

Artificial Intelligence²

Roman Barták

Department of Theoretical Computer Science and Mathematical Logic

Introduction

We construct rational agents.

An **agent** is an entity that perceives its **environment** through **sensors** and acts upon that environment through **actuators**.

A **rational agent** is an agent maximizing its expected performance measure.



In AI 1 we dealt mainly with a logical approach to agent design (no uncertainty).

We ignored

- interface to environment (sensors, actuators)
- uncertainty
- the possibility of self-improvement (learning)

Introduction

- motivation and background on probability

Probabilistic reasoning

- uncertainty, probabilistic reasoning, Bayesian networks, Hidden Markov Models

Rational decisions

- utility theory, Markov Decision Processes, game theory, mechanism design

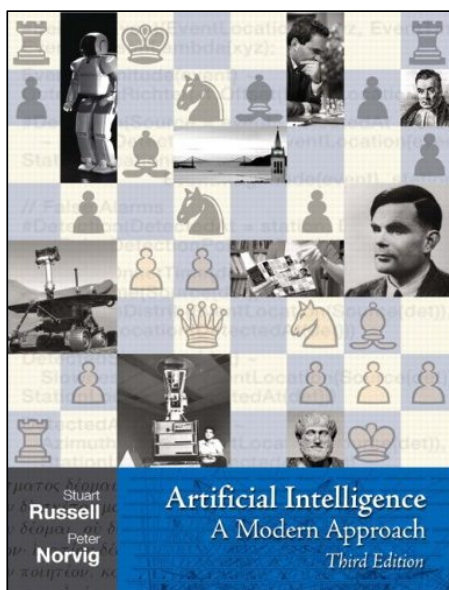
Machine learning

- decision trees, regression, SVM, reinforcement learning



Resources

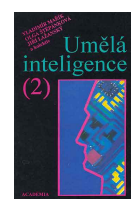
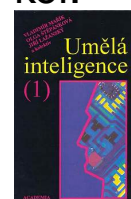
Artificial Intelligence: A Modern Approach



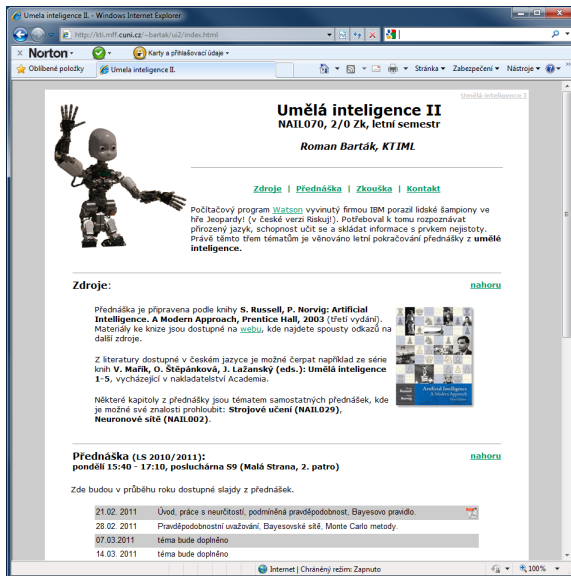
- S. Russell and P. Norvig
- Prentice Hall, 2010 (3rd ed.)
- <http://aima.cs.berkeley.edu/>

Umělá inteligencia 1-6

- Vladimír Mařík, Olga Štěpánková, Jiří Lažanský a kol.
- Academia



<http://ktiml.mff.cuni.cz/~bartak/ui2>



You can find there:

- slides
- links and resources
- contacts
- quiz
- ...

Seminar on Artificial Intelligence II

- how to apply AI techniques in practice

Machine learning

- how can computers learn new things

Multi-agent systems

- how to handle multiple agents

Probabilistic graphical models

- how to do Bayesian inference efficiently etc.

Human-like artificial agents

- how to design agents for virtual environments

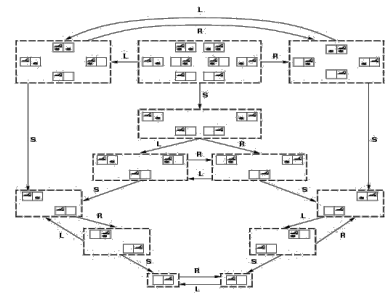
Practical course on robotics

- how to design hardware agents

Can we handle uncertain information in the pure logical approach?

belief states

- represent sets of all possible world states for the agent



Drawbacks

- a logical agent must consider every logically possible explanation for the observations, no matter how unlikely (large and complex representations)
- a correct contingent plan must consider arbitrary likely contingencies (big plans)
- sometimes there is no plan that is guaranteed to achieve the goal, yet the agent must act

Example

Diagnosing a dental patient's toothache

Let us try to apply propositional logic:

Toothache \Rightarrow Cavity

Hmm, is it really true?

