

## AAAI 2018 Tutorial Integrating Learning into Reasoning

Brendan Juba      Loizos Michael  
Washington Univ. St. Louis      Open University of Cyprus  
[bjuba@wustl.edu](mailto:bjuba@wustl.edu)      [loizos@ouc.ac.cy](mailto:loizos@ouc.ac.cy)  
[www.cse.wustl.edu/~bjuba](http://www.cse.wustl.edu/~bjuba)      [cognition.ouc.ac.cy/loizos](http://cognition.ouc.ac.cy/loizos)

B. Juba supported by an AFOSR Young Investigator Award and NSF Award CCF-1718380

AAAI 2018 (New Orleans, Louisiana, U.S.A.)      February 03, 2018

## It Has Been Said That...



*"For every belief comes either  
through **syllogism** or from **induction**"*

Aristotle, Organon (Prior Analytics II, §23)

Artificial Intelligence research today?  
**Syllogism XOR Induction**

2

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

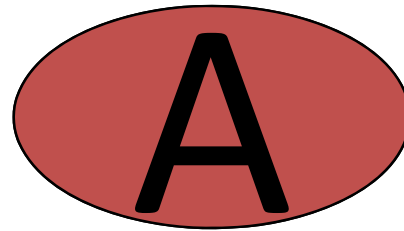
## High-Level Tutorial Roadmap

- A. Why KRR Should Embrace Learning
- B. Introduction to PAC-Semantics
- C. Integrating Deduction and Induction
- D. Reasoning Non-Monotonically
- E. Overall Summary and Conclusions

3

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



4

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## High-Level Tutorial Roadmap

- A. **Why KRR Should Embrace Learning**
- B. Introduction to PAC-Semantics
- C. Integrating Deduction and Induction
- D. Reasoning Non-Monotonically
- E. Overall Summary and Conclusions

5

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

- A. **Why KRR Should Embrace Learning**

Review of Standard KRR Problems

6

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Representation and Reasoning

- Propositions ( $p, q, \dots$ ). Connectives ( $\wedge, \neg, \dots$ ).
  - Implications:  $\varphi \Rightarrow x$ . Equivalences:  $\varphi \Leftrightarrow x$ .
- Reasoning *semantics* through entailment  $\models$ .
- Proof procedures  $\vdash$  to *compute* entailment.
- Given formulas in *KB* and an input *O*, *deduce* whether a result *R* is entailed ( $KB \cup O \models R$ ).
- Given formulas in *KB* and an input *O*, *abduce* an explanation *E* that entails *O* ( $KB \cup E \models O$ ).

7

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Relational Representation and IQEs

- Predicates ( $p, q, \dots$ ). Variables ( $x, y, \dots$ ). **Tokens ( $t_i$ )**. Connectives ( $\wedge, \neg, \dots$ ), Quantifiers ( $\forall, \exists$ ).
- From the class of implications / equivalences, consider only those whose body comprises **independently quantified expressions (IQEs)**.
  - $\exists y \exists z [ num(y) \wedge num(z) \wedge larger(z,y) \wedge \neg div(y,z) ]$
  - No tokens (they carry no meaning), small arity.
  - $\forall xs [ formula\ over\ IQEs \Leftrightarrow head\_predicate(xs) ]$  (see later: these restrictions support learnability)

8

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Non-Monotonic Reasoning

- Non-monotonicity typically viewed as property of **extending input *O* for fixed *KB***, and having result *R* become “smaller”.
- Useful also when **extending *KB*** as a prerequisite to elaboration tolerance.
- Will use logic-based argumentation for NMR.
  - Most (all?) major NMR logics have been reformulated in terms of argumentation.
  - Compatible with human cognition [Kakas+ '16].

9

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Formal Argumentation in AI

- Abstract Argumentation framework  $\langle Arg, Att \rangle$ 
  - Arg* is a set of arguments (no internal structure)
  - Att* is a binary relation on *Arg* (lifted on sets)
- Goal: Find  $S \subseteq Arg$  that defends all its attacks.
  - Several ways to make this precise [Dung '95].
- Structured argumentation (e.g., ABA, ASPIC+): argument is a classical proof from inputs and KB rules. The overall reasoning is not classical.

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Activity: Logic-Based Arguments

treated **implies** will\_survive.  
 fatal **implies**  $\neg$ will\_survive.  
 viral\_meningitis **implies** meningitis.  
 bacterial\_meningitis **implies** meningitis.  
 bacterial\_meningitis **implies** fatal.  
 fatal **implies**  $\neg$ treatable.  
 meningitis **implies**  $\neg$ fatal.  
 meningitis **implies** treatable.  
 $\neg$ fatal **implies** will\_survive.  
 true **implies**  $\neg$ meningitis.

