R \cap B \mid Y) = \Pr(R \mid Y) \Pr(B \mid Y)^{[1]}$$

but not conditionally independent given not  $Y$  because:

$$\Pr(R \cap B \mid \text{not } Y) \neq \Pr(R \mid \text{not } Y) \Pr(B \mid \text{not } Y).$$

[from Wikipedia, "Conditional Independence"]

# Graphical models (*Bayesian Networks*)

Structure and numbers, instead of just numbers

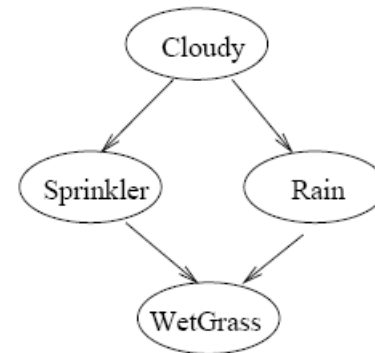
- A structured representation of a joint probability

Each model is an *oriented* graph

The nodes are random variables

The arcs represent *dependence*

C	P(S=F)	P(S=T)
F	0.5	0.5
T	0.9	0.1



C	P(R=F)	P(R=T)
F	0.8	0.2
T	0.2	0.8

Note that a complete specification of a joint probability would require  $2^4 = 16$  values

The values in figure are just 9

S	R	P(W=F)	P(W=T)
F	F	1.0	0.0
T	F	0.1	0.9
F	T	0.1	0.9
T	T	0.01	0.99

# From graphical models to joint probability

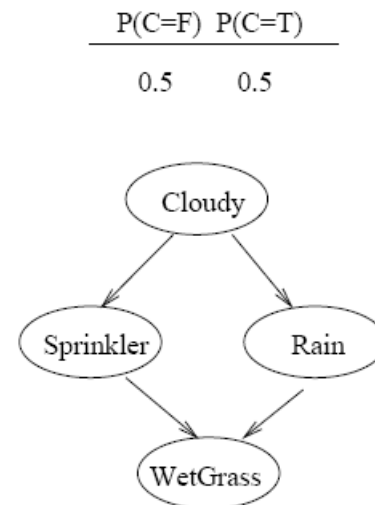
## Joint probability

It can be expressed as a product of conditional probabilities

(due to the extension of the *chain rule*)

Example:

C	P(S=F)	P(S=T)
F	0.5	0.5
T	0.9	0.1



C	P(R=F)	P(R=T)
F	0.8	0.2
T	0.2	0.8

$$P(C, S, R, W) = P(C)P(S | C)P(R | S, C)P(W | R, S, C)$$

In a graphical model, the joint distribution is

$$P(X_1, X_2, \dots, X_n) = \prod_i P(X_i | \text{parents}(X_i))$$

Where  $\text{parents}(X_i)$  the nodes from which there is an entry arc to  $X_i$

In the example:

$$P(C, S, R, W) = P(C)P(S | C)P(R | C)P(W | R, S)$$

Conditional independence assumptions:  $\langle R \perp S | C \rangle, \langle W \perp C | R, S \rangle$

S	R	P(W=F)	P(W=T)
F	F	1.0	0.0
T	F	0.1	0.9
F	T	0.1	0.9
T	T	0.01	0.99



# Graphical models and conditional independence

- *D-separation (Dependency-separation)*

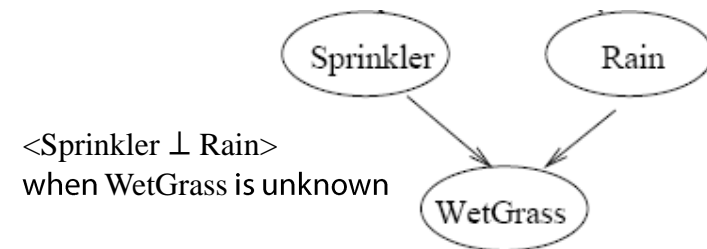
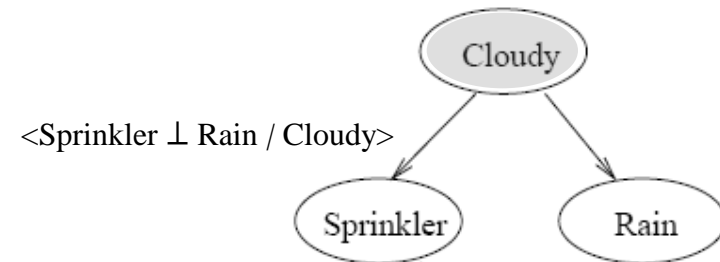
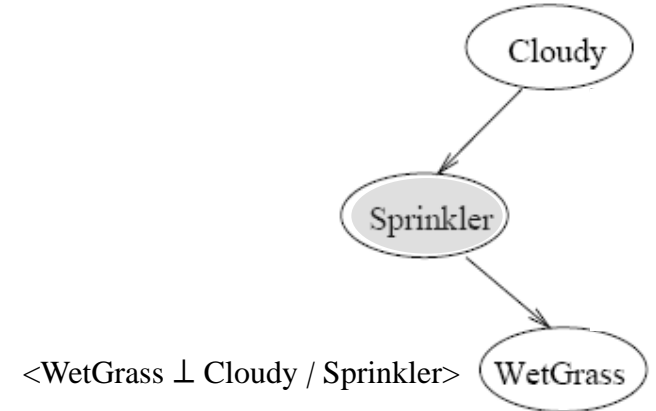
i.e. how to read a graphical model

