



Knowledge Representation

Introduction to Artificial Intelligence

Chandra Gummaluru
University of Toronto

Version W22.1

- The following is based on material developed by many individuals, including (but not limited to):
 - Sheila McIlraith
 - Bahar Aameri
 - Fahiem Bacchus
 - Sonya Allin

- Part of being an intelligent agent involves being able to infer implicit facts based on known or assumed ones.

Example: Avatar Aang's Love Life

- Suppose we knew the following:
 - Anyone who likes someone is sad if they do not like them back.
 - Aang likes Katara.
 - Aang is an fire bender.
 - Aang is sad.
 - Katara does not like all fire benders.
- Does Katara like Aang?
- Humans develop this ability through experience. Our goal is to instil artificial agents with the same ability.
- Without reasoning, we would have to explicitly remember every fact we've learned.

- To achieve this ability, we do two things:
 - ① **Represent** (encode) known statements in our brain.
 - ② **Reason** (infer) new statements from the known ones.
- Thus, to achieve this ability artificially, we need to do things:
 - ① Develop a **formal languages** to represent statements.
 - ② Develop a **reasoning mechanism** for the formal system.
- There are many formal languages and reasoning mechanisms we could use.
- We will consider a representation called **first-order logic** (FOL), and a reasoning mechanism called **resolution**.

Informal Languages versus Formal Languages

- Roughly speaking, the goal of any formal language is to facilitate the expression of knowledge but with strict rules to avoid any ambiguity.
- **Example:** Ambiguity in English
 - What is the appropriate response to the request, “call me an ambulance,”?
 - “okay.”
 - “uh...you’re an ambulance...”.
- The ambiguity arises from the fact that, in English (and other languages), there are multiple interpretations of many words/phrases.

- To avoid such ambiguity, a formal language must define the notion used to build its statements, as well as a system for interpreting those statements.
- The notion is called the **syntax** and the interpretations are the **semantics**.
- So, a formal language needs to provide syntax and a way to introduce semantics. However, it does not provide the semantics themselves.

- FOL syntax consists of the following components:
 - ① **variables**, where each variable is, by definition, a **term**.
 - E.g., x .
 - ② **functions**, which each map zero or more terms to a single term.
 - E.g., $\text{nation}(x)$, which should refer to x 's nation.
 - ③ **predicates**, which each map many terms to true/false.
 - E.g., $\text{likes}(x, y)$, which should mean that x likes y .
- An FOL **vocabulary** is a triple, $\mathcal{L} = (\mathcal{V}, \mathcal{F}, \mathcal{P})$, where \mathcal{V} , \mathcal{F} , and \mathcal{P} are sets of variables, functions, and predicates, respectively.

- Let \mathcal{L} be a vocabulary.
- For any n -ary \mathcal{L} -predicate P , and \mathcal{L} -terms, t_1, \dots, t_n , the expression, $P(t_1, \dots, t_n)$ is called an atomic \mathcal{L} -formula.
 - E.g., $\text{likes}(x, \text{nation}(y))$.
- Atomic formula represents the most fundamental statements.

- Non-atomic \mathcal{L} -formulae are defined recursively as follows:
 - **negation:** $\neg f$, where f is any \mathcal{L} -formula.
 - **disjunction:** $f_1 \vee f_2$, where f_1 and f_2 are \mathcal{L} -formulae.
 - **conjunction:** $f_1 \wedge f_2$, where f_1 and f_2 are \mathcal{L} -formulae.
 - **implication:** $f_1 \rightarrow f_2$, where f_1 and f_2 are \mathcal{L} -formulae.
 - **existential:** $\exists x f$, where x is a variable and f is any \mathcal{L} -formula.
 - **universal:** $\forall x f$, where x is a variable and f is any \mathcal{L} -formula.

Example: Aang's Love Life

- Suppose we have a vocabulary, \mathcal{L} , with:
 - constant symbols: Katara, Aang, Zuko, Azula, fire, water, air, earth
 - variable symbols; x, y
 - function symbols; $\text{sibling}(x)$, and $\text{nation}(x)$,
 - predicate symbols; $\text{person}(x)$, $\text{sad}(x)$, $\text{likes}(x, y)$, and $\text{bender}(x, y)$.
- We can express each fact using our vocabulary:

① Anyone who is likes someone is sad if they don't like them back:

$$\forall x (\text{person}(x) \rightarrow (\exists y (\text{person}(y) \wedge \text{likes}(x, y) \wedge \neg \text{likes}(y, x)) \rightarrow \text{sad}(x)))$$

- ② Aang likes Katara: $\text{likes}(\text{Aang}, \text{Katara})$.
- ③ Aang is a fire bender: $\text{bender}(\text{Aang}, \text{fire})$.
- ④ Katara does not like all fire benders: $\exists x \text{bender}(x, \text{fire}) \wedge \neg \text{likes}(\text{Katara}, x)$.
- ⑤ Aang is sad: $\text{sad}(\text{Aang})$.