**total ordering (for simplicity)**  
 $\uparrow$   
**more preferred**

**What inference should we draw on each input?**  
 {}, {viral\_meningitis}, {bacterial\_meningitis}, {bacterial\_meningitis,treated}

11

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning in Temporal Settings

- Frame Problem** (commonsense law of inertia): Properties persist unless caused to stop.
  - If you see a bird flying, it is flying a bit later.
  - Persistence rules are weaker than causal rules.
- Ramification Problem**: Production of indirect effects as a way to satisfy state constraints.
  - If you shoot a bird, it stops being able to fly.
  - Encode ramifications as causal rules (since constraints could also qualify causal change).

12

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning in Temporal Settings

- **Qualification Problem:** Effects of actions are blocked if they would violate state constraints.
  - If you scare a dead bird, it does not fly away.
  - This constraint does not produce ramifications.
  - Encode constraints as preclusion / block rules.
- **State Default Problem [Kakas+ '08]:** If a state constraint is violated, the exception persists.
  - If you see a flying penguin, it remains one.
  - Persistence rules are stronger than constraints.

13

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Activity: Story Understanding

What inferences follow from this story?

*Papa Joe woke up early at dawn, and went off to the forest. He walked for hours, until the sight of a turkey in the distance made him stop. A bird on a tree nearby was cheerfully chirping away, building its nest. He carefully aimed at the turkey, and pulled the trigger of his shotgun. Undisturbed, the bird nearby continued chirping.*

Q1: What is the condition of the turkey?

- (a) Alive and unharmed. (b) Dead. (c) Injured.

14

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## ... Through Argumentation

scene 1: forest( $t_1$ ) at 1  
 scene 2: bird( $t_3$ ) at 4.  
 gun(f) at 7.  
 tokens have no meaning  
 Uniform treatment for causality and change.  
 Addition of time rule bodies are IQEs  
 RAC: inertia, endogenous arguments  
 Try it online! <http://cognition.ouc.ac.cy/star>  
 [Diakidoy+ '14,15]

15

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## A. Why KRR Should Embrace Learning

### The Task of Knowledge Acquisition

16

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## The Declaration of Independence (of KRR from Machine Learning)

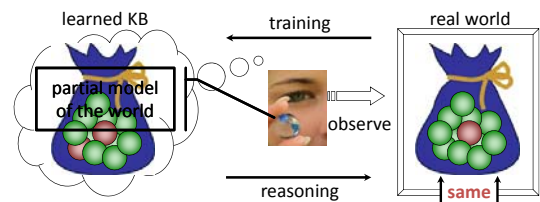
*We hold these truths to be self-evident:  
 an appropriate knowledge base is given  
 reasoning is evaluated against the KB  
 chaining of rules is trivially beneficial  
 KB rules can be applied in any order  
 acquisition of KB can be done a priori*

17

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## What if KB not Given but Learned?



- Then, KB is **Probably Approximately Correct**:
  - improbable the samples were unrepresentative
  - future predictions will be approximately correct

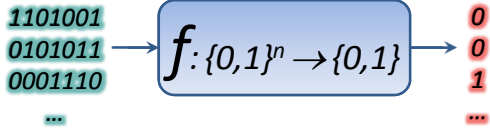
18

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Program, Memorize, or Induce?

• Examples: Target Concept: Labels:



- **Concept Class C:** bias, possible target concepts.
- **Hypothesis Class H:** all acceptable hypotheses.
- **Complexity:** examples/time/memory to learn.

19

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Typical Boolean Target Concepts

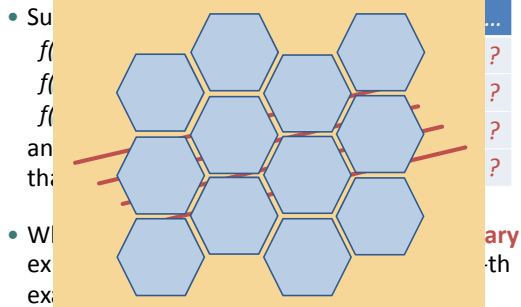
- **Conjunction function:**  $x_3 \wedge x_5 \wedge x_7$ 
  - **1011101** the conjunction of  $x_3, x_5, x_7$  is **1**
- **Parity function:**  $x_2 \oplus x_3 \oplus x_5 \oplus x_6$ 
  - **0010101** the parity of  $x_2, x_3, x_5, x_6$  is **0**
- **Decision list:**  $\langle x_4, 1 \rangle, \langle x_2, 0 \rangle, \langle x_7, 1 \rangle, \langle \text{default}, 0 \rangle$ 
  - **1010011** false:  $x_4, x_2$  first true:  $x_7$  value: **1**
- **Linear threshold** (hyperplane, perceptron):  
 $\text{weights} = \langle 0.3, 0.1, 0.5, 0.2, 0.7, 0.9 \rangle, \text{threshold} = 1$ 
  - **0110100**  $110100 \cdot \text{weights} = 0.6 < \text{threshold}$

20

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Activity: Is Learning Even Possible?



