

# Using FOL Theorem Provers

Marijn J.H. Heule

Carnegie  
Mellon  
University

## Introduction

## Herbrand's Theorem

## Using FOL Theorem Provers

# Introduction

Herbrand's Theorem

Using FOL Theorem Provers

# Introduction: Recap FOL

## Undecidable

- ▶ No single algorithm decides on **all** formulas
- ▶ E.g., no decision procedure for non-linear integer arithmetic

# Introduction: Recap FOL

## Undecidable

- ▶ No single algorithm decides on **all** formulas
- ▶ E.g., no decision procedure for non-linear integer arithmetic

## Refutationally complete

- ▶ Decision procedures for **valid/unsatisfiable** formulas
- ▶ E.g., pure FOL formulas with universal sentences
- ▶ The algorithm could run forever on satisfiable formulas

# Introduction: Recap FOL

## Undecidable

- ▶ No single algorithm decides on **all** formulas
- ▶ E.g., no decision procedure for non-linear integer arithmetic

## Refutationally complete

- ▶ Decision procedures for **valid/unsatisfiable** formulas
- ▶ E.g., pure FOL formulas with universal sentences
- ▶ The algorithm could run forever on satisfiable formulas

## Decidable

- ▶ Decision procedures for **all** formulas in a FOL fragment
- ▶ E.g., Presburger arithmetic, uninterpreted functions

# Introduction: Differences FOL and Propositional Logic

## Decidable

- ▶ Propositional logic is decidable
- ▶ Only some fragments of FOL are decidable

# Introduction: Differences FOL and Propositional Logic

## Decidable

- ▶ Propositional logic is decidable
- ▶ Only some fragments of FOL are decidable

## Complexity of resolution

- ▶ Variable elimination determines satisfiability in Prop
- ▶ Clauses can produce infinitely many resolvents in FOL
- ▶ E.g.  $\forall x. \neg P(x) \vee P(f(x))$

# Introduction: Differences FOL and Propositional Logic

## Decidable

- ▶ Propositional logic is decidable
- ▶ Only some fragments of FOL are decidable

## Complexity of resolution

- ▶ Variable elimination determines satisfiability in Prop
- ▶ Clauses can produce infinitely many resolvents in FOL
- ▶ E.g.  $\forall x. \neg P(x) \vee P(f(x))$

## Properties of resolution

- ▶ In Prop, resolving on multiple clashing pairs  $\rightarrow$  tautology
- ▶ In FOL, literals may be merged
- ▶ E.g.  $\forall y. \neg S(a, y) \vee \neg S(y, y)$  and  $\forall z. S(z, z) \vee S(a, z)$

# Introduction

## Herbrand's Theorem

## Using FOL Theorem Provers

## The Big Picture

We turn a FOL formula of universal sentences into a (typically infinite) propositional formula and have the following:

- ▶ The propositional formula is satisfiable iff every finite subset of the formula is satisfiable ([Compactness Lemma](#)).
- ▶ The satisfying assignment of the (possibly infinite) propositional formula can be turned into a (possibly infinite) model of the FOL formula ([Herbrand's Theorem](#)).
- ▶ If a finite set of the propositional formula is unsatisfiable, then the FOL formula is unsatisfiable ([Lifting Lemma](#)).

# Compactness Lemma (Propositional Logic)

## Theorem (Compactness)

Let  $\tau$  be any partial truth assignment, and suppose that **every finite subset of  $\Gamma$  is satisfied** by a truth assignment that extends  $\tau$ . Then for every propositional variable  $p_i$ , either **every finite subset of  $\Gamma$  is satisfied by a truth assignment that extends  $\tau[p_i \mapsto \top]$**  or **every finite subset of  $\Gamma$  is satisfied by a truth assignment that extends  $\tau[p_i \mapsto \perp]$** .

# Compactness Lemma (Propositional Logic)

## Theorem (Compactness)

Let  $\tau$  be any partial truth assignment, and suppose that **every finite subset of  $\Gamma$  is satisfied** by a truth assignment that extends  $\tau$ . Then for every propositional variable  $p_i$ , either **every finite subset of  $\Gamma$  is satisfied by a truth assignment that extends  $\tau[p_i \mapsto \top]$**  or **every finite subset of  $\Gamma$  is satisfied by a truth assignment that extends  $\tau[p_i \mapsto \perp]$** .