- In FOL, the semantics are provided by what we refer to as a **structure**. The model, \mathcal{S} , consists of the following components:
 - ① a **domain of discourse**, D , which is a set of relevant elementary objects.
 - ② **specializations of functions**, $f^{\mathcal{S}} : D^n \rightarrow D$, for each n -ary function, f , so that $f^{\mathcal{S}}$ assigns f within D .
 - ③ **specializations of predicates**, $p^{\mathcal{S}} \subseteq D^n$, for each n -ary predicate, p , so that $p(t_1, \dots, t_n)$ is true if and only if $(t_1, \dots, t_n) \in p^{\mathcal{S}}$.

Example: Avatar Aang's Love Life

- A structure, \mathcal{S} of \mathcal{L} is shown below:
 - $D = \{\mathbf{Aang}, \mathbf{Katara}, \mathbf{Zuko}, \mathbf{Azula}, \mathbf{water}, \mathbf{fire}, \mathbf{air}, \mathbf{earth}\}$
 - $\text{Katara}^{\mathcal{S}} = \mathbf{Katara}$, $\text{Aang}^{\mathcal{S}} = \mathbf{Aang}$, $\text{Zuko}^{\mathcal{S}} = \mathbf{Zuko}$, $\text{Azula}^{\mathcal{S}} = \mathbf{Azula}$
 - $\text{sibling}^{\mathcal{S}}(\mathbf{Zuko}) = \mathbf{Azula}$, $\text{sibling}^{\mathcal{S}}(\mathbf{Azula}) = \mathbf{Zuko}$
 - $\text{nation}^{\mathcal{S}}(\mathbf{Aang}) = \mathbf{air}$, $\text{nation}^{\mathcal{S}}(\mathbf{Katara}) = \mathbf{water}$, $\text{nation}^{\mathcal{S}}(\mathbf{Zuko}) = \mathbf{fire}$,
 $\text{nation}^{\mathcal{S}}(\mathbf{Azula}) = \mathbf{fire}$
 - $\text{bender}^{\mathcal{S}} = \{\langle \mathbf{Aang}, \mathbf{water} \rangle, \langle \mathbf{Aang}, \mathbf{fire} \rangle, \langle \mathbf{Aang}, \mathbf{air} \rangle, \langle \mathbf{Aang}, \mathbf{earth} \rangle, \langle \mathbf{Katara}, \mathbf{water} \rangle, \langle \mathbf{Zuko}, \mathbf{fire} \rangle, \langle \mathbf{Azula}, \mathbf{fire} \rangle\}$
 - $\text{likes}^{\mathcal{S}} = \{\langle \mathbf{Katara}, \mathbf{water} \rangle, \langle \mathbf{Aang}, \mathbf{Katara} \rangle, \langle \mathbf{Katara}, \mathbf{Aang} \rangle\}$
 - $\text{person}^{\mathcal{S}} = \{\langle \mathbf{Aang} \rangle, \langle \mathbf{Katara} \rangle, \langle \mathbf{Zuko} \rangle, \langle \mathbf{Azula} \rangle\}$
 - $\text{sad}^{\mathcal{S}} = \{\langle \mathbf{Aang} \rangle, \langle \mathbf{Zuko} \rangle\}$

- An occurrence of a variable, x , in an FOL formula, f , is **bound** if and only if it is in a sub-formula of the form $\forall x f'$ or $\exists x f'$. Otherwise, x is **free**.
- If a formula f contains free variables, then its value depends on what values we assign to those free variables.
- We define an **assignment function**, $\sigma : \mathcal{V} \rightarrow D$, so that $\sigma(x)$ is the element in the universe represented by the variable x .
- The value of any \mathcal{L} -term is defined recursively with an **extended assignment function**, $\bar{\sigma}$ so that $\bar{\sigma}(x) = \sigma(x)$ and $\bar{\sigma}(f(t_1, \dots, t_n)) = f^S(\bar{\sigma}(t_1), \dots, \bar{\sigma}(t_n))$.