21

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Evaluation of Reasoning Process

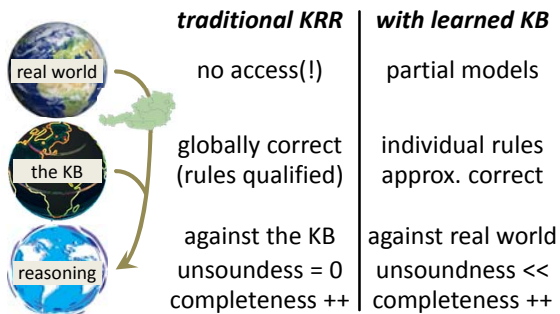
- Evaluate reasoning process against given KB:
  - improve completeness, insisting on **full soundness**
  - **okay** since KB is considered the **golden standard**
  - **not okay** when KB is only *approximately* correct
- Evaluate reasoning process when KB *learned*:
  - improve completeness, without compromising soundness much more than what is *necessary*
  - soundness and completeness wrt an **"ideal KB"** (access only to its partial models during training)

22

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Real World, the KB, and Reasoning



23

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

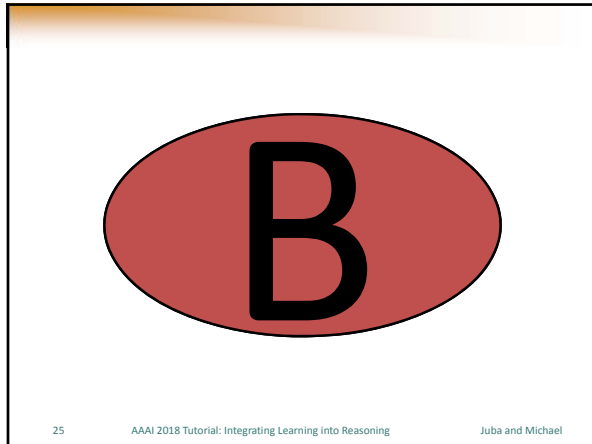
## Is Chaining *Trivially* Beneficial?

- **Yes**, if rules in the KB are given / programmed.
  - Given: *"if you have fever then you are sick"*
  - Given: *"if you are sick then visit a doctor"*
  - Infer "visit a doctor" given "you have fever"?
- **Prima facie no**, if rules in the KB are learned.
  - Learning can render rule chaining superfluous (cf. shortcuts, heuristics, fast thinking, hunch).
  - Learned: *"if you have fever then visit a doctor"*

24

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



## High-Level Tutorial Roadmap

- A. Why KRR Should Embrace Learning
- B. Introduction to PAC-Semantics**
- C. Integrating Deduction and Induction
- D. Reasoning Non-Monotonically
- E. Overall Summary and Conclusions

26 AAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

## B. Introduction to PAC-Semantics

### Semantics for Propositional Logic

27 AAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Motivating example: birds.com Analytics

birds.com™

The data

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
107	54	Mouse
107	55	Grubs
⋮	⋮	⋮

~[FOOD=FISH] ?

To determine: "true... with high probability for data source?"

28 AAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

- ## PAC-semantics for propositional logics
- Fixed, finite set of propositional variables: [FOOD=SEED], [FOOD=GRUBS], [FOOD=MOUSE], [FOOD=WORM], [FOOD=FISH], ..., [FLIES], [SWIMS], [HAS\_BILL], [HAS\_BEAK], [COLOR=BLACK], [COLOR=RED], [COLOR=BROWN], [COLOR=WHITE], [COLOR=GREEN], [COLOR=BLUE], ...
  - Probability distribution D over Boolean valuations for the propositional variables
    - **NOTE:** generally not uniform, not independent
- 29 AAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

- ## PAC-semantics for propositional logics
- Probability distribution D over Boolean valuations for the propositional variables
    - Usually propositional variables capture attributes of data entry, sensor values, etc.
    - D captures range of possible combinations of values and their relative frequency
- Definition.** [Valiant '00] A formula  $\varphi(x_1, \dots, x_n)$  is  $(1-\epsilon)$ -valid under D if  $\Pr_D[\varphi(x_1, \dots, x_n)] \geq 1-\epsilon$ .
- $\epsilon$  may or may not be small (cf. [Adams '75])
- 30 AAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

## PAC-semantics for propositional logics

**Definition.** [Valiant '00] A formula  $\varphi(x_1, \dots, x_n)$  is  $(1-\epsilon)$ -valid under  $D$  if  $\Pr_D[\varphi(x_1, \dots, x_n)] \geq 1-\epsilon$ .