- not all patients with toothaches have cavities; some of them have gum disease, an abscess, or other problems

Toothache \Rightarrow Cavity \vee GumProblem \vee Abscess \vee ...

We could try turning the rule into a causal rule:

Cavity \Rightarrow Toothache

But this is not right either – not all cavities cause pain

The only way to fix the rule is to make it logically exhaustive!



Why does logic fail to cope with a domain like medical diagnosis?



- **laziness:** it is too much work to list the complete set of antecedent or consequents and too hard to use such rules
- **theoretical ignorance:** medical science has no complete theory for the domain
- **practical ignorance:** even if we know all the rules we might be uncertain because not all the necessary tests have been or can be run

We need another tool to deal with degrees of belief – **probability theory.**

A logical agent believes each sentence to be true or false or has no opinion. A probabilistic agent may have a numerical degree of belief between 0 (certainly false) and 1 (certainly true).

Basic probability notation

Like logical assertions, probabilistic assertions are about possible worlds – **sample space** Ω .

– the possible worlds are **mutually exclusive** and **exhaustive**

Each possible world ω is associated with a numerical probability $P(\omega)$ such that:

$$0 \leq P(\omega) \leq 1$$

$$\sum_{\omega \in \Omega} P(\omega) = 1$$

Example:

If we are about to roll two (distinguishable) dice, there are 36 possible worlds to consider: (1,1), (1,2), ..., (6,6)

$$P(\omega) = 1/36$$



The sets of possible worlds are called **events**.

Example: „doubles are rolled“ is an event

The probability of event is the sum of probabilities of possible worlds in the event.

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

Example:

$$P(\text{doubles}) =$$

$$1/36 + 1/36 + 1/36 + 1/36 + 1/36 + 1/36 = 1/6$$

These probabilities are called **unconditional** or **prior probabilities** („priors“ for short).

Conditional probability

Frequently, we have some information (**evidence**) and we are interested in probability of some event.

For example, what is the probability of double if we already know that first die rolled to 5?

$$P(\text{doubles} \mid \text{Die}_1 = 5) = 1/36 / (6 \cdot 1/36) = 1/6$$

This is called **conditional** or **posterior probability**

$$P(a \mid b) = P(a \wedge b) / P(b), \text{ whenever } P(b) > 0$$

This can be also written in a different form called the **product rule**

$$P(a \wedge b) = P(a \mid b) \cdot P(b)$$

Beware! If we have more evidence then the conditional probability needs to assume it.

$$P(\text{doubles} \mid \text{Die}_1 = 5, \text{Die}_2 = 5) = 1$$

In a factored representation, a possible world is represented by a set of variable/value pairs.

Variables in probability theory are called **random variables**. Every random variable has a domain – the set of possible values it can take on (similarly to a CSP).

Die₁ – represents a value on the first die 1 (1,...,6)

Cavity – describes whether the patient has or has not cavity (true, false)

A possible world is fully identified by values of all random variables.

$P(\text{Die}_1 = 5, \text{Die}_2 = 5)$

Probability distribution

Probability of all possible worlds can be described using a table called a **full joint probability distribution** – the elements are indexed by values of random variables.

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

Given the table, we can calculate probabilities of values of any random variable:

$$P(\text{toothache}=\text{true}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$$

$$P(\text{toothache}=\text{false}) = 0.072 + 0.008 + 0.144 + 0.576 = 0.8$$

We will describe the table in a short way as:

$$\mathbf{P}(\text{Toothache}) = \langle 0.2, 0.8 \rangle$$

$$P(\neg a) = 1 - P(a)$$

inclusion-exclusion principle

$$P(a \vee b) = P(a) + P(b) - P(a \wedge b)$$

chain rule

$$P(A, B, C, D)$$

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

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

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

Inference using full joint distributions

How to answer questions?

Knowledge base is

represented using full joint distribution.

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

To compute posterior probability of a query proposition given observed evidence, we add up probabilities of possible worlds in which the proposition is true (**marginalization** or **summing out**).