## Proof.

Suppose otherwise. Then there is a finite subset  $\Gamma_0$  of  $\Gamma$  that is **not satisfied** by any truth assignment that extends  $\tau[p_i \mapsto \top]$  and there is a finite subset  $\Gamma_1$  of  $\Gamma$  that is **not satisfied** by any truth assignment that extends  $\tau[p_i \mapsto \perp]$ . Then there is no truth assignment extending  $\tau$  that satisfies  $\Gamma_0 \cup \Gamma_1$ , because any such truth assignment has to assign  $\top$  or  $\perp$  to  $p_i$ . □

# Herbrand's Theorem

## Theorem

*Let  $\Gamma$  be a set of universal first-order formulas in a language that contains at least one constant. Suppose every finite set of ground instances of formulas in  $\Gamma$  is satisfiable as a set of propositional formulas. Then there exists a model of  $\Gamma$ .*

# Herbrand's Theorem

## Theorem

*Let  $\Gamma$  be a set of universal first-order formulas in a language that contains at least one constant. Suppose every finite set of ground instances of formulas in  $\Gamma$  is satisfiable as a set of propositional formulas. Then there exists a model of  $\Gamma$ .*

## Theorem (Herbrand's Theorem)

*Let  $A = \forall x_1 \dots x_n. A_1 \wedge \dots \wedge A_m$  be a clause normal form formula with constant and function symbols from  $\Sigma$ . Let  $\Sigma'$  be the set of ground terms that can be made from symbols in  $\Sigma$ .*

*A is unsatisfiable if and only if there is a finite set  $\Gamma$  where:*

- ▶ *Each element in  $\Gamma$  is a clause  $A_i[t_1/x_1, \dots, t_n/x_n]$  where  $1 \leq i \leq m$  and  $t_1 \dots t_n \in \Sigma'$ .*
- ▶ *If each distinct literal in  $\Gamma$  is interpreted as a unique prop variable, then  $\Gamma$  is unsatisfiable in prop logic.*

## Herbrand Model

Herbrand universe: set of all ground terms

### Example

Consider  $\Gamma := \forall x, y. (f(x) \neq c) \wedge (f(x) \neq f(y)) \vee x = y$ .

The Herbrand universe of  $\Gamma$  is  $|\mathfrak{M}| := \{c, f(c), f(f(c)), \dots\}$ .

The Herbrand universe can be infinite while  $\Gamma$  is finite.

## Herbrand Model

Herbrand universe: set of all ground terms

### Example

Consider  $\Gamma := \forall x, y. (f(x) \neq c) \wedge (f(x) \neq f(y)) \vee x = y$ .

The Herbrand universe of  $\Gamma$  is  $|\mathfrak{M}| := \{c, f(c), f(f(c)), \dots\}$ .

The Herbrand universe can be infinite while  $\Gamma$  is finite.

Herbrand model: interpret each constant/function as itself

- ▶  $c^{\mathfrak{M}}$  is interpreted as  $c$
- ▶  $f^{\mathfrak{M}}(t_1, \dots, t_n)$  is interpreted as  $f(t_1, \dots, t_n)$

# Transforming FOL to Propositional Logic

## Example

Let  $\Gamma := \forall x, y. \neg R(f(x), c) \wedge (\neg R(f(x), f(y)) \vee R(x, y))$ .

The Herbrand universe of  $\Gamma$  is  $|\mathfrak{M}| := \{c, f(c), f(f(c)), \dots\}$ .

$\Gamma$  can be transformed into the (infinite) propositional formula:

$$\neg p_{R(f(c),c)} \wedge (\neg p_{R(f(c),f(c))} \vee p_{R(c,c)}) \wedge \quad x = c, y = c$$

$$\neg p_{R(f(f(c)),c)} \wedge (\neg p_{R(f(f(c)),f(y))} \vee p_{R(f(c),c)}) \wedge \quad x = f(c), y = c$$

$$\neg p_{R(f(c),c)} \wedge (\neg p_{R(f(c),f(f(c)))} \vee p_{R(c,f(c))}) \wedge \quad x = c, y = f(c)$$

...

# Transforming FOL to Propositional Logic

## Example

Let  $\Gamma := \forall x, y. \neg R(f(x), c) \wedge (\neg R(f(x), f(y)) \vee R(x, y))$ .