- $\epsilon$  may or may not be small (cf. [Adams '75])

WHAT INFERENCE PRESERVE  $(1-\epsilon)$ -VALIDITY?

**Theorem.** [Juba '13] Suppose  $h_1, \dots, h_k$  are, respectively,  $1-\epsilon_1, \dots, 1-\epsilon_k$ -valid, and  $\{h_1, \dots, h_k\} \models \varphi$  (classical entailment). Then  $\varphi$  is  $1-(\epsilon_1 + \dots + \epsilon_k)$ -valid.



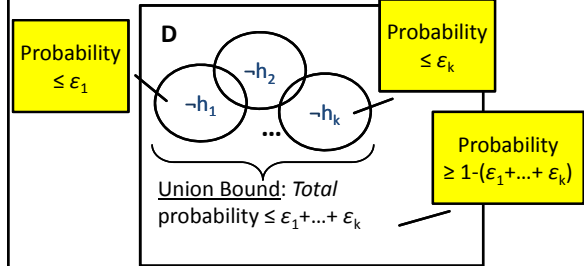
31

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Classical entailment to PAC entailment: the union bound

Consider the events  $[-h_1], \dots, [-h_k]$



32

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Classical entailment to PAC entailment: the union bound

- Consider the events  $[-h_1], \dots, [-h_k]$
- Union bound:** the event  $[-h_1 \vee \dots \vee -h_k] = \neg(h_1 \wedge \dots \wedge h_k)$  has total probability  $\leq \epsilon_1 + \dots + \epsilon_k$
- Classical entailment:**  $[h_1 \wedge \dots \wedge h_k] \subseteq [\varphi]$
- Therefore:  $\Pr_D[\varphi] \geq \Pr_D[h_1 \wedge \dots \wedge h_k] \geq 1 - (\epsilon_1 + \dots + \epsilon_k)$  ■
- Summary:** all classical inferences preserved, but each additional premise may incur a cost

33

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## PAC-semantics does not mean Probability Logic (e.g. Nilsson'86)

- |                                                                                                                                                                                                                                                                                                                             |                                                                                                                                                                                                                                                                                                                         |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>PAC-semantics</b></p> <ul style="list-style-type: none"> <li>Classical Boolean in object language</li> <li>Probability bounds in interpretation             <ul style="list-style-type: none"> <li>Classical proof of a formula guarantees it holds with some probability under <math>D</math></li> </ul> </li> </ul> | <p><b>Logics of Probability</b></p> <ul style="list-style-type: none"> <li>Probability bounds in object language</li> <li>Classical (Tarskian) semantics             <ul style="list-style-type: none"> <li>Classical proof of a probability bound on <math>D</math> that is true with certainty</li> </ul> </li> </ul> |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

34

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Motivating example:

### birds.com Analytics

birds.com™

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
107	54	Mouse
107	55	Grubs
⋮	⋮	⋮

$\neg[\text{FOOD}=\text{FISH}]$

Seems  $\approx 7/8$ -valid...

...JUSTIFIED ON WHAT GROUNDS??

35

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## B. Introduction to PAC-Semantics

### Learning and Inductive Inference

36

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Knowledge acquisition from data: birds.com Analytics

birds.com™

The data

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
107	54	Mouse
107	55	Grubs
⋮	⋮	⋮

To determine:  
(1-ε)-valid for D?  
(D is arbitrary)

Assume: Each entry (row) is an "example," drawn independently from D

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

### The "i.i.d. data" assumption

- **Assume:** data consists of *examples* – valuations of the propositional variables drawn independently from common D ("i.i.d.")
  - Will enable us to draw conclusions about D
- **Recall:** if A,B,C,... are (mutually) *independent*, then  $\Pr[A \wedge B \wedge C \wedge \dots] = \Pr[A] \Pr[B] \Pr[C] \dots$
- Likewise, if *random variables* W,X,Y,... are *independent* then  $E[WXYZ\dots] = E[W] E[X] E[Y] \dots$  (and  $E[f(W)g(X)h(Y)\dots] = E[f(W)]E[g(X)]E[h(Y)]\dots$ )

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

### Inference from i.i.d. examples

- **Assume:** data consists of *examples* – valuations of the propositional variables drawn independently from common D ("i.i.d.")
  - Will enable us to draw conclusions about D
- Suppose  $\phi$  is not (1-ε)-valid under D.
- Draw  $X^{(1)} = (X_1^{(1)}, X_2^{(1)}, \dots, X_n^{(1)})$ ,  
 $X^{(2)} = (X_1^{(2)}, X_2^{(2)}, \dots, X_n^{(2)})$ ,  
 $\dots$ ,  
 $X^{(m)} = (X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)})$  } Data set of m "examples" independently from D.

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

### Inference from i.i.d. examples

- Suppose  $\phi$  is not (1-ε)-valid under D.
- Draw  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$  independently from D.
- $\Pr[\phi(X^{(1)}), \phi(X^{(2)}), \dots, \text{and } \phi(X^{(m)})]$   
 $= \Pr[\phi(X^{(1)})] \Pr[\phi(X^{(2)})] \dots \Pr[\phi(X^{(m)})]$
- *By hypothesis*, each  $\Pr[\phi(X^{(i)})] \leq 1 - \epsilon$
- So, the probability that we fail to observe that  $\phi$  is false for some example  $X^{(i)}$  is at most  $(1 - \epsilon)^m \leq e^{-\epsilon m}$  (...since for all x,  $1 + x \leq e^x$ )
- Less than any given  $\delta$  if  $m \geq \frac{1}{\epsilon} \ln \frac{1}{\delta}$

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

### Inference from i.i.d. examples

- So we have shown  
**Theorem.** If  $\phi$  is consistent with  $m \geq \frac{1}{\epsilon} \ln \frac{1}{\delta}$  examples drawn independently from D, then with probability  $1 - \delta$ ,  $\phi$  is (1-ε)-valid under D. ("Probably Approximately Correct")
- Only guaranteed to work if  $\phi$  is actually always true. What if  $\phi$  is only (1-ε')-valid under D, for some  $\epsilon' < \epsilon$ ? (e.g., only  $\frac{7}{8}$ -valid, in the case of birds.com?)

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

### Chernoff/Hoeffding bounds: sample averages are good estimates

- Let  $Y^{(1)}, Y^{(2)}, \dots, Y^{(m)}$  be independent random variables taking values in  $[0, 1]$ . Let  $\mu = E[(\frac{1}{m})(Y^{(1)} + Y^{(2)} + \dots + Y^{(m)})]$ .
- **Hoeffding bound:** for any  $\gamma > 0$ ,  
 $\Pr[(\frac{1}{m})(Y^{(1)} + Y^{(2)} + \dots + Y^{(m)}) > \mu + \gamma] < e^{-2m\gamma^2}$   
 $\Pr[(\frac{1}{m})(Y^{(1)} + Y^{(2)} + \dots + Y^{(m)}) < \mu - \gamma] < e^{-2m\gamma^2}$
- **Chernoff bound:** for any  $\gamma > 0$ ,  
 $\Pr[(\frac{1}{m})(Y^{(1)} + Y^{(2)} + \dots + Y^{(m)}) > (1 + \gamma)\mu] < e^{-m\gamma^2/3}$   
 $\Pr[(\frac{1}{m})(Y^{(1)} + Y^{(2)} + \dots + Y^{(m)}) < (1 - \gamma)\mu] < e^{-m\gamma^2/2}$

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Inference using Chernoff/Hoeffding

- We draw  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$  independently from  $D$
- For  $p = \frac{1}{m}([\varphi(X^{(1)})] + [\varphi(X^{(2)})] + \dots + [\varphi(X^{(m)})])$ , how large must  $m$  be to conclude that with probability  $1-\delta$ ,  $\varphi$  is  $(p \pm \gamma)$ -valid under  $D$ ?
- Use the Hoeffding bound: for any  $\gamma > 0$ ,  
 $\Pr[\frac{1}{m}(\gamma^{(1)} + \gamma^{(2)} + \dots + \gamma^{(m)}) > \mu + \gamma] < e^{-2m\gamma^2}$   
 $\Pr[\frac{1}{m}(\gamma^{(1)} + \gamma^{(2)} + \dots + \gamma^{(m)}) < \mu - \gamma] < e^{-2m\gamma^2}$

**TRY IT!**

43

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Inference using Chernoff/Hoeffding

- We draw  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$  independently from  $D$
- For  $p = \frac{1}{m}([\varphi(X^{(1)})] + [\varphi(X^{(2)})] + \dots + [\varphi(X^{(m)})])$ , how large must  $m$  be to conclude that with probability  $1-\delta$ ,  $\varphi$  is