Example: Avatar Aang's Love Life

- Suppose we wanted to compute the value of $\text{nation}(\text{sibling}(x))$, under $\sigma(x) = \mathbf{Azula}$.
- We proceed as follows:

$$\begin{aligned}\bar{\sigma}(\text{nation}(\text{sibling}(x))) &= \text{nation}^{\mathcal{S}}(\bar{\sigma}(\text{sibling}(x))) \\ &= \text{nation}^{\mathcal{S}}(\text{sibling}^{\mathcal{S}}(\bar{\sigma}(x))) \\ &= \text{nation}^{\mathcal{S}}(\text{sibling}^{\mathcal{S}}(\sigma(x))) \\ &= \text{nation}^{\mathcal{S}}(\text{sibling}^{\mathcal{S}}(\mathbf{Azula})) \\ &= \text{nation}^{\mathcal{S}}(\mathbf{Zuko}) \\ &= \mathbf{fire}\end{aligned}$$

- Notice that the value depends on both \mathcal{S} and σ .

- We write $\mathcal{S} \models f[\sigma]$ to denote that \mathcal{S} **satisfies** the formula, f , under σ .
- By definition, $\mathcal{S} \models P(t_1, \dots, t_n)[\sigma]$ if and only if $\langle \bar{\sigma}(t_1), \dots, \bar{\sigma}(t_n) \rangle \in P^{\mathcal{S}}$.
- For other \mathcal{L} -formulae, \models is defined recursively as follows:
 - $\mathcal{S} \models (t_1 = t_2)[\sigma]$ if and only if $\bar{\sigma}(t_1) = \bar{\sigma}(t_2)$
 - $\mathcal{S} \models \neg f[\sigma]$ if and only if $\mathcal{S} \not\models f[\sigma]$
 - $\mathcal{S} \models (f_1 \vee f_2)[\sigma]$ if and only if $\mathcal{S} \models f_1[\sigma]$ or $\mathcal{S} \models f_2[\sigma]$
 - $\mathcal{S} \models (f_1 \wedge f_2)[\sigma]$ if and only if $\mathcal{S} \models f_1[\sigma]$ and $\mathcal{S} \models f_2[\sigma]$
 - $\mathcal{S} \models (f_1 \rightarrow f_2)[\sigma]$ if and only if $\mathcal{S} \models (f_1 \vee \neg f_2)[\sigma]$
 - $\mathcal{S} \models (\forall x f)[\sigma]$ if and only if $\mathcal{S} \models f[\sigma[x/m]]$ for all $m \in D$
 - $\mathcal{S} \models (\exists x f)[\sigma]$ if and only if $\mathcal{S} \models f[\sigma[x/m]]$ for some $m \in D$
- Here $\sigma[x/m]$ is defined assuming x is fixed and so that

$$\sigma[x/m](y) = \begin{cases} \sigma(y), & y \neq x \\ m, & y = x \end{cases}$$

Example: Avatar Aang's Love Life

- Suppose we wanted to compute the value of $\exists x \text{bender}(x, \text{fire}) \wedge \neg \text{likes}(\text{Katara}, x)$.
- We proceed as follows:

$$\mathcal{S} \models (\exists x (\text{bender}(x, \text{fire}) \wedge \neg \text{likes}(\text{Katara}, x))) [\sigma]$$

$$\Leftrightarrow \mathcal{S} \models (\text{bender}(x, \text{fire}) \wedge \neg \text{likes}(\text{Katara}, x)) [\sigma[x/m]] \text{ for some } m \in D$$

$$\Leftrightarrow \mathcal{S} \models \text{bender}(x, \text{fire})[\sigma[x/m]] \text{ and } \mathcal{S} \not\models \text{likes}(\text{Katara}, x)[\sigma[x/m]] \text{ for some } m \in D$$

$$\Leftrightarrow \langle m, \text{fire} \rangle \in \text{bender}^{\mathcal{S}} \text{ and } \langle \text{Katara}, m \rangle \notin \text{likes}^{\mathcal{S}} \text{ for some } m \in D$$

The above holds for $m = \mathbf{Azula}$.

- A **sentence** is any formula consisting of only bound variables.
- The value of a sentence, s , is independent of the variable assignments, i.e., for any σ and σ' ,

$$\mathcal{S} \models s[\sigma] \text{ if and only if } \mathcal{S} \models s[\sigma'].$$

- Thus, for any sentence, s , we simply write $\mathcal{S} \models s$ to denote that \mathcal{S} satisfies s .

- Let Φ be a set of sentences.
- We say that Φ is **satisfiable** if there exists some structure \mathcal{S} such that $\mathcal{S} \models s$ for every $s \in \Phi$; such a structure is called a **model** of Φ .
- Say we are given another set of sentences Φ' :
 - ① If every model of Φ is a model of Φ' , we say that Φ' is a consequence of Φ .
 - ② If no model of Φ is a model of Φ' , we say that Φ' is a contradiction of Φ .
 - ③ If some models of Φ are models of Φ' , but other models of Φ are not models of Φ' , then Φ' is neither a consequence of contradiction of Φ .
- Proving logical consequences / contradictions can only be done through a syntactic reasoning mechanism.