The Herbrand universe of  $\Gamma$  is  $|\mathfrak{M}| := \{c, f(c), f(f(c)), \dots\}$ .

$\Gamma$  can be transformed into the (infinite) propositional formula:

$$\neg p_{R(f(c),c)} \wedge (\neg p_{R(f(c),f(c))} \vee p_{R(c,c)}) \wedge \quad x = c, y = c$$

$$\neg p_{R(f(f(c)),c)} \wedge (\neg p_{R(f(f(c)),f(y))} \vee p_{R(f(c),c)}) \wedge \quad x = f(c), y = c$$

$$\neg p_{R(f(c),c)} \wedge (\neg p_{R(f(c),f(f(c)))} \vee p_{R(c,f(c))}) \wedge \quad x = c, y = f(c)$$

...

The equality relation needs special attention

- ▶ Replace it with a new relation, say  $E(x, y)$
- ▶ Enforce symmetry, transitivity, and congruence

## Proof of Herbrand's Theorem

Suppose every finite set of ground instances of the formulas in  $\Gamma$  is **satisfiable** as a set of propositional formulas. Then, by the **compactness theorem** for propositional logic, there is a truth assignment  $\tau$  that **satisfies** all the ground instances of the formulas in  $\Gamma$ .

## Proof of Herbrand's Theorem

Suppose every finite set of ground instances of the formulas in  $\Gamma$  is **satisfiable** as a set of propositional formulas. Then, by the **compactness theorem** for propositional logic, there is a truth assignment  $\tau$  that **satisfies** all the ground instances of the formulas in  $\Gamma$ .

We look to  $\tau$  to determine the truth of the atomic formulas in the language of  $\Gamma$ . Let  $\mathfrak{M}$  be the **Herbrand model** of  $\Gamma$ .

For every  $R$ , we define  $R^{\mathfrak{M}}$  to hold iff  $\tau(R(t_1, \dots, t_n))$  is true. In other words, **first-order evaluation** in the model  $\mathfrak{M}$  is the same as propositional evaluation under the truth assignment  $\tau$ .

## Proof of Herbrand's Theorem

Suppose every finite set of ground instances of the formulas in  $\Gamma$  is **satisfiable** as a set of propositional formulas. Then, by the **compactness theorem** for propositional logic, there is a truth assignment  $\tau$  that **satisfies** all the ground instances of the formulas in  $\Gamma$ .

We look to  $\tau$  to determine the truth of the atomic formulas in the language of  $\Gamma$ . Let  $\mathfrak{M}$  be the **Herbrand model** of  $\Gamma$ .

For every  $R$ , we define  $R^{\mathfrak{M}}$  to hold iff  $\tau(R(t_1, \dots, t_n))$  is true. In other words, **first-order evaluation** in the model  $\mathfrak{M}$  is the same as propositional evaluation under the truth assignment  $\tau$ .

Since the universe of  $\mathfrak{M}$  consists of ground terms, a formula  $\forall x_1, \dots, x_n. A$  in  $\Gamma$  is true in  $\mathfrak{M}$  iff for every  $t_1, \dots, t_n$ ,  $A$  is true under the assignment  $\{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$ . By the definition of  $\mathfrak{M}$ , this is the case iff  $\tau(A(t_1, \dots, t_n))$  is true. This holds as  $\tau$  satisfies every closed instance of a formula in  $\Gamma$ .

## Lifting Lemma

If a finite set of the propositional formula is unsatisfiable, then the FOL formula is unsatisfiable.

## Lifting Lemma

If a finite set of the propositional formula is unsatisfiable, then the FOL formula is unsatisfiable.

### Lemma (Lifting Lemma)

Let  $A = \forall x_1, \dots, x_n. A_1 \wedge \dots \wedge A_m$  be a clause normal form formula with constant and function symbols from  $\Sigma$ . Let  $\Sigma'$  be the set of closed terms that can be made from symbols in  $\Sigma$ .

Let  $\Gamma$  be a set where each element is a clause  $A_i[t_1/x_1, \dots, t_n/x_n]$  where  $1 \leq i \leq m$  and  $t_1 \dots t_n \in \Sigma'$ .

If each distinct literal in  $\Gamma$  is interpreted as a unique propositional variable, then any propositional resolution refutation of  $\Gamma$  can be transformed into a first-order resolution refutation of  $A$ .