In a graphical model

Two nodes  $X$  and  $Y$  are conditional independent given a set of nodes  $\{Z_k\}$  when all paths are blocked (see below)

A path between  $X$  e  $Y$  is blocked if:

- 1) It is either a sequence  $X \rightarrow \dots Z_i \dots \rightarrow Y$  or a *fork*  $X \leftarrow \dots Z_i \dots \rightarrow Y$  ( $Z_i \in \{Z_k\}$ )
- 2) It is a *join*  $X \rightarrow \dots N \dots \leftarrow Y$  where neither  $N$  nor all the *descendants* of  $N$  belong to  $\{Z_k\}$



# Explaining Away

A few more words on condition 2) of *D-separation*

## Graphical model, with a *join*

Joint probability, from the graph:

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

Marginal probability w.r.t  $X$  and  $Y$  ( $Z$  unknown):

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

Therefore  $X$  e  $Y$  are *marginally independent*

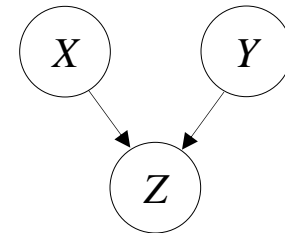
But when  $Z$  is known, then  $X$  and  $Y$  are *dependent*:

$$P(X, Y | Z=v) = \frac{P(X, Y, Z=v)}{P(Z=v)} = \frac{P(X)P(Y)P(Z=v|X, Y)}{\sum_{X, Y} P(X)P(Y)P(Z=v|X, Y)}$$

It is not a paradox.

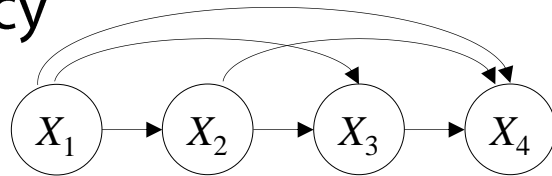
Example:

$X$  and  $Y$  are two tosses of the same coin,  $Z=1$  if the result is the same,  $Z=0$  otherwise.



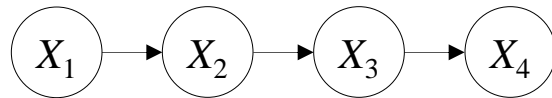
# Example of graphical models

## ■ Complete dependency



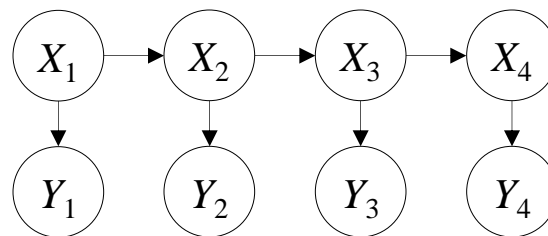
$$P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2)P(X_4 | X_1, X_2, X_3)$$

## ■ Markovian model



$$P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2 | X_1)P(X_3 | X_2)P(X_4 | X_3) = P(X_1) \prod_{i=2}^n P(X_i | X_{i-1})$$

## ■ 'Hidden' Markovian model



Typically, nodes  $X_i$  are *hidden*, in the sense of *non-observable*

$$\begin{aligned} P(X_1, X_2, X_3, X_4, Y_1, Y_2, Y_3, Y_4) &= P(X_1)P(Y_1 | X_1)P(X_2 | X_1)P(Y_2 | X_2)P(X_3 | X_2)P(Y_3 | X_3)P(X_4 | X_3)P(Y_4 | X_4) \\ &= P(X_1)P(Y_1 | X_1) \prod_{i=2}^n P(X_i | X_{i-1})P(Y_i | X_i) \end{aligned}$$