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

$$P(Y) = \sum_{z \in Z} P(Y, z)$$

Example of probabilistic inference

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

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

$$P(Y) = \sum_{z \in Z} P(Y, z)$$

$$\begin{aligned} P(\text{toothache}) & (= P(\text{Toothache}=\text{true})) \\ &= 0.108 + 0.012 + 0.016 + 0.064 = 0.2 \end{aligned}$$

$$\begin{aligned} P(\text{cavity} \vee \text{toothache}) \\ &= 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28 \end{aligned}$$

$$\begin{aligned} P(\neg \text{cavity} | \text{toothache}) \\ &= P(\neg \text{cavity} \wedge \text{toothache}) / P(\text{toothache}) \\ &= (0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064) \\ &= 0.4 \end{aligned}$$

Normalization

$$\begin{aligned} P(\neg \text{cavity} | \text{toothache}) \\ &= P(\neg \text{cavity} \wedge \text{toothache}) / P(\text{toothache}) \\ &= (0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064) = 0.4 \end{aligned}$$

$$\begin{aligned} P(\text{cavity} | \text{toothache}) \\ &= P(\text{cavity} \wedge \text{toothache}) / P(\text{toothache}) \\ &= (0.108 + 0.012) / (0.108 + 0.012 + 0.016 + 0.064) = 0.6 \end{aligned}$$

Notice that denominators are identical in both formulas!

We even do not need to know the exact value of denominator:

$$P(\neg \text{cavity} | \text{toothache}) + P(\text{cavity} | \text{toothache}) = 1$$

We can use a **normalization constant** α instead, computed such that the evaluated distribution adds up to 1.

$$\begin{aligned} P(\text{Cavity} | \text{toothache}) &= \alpha P(\text{Cavity}, \text{toothache}) \\ &= \alpha [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg \text{catch})] \\ &= \alpha [\langle 0.108, 0.016 \rangle + \langle 0.012, 0.064 \rangle] \\ &= \alpha [\langle 0.12, 0.08 \rangle] = [\langle 0.6, 0.4 \rangle] \end{aligned}$$

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

In a typical case, we know values e of random variables E from the **observation** and we are looking for probability distribution of random variables Y from the **query**.

The other random variables are **hidden** $H = X - Y - E$.

$$P(Y \mid E=e) = \alpha P(Y, E=e) = \alpha \sum_h P(Y, E=e, H=h)$$

Some drawbacks of inference by enumeration:

- the **worst-case time complexity is $O(d^n)$** , where d is the number of values in domains of each random variable
- **to store full joint probability distribution we need $O(d^n)$ space**
- last but not least, it is not easy to obtain probabilities for all possible worlds



Let us expand the full joint distribution by adding a fourth variable Weather with the domain {cloudy, sunny, rain, snow} – the new full joint distribution has $2 \times 2 \times 2 \times 4 = 32$ elements (possible worlds).

$$\begin{aligned} P(\text{toothache}, \text{catch}, \text{cavity}, \text{cloudy}) \\ = P(\text{cloudy} \mid \text{toothache}, \text{catch}, \text{cavity}) * P(\text{toothache}, \text{catch}, \text{cavity}) \end{aligned}$$

Do one's dental problems influence the weather?

$$P(\text{cloudy} \mid \text{toothache}, \text{catch}, \text{cavity}) = P(\text{cloudy})$$

We can write in general:

$$\begin{aligned} P(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) \\ = P(\text{Toothache}, \text{Catch}, \text{Cavity}) * P(\text{Weather}) \end{aligned}$$

Hence the full joint distribution can be constructed from two smaller tables, one with 8 elements and one with 4 elements.

This property is called **(absolute) independence**:

$$P(X \mid Y) = P(X) \text{ or } P(Y \mid X) = P(Y) \text{ or } P(X, Y) = P(X).P(Y)$$



Conditional independence

Full independence allows us **reducing the size** of the domain representation, but unfortunately full independence is rare and even independent subsets can be quite large.

When one has cavity, does catch depend on toothache?

$$P(\text{catch} \mid \text{toothache}, \text{cavity}) = P(\text{catch} \mid \text{cavity})$$

$$P(\text{catch} \mid \text{toothache}, \neg \text{cavity}) = P(\text{catch} \mid \neg \text{cavity})$$

Random variables Catch and Toothache are independent if we know the value of Cavity.

$$P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$$

This property is called **conditional independence**:

$$P(X \mid Y, Z) = P(X \mid Y) \text{ or } P(Z \mid X, Y) = P(Z \mid Y) \text{ or}$$

$$P(Z, X \mid Y) = P(Z \mid Y) P(X \mid Y)$$

Exploiting conditional independence

Conditional independence can be used to further reduce the size of domain representation.

$$P(\text{Toothache}, \text{Catch}, \text{Cavity})$$

$$= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch}, \text{Cavity})$$

$$= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$$

$$= P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$$

The full joint distribution can be constructed from three smaller tables of sizes $2 + 2 + 1 = 5$ (only independent elements are represented).

Let us go back to **diagnostic problems**.

Usually we are looking for disease (the source of problems) based on symptoms (observations).

- we are interested in the **diagnostic direction** expressed as conditional probability

$$P(\text{disease} | \text{symptoms})$$

However, from past experience we often have other information:

- the probability of disease $P(\text{disease})$
- the probability of symptoms $P(\text{symptoms})$
- the **causal relation** expressed as conditional probability $P(\text{symptoms} | \text{disease})$



How can this information be exploited to get the probability of the diagnostic direction?

Recall the product rule

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

We can deduce a so called **Bayes' rule** (law or theorem):

$$P(a | b) = P(b | a) P(a) / P(b)$$

in general:

$$P(Y | X) = P(X | Y) P(Y) / P(X) = \alpha P(X | Y) P(Y)$$

It looks like two steps backward as now we need to know $P(X | Y)$, $P(Y)$, $P(X)$.

But these are the values that we frequently have.

$$P(\text{cause} | \text{effect}) = P(\text{effect} | \text{cause}) P(\text{cause}) / P(\text{effect})$$

- $P(\text{effect} | \text{cause})$ describes the **causal direction**
- $P(\text{cause} | \text{effect})$ describes the **diagnostic relation**

Medical diagnosis

- from past cases we know $P(\text{symptoms} | \text{disease})$, $P(\text{disease})$, $P(\text{symptoms})$
- for a new patient we know symptoms and looking for diagnosis $P(\text{disease} | \text{symptoms})$

Example:

- meningitis causes a stiff neck 70% of the time
- the prior probability of meningitis is 1/50 000
- the prior probability of stiff neck is 1%

What is the probability that a patient having a stiff neck has meningitis?

$$P(m|s) = P(s|m).P(m) / P(s) = 0.7 * 1/50000 / 0.01 = 0.0014$$

Why the conditional probability for the diagnostic direction is not stored directly?

- diagnostic knowledge is often more fragile than causal knowledge
- for example, if there is a sudden epidemic of meningitis, the unconditional probability of meningitis $P(m)$ will go up so $P(m|s)$ should also go up while the causal relation $P(s|m)$ is unaffected by the epidemic, as it reflects how meningitis works

What if there are more observations?

We can exploit conditional independence as follows

$$\begin{aligned} &P(\text{Toothache}, \text{Catch}, \text{Cavity}) \\ &= P(\text{Toothache} | \text{Cavity}) P(\text{Catch} | \text{Cavity}) P(\text{Cavity}) \end{aligned}$$

If all the effects are conditionally independent given the cause variable, we get:

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

Such a probability distribution is called a **naive Bayes model** (it is often used even in cases where the “effect” variables are not actually conditionally independent given the value of the cause variable).

Wumpus is back!

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

We have a maze with pits that are detected in neighboring squares via breeze (Wumpus and gold will not be assumed now).

Where does the agent should go, if there is breeze at (1,2) and (2,1)?

Pure logical inference can conclude nothing about which square is most likely to be safe!



To which square does the agent should go?

Wumpus: probabilistic model

Boolean variables:

$P_{i,j}$ – pit at square (i,j)

$B_{i,j}$ – breeze at square (i,j)

(only for the observed squares $B_{1,1}$, $B_{1,2}$ a $B_{2,1}$).

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

Full joint probability distribution

$$P(P_{1,2}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$$

$$= P(B_{1,1}, B_{1,2}, B_{2,1} | P_{1,2}, \dots, P_{4,4}) * P(P_{1,2}, \dots, P_{4,4})$$

$$P(P_{1,2}, \dots, P_{4,4}) = \prod_{i,j} P(P_{i,j})$$

$$P(P_{1,2}, \dots, P_{4,4}) = 0.2^n * 0.8^{16-n}$$

product rule

pits are spread independently

probability of pit is 0.2 and there are n pits

Wumpus: query and simple reasoning

Assume that we have **evidence**:

$$b = b_{1,1} \wedge b_{1,2} \wedge b_{2,1}$$

$$\text{known} = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$$

We are interested in answering queries such as $P(P_{1,3} \mid \text{known}, b)$.

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

Answer can be computed by enumeration of the full joint probability distribution.

Let Unknown be the variables $P_{i,j}$ except $P_{1,3}$ and Known:

$$P(P_{1,3} \mid \text{known}, b)$$

$$= \sum_{\text{unknown}} P(P_{1,3}, \text{unknown}, \text{known}, b)$$

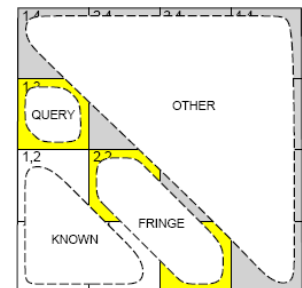
But it means to explore all possible values of variables Unknown and there are $2^{12} = 4096$ terms!

Can we do it better (faster)?

Wumpus: conditional independence

Observation:

The observed breezes are conditionally independent of the other variables given the known (white), frontier (yellow), and query variables.



We split the set of hidden variables into fringe and other variables:

$$\text{Unknown} = \text{Fringe} \cup \text{Other}$$

From conditional independence we get:

$$P(b \mid P_{1,3}, \text{known}, \text{unknown}) = P(b \mid P_{1,3}, \text{known}, \text{fringe})$$

Now, let us exploit this formula.

$$P(P_{1,3} \mid \text{known}, b)$$

$$= \alpha \sum_{\text{unknown}} P(P_{1,3}, \text{known}, \text{unknown}, b)$$

product rule $P(X,Y) = P(X|Y) P(Y)$

$$= \alpha \sum_{\text{unknown}} P(b \mid P_{1,3}, \text{known}, \text{unknown}) * P(P_{1,3}, \text{known}, \text{unknown})$$

$$= \alpha \sum_{\text{fringe}} \sum_{\text{other}} P(b \mid P_{1,3}, \text{known}, \text{fringe}, \text{other}) * P(P_{1,3}, \text{known}, \text{fringe}, \text{other})$$

$$= \alpha \sum_{\text{fringe}} \sum_{\text{other}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) * P(P_{1,3}, \text{known}, \text{fringe}, \text{other})$$

$$= \alpha \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) * \sum_{\text{other}} P(P_{1,3}, \text{known}, \text{fringe}, \text{other})$$

$$= \alpha \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) * \sum_{\text{other}} P(P_{1,3}) P(\text{known}) P(\text{fringe}) P(\text{other})$$

$$= \alpha P(\text{known}) P(P_{1,3}) \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe}) \sum_{\text{other}} P(\text{other})$$

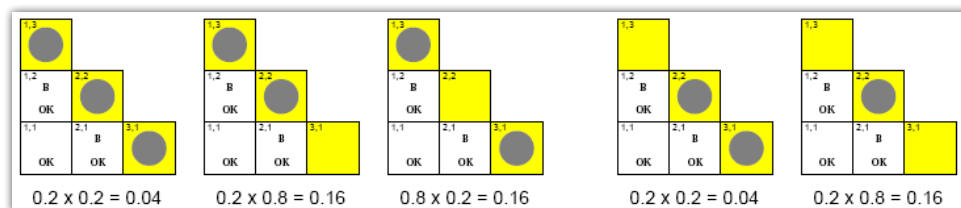
$$= \alpha' P(P_{1,3}) \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe})$$

$$\alpha' = \alpha \cdot P(\text{known})$$

$$\sum_{\text{other}} P(\text{other}) = 1$$

$$P(P_{1,3} \mid \text{known}, b) = \alpha' P(P_{1,3}) \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe})$$

Let us explore possible models (values) of Fringe that are compatible with observation b.



$$P(P_{1,3} \mid \text{known}, b)$$

$$= \alpha' \langle 0.2 (0.04 + 0.16 + 0.16), 0.8 (0.04 + 0.16) \rangle$$

$$= \langle 0.31, 0.69 \rangle$$

$$P(P_{2,2} \mid \text{known}, b) = \langle 0.86, 0.14 \rangle$$

Definitely avoid the square (2,2)!

Probability theory is a formal mechanism to handle **uncertainty**.

Full joint distribution describes probabilities of all possible worlds.

Answers to queries can be obtained by summing out probabilities of possible worlds consistent with the observation.

However, larger problems will require a better approach.

We are going to exploit **independence** and **conditional independence**.



© 2016 Roman Barták

Department of Theoretical Computer Science and Mathematical Logic
bartak@ktiml.mff.cuni.cz