## Introduction

## Herbrand's Theorem

## Using FOL Theorem Provers

## Using FOL Theorem Provers: Aunt Agatha

Someone who lives in Dreadbury Mansion killed Aunt Agatha. Agatha, the butler, and Charles live in Dreadbury Mansion, and are the only people who live therein. A killer always hates his victim, and is never richer than his victim. Charles hates no one that Aunt Agatha hates. Agatha hates everyone except the butler. The butler hates everyone not richer than Aunt Agatha. The butler hates everyone Aunt Agatha hates. No one hates everyone. Agatha is not the butler.

Who killed Aunt Agatha?

# Using FOL Theorem Provers: Aunt Agatha

Someone who lives in Dreadbury Mansion killed Aunt Agatha. Agatha, the butler, and Charles live in Dreadbury Mansion, and are the only people who live therein. A killer always hates his victim, and is never richer than his victim. Charles hates no one that Aunt Agatha hates. Agatha hates everyone except the butler. The butler hates everyone not richer than Aunt Agatha. The butler hates everyone Aunt Agatha hates. No one hates everyone. Agatha is not the butler.

Who killed Aunt Agatha?

```
def aunt_agatha_hypotheses : List F0Form := [
  fo!{∃ x. lives_at_dreadbury(%x) ∧ killed(%x, agatha)},
  fo!{∀ x. lives_at_dreadbury(%x) ↔ (%x = agatha ∨ %x = butler ∨ %x = charles)},
  fo!{∀ x. ∀ y. killed(%x, %y) → hates(%x, %y)},
  fo!{∀ x. ∀ y. killed(%x, %y) → ¬ richer(%x, %y)},
  fo!{∀ x. hates(charles, %x) → ¬ hates(agatha, %x)},
  fo!{∀ x. ¬ hates(agatha, %x) ↔ %x = butler},
  fo!{∀ x. ¬ richer(%x, agatha) → hates(butler, %x)},
  fo!{∀ x. hates(agatha, %x) → hates(butler, %x)},
  fo!{∀ x. ∃ y. ¬ hates(%x, %y)},
  fo!{¬ agatha = butler}
]
```

# Using FOL Theorem Provers: Aunt Agatha in Lean

Examples/using\_first\_order\_theorem\_provers/aunt\_agatha.lean

```
def aunt_agatha_hypotheses : List F0Form := [
  fo!{∃ x. lives_at_dreadbury(%x) ∧ killed(%x, agatha)},
  fo!{∀ x. lives_at_dreadbury(%x) ↔ (%x = agatha ∨ %x = butler ∨ %x = charles)},
  fo!{∀ x. ∀ y. killed(%x, %y) → hates(%x, %y)},
  fo!{∀ x. ∀ y. killed(%x, %y) → ¬ richer(%x, %y)},
  fo!{∀ x. hates(charles, %x) → ¬ hates(agatha, %x)},
  fo!{∀ x. ¬ hates(agatha, %x) ↔ %x = butler},
  fo!{∀ x. ¬ richer(%x, agatha) → hates(butler, %x)},
  fo!{∀ x. hates(agatha, %x) → hates(butler, %x)},
  fo!{∀ x. ∃ y. ¬ hates(%x, %y)},
  fo!{¬ agatha = butler}
]
```

Did the butler kill aunt Agatha?

## Using FOL Theorem Provers: Raymond Smullyan

In an article called “The Asylum of Doctor Tarr and Professor Fether,” Raymond Smullyan tells of an investigation of 11 insane asylums by Inspector Craig of Scotland Yard.

In each of these asylums, every inhabitant  $x$  is either a **doctor** ( $\text{Doctor}(x)$ ) or a **patient** ( $\neg\text{Doctor}(x)$ ), and every inhabitant is either **sane** ( $\text{Sane}(x)$ ) or **insane** ( $\neg\text{Sane}(x)$ ).

The sane inhabitants are totally sane and the insane inhabitants are totally insane, in the following sense:  
for any **proposition**  $P$ ,  
a sane inhabitant **believes**  $P$  if and only if  $P$  is true,  
and an insane inhabitant believes  $P$  if and only if  $P$  is false.

## Using FOL Theorem Provers: The 8th Asylum

1. Some of the inhabitants are **teachers** of other inhabitants.  
Each inhabitant has at least one teacher.

## Using FOL Theorem Provers: The 8th Asylum

1. Some of the inhabitants are **teachers** of other inhabitants.  
Each inhabitant has at least one teacher.  
 $\forall x. \exists y. \text{Teaches}(y, x)$
2. There is one inhabitant who **trusts** all the **patients** but  
does not trust any of the **doctors**.

## Using FOL Theorem Provers: The 8th Asylum

1. Some of the inhabitants are **teachers** of other inhabitants.  
Each inhabitant has at least one teacher.  
 $\forall x. \exists y. \text{Teaches}(y, x)$
2. There is one inhabitant who **trusts** all the **patients** but does not trust any of the **doctors**.  
 $\exists x. \forall y. \neg \text{Doctor}(y) \leftrightarrow \text{Trusts}(x, y)$
3. No inhabitant  $x$  is willing to be a **teacher** of an inhabitant  $y$  unless  $x$  **believes** that  $y$  **trusts** himself.

## Using FOL Theorem Provers: The 8th Asylum

1. Some of the inhabitants are **teachers** of other inhabitants.  
Each inhabitant has at least one teacher.  
 $\forall x. \exists y. \text{Teaches}(y, x)$
2. There is one inhabitant who **trusts** all the **patients** but does not trust any of the **doctors**.  
 $\exists x. \forall y. \neg \text{Doctor}(y) \leftrightarrow \text{Trusts}(x, y)$
3. No inhabitant  $x$  is willing to be a **teacher** of an inhabitant  $y$  unless  $x$  **believes** that  $y$  **trusts** himself.  
 $\forall x, y. \text{Teaches}(x, y) \rightarrow (\text{Sane}(x) \leftrightarrow \text{Trusts}(y, y))$
4. For any inhabitant  $x$  there is an inhabitant  $y$  who **trusts** all and only those inhabitants who have at least one **teacher** who is **trusted** by  $x$ .

# Using FOL Theorem Provers: The 8th Asylum

1. Some of the inhabitants are **teachers** of other inhabitants.  
Each inhabitant has at least one teacher.

$$\forall x. \exists y. \text{Teaches}(y, x)$$

2. There is one inhabitant who **trusts** all the **patients** but does not trust any of the **doctors**.

$$\exists x. \forall y. \neg \text{Doctor}(y) \leftrightarrow \text{Trusts}(x, y)$$

3. No inhabitant  $x$  is willing to be a **teacher** of an inhabitant  $y$  unless  $x$  **believes** that  $y$  **trusts** himself.

$$\forall x, y. \text{Teaches}(x, y) \rightarrow (\text{Sane}(x) \leftrightarrow \text{Trusts}(y, y))$$

4. For any inhabitant  $x$  there is an inhabitant  $y$  who **trusts** all and only those inhabitants who have at least one **teacher** who is **trusted** by  $x$ .

$$\forall x. \exists y. \forall z. \text{Trusts}(y, z) \leftrightarrow \exists w. \text{Teaches}(w, z) \wedge \text{Trusts}(x, w)$$

# Using FOL Theorem Provers: The 8th Asylum in Lean

Examples/using\_first\_order\_theorem\_provers/asylum\_eight.lean

```
def asylum_eight_hypotheses : List F0Form := [
  fo!{∀ x. ∃ y. Teaches(%y, %x)},
  fo!{∀ x. ∀ y. Teaches(%x, %y) → (Sane(%x) ↔ Trusts(%y, %y))},
  fo!{∀ x. ∃ y. ∀ z. Trusts(%y, %z) ↔ ∃ w. Teaches(%w, %z) ∧ Trusts(%x, %w)},
  fo!{∃ x. ∀ y. ¬ Doctor(%y) ↔ Trusts(%x, %y)}
]
```

```
def asylum_eight_conclusion :=
  fo!{∃ x. Doctor(%x) ↔ ¬ Sane(%x)}

#eval (do
  discard <| callVampireTptp asylum_eight_hypotheses asylum_eight_conclusion
  (verbose := true)
  : IO Unit)
```

## Try at Home

Refute the 8th Asylum using resolution

- ▶ Skolemize the formula
- ▶ Turn it into CNF
- ▶ Negate the conclusion

Encode the Last Asylum in Lean

- ▶ The description is the Section 17.3 of the textbook
- ▶ Solve it using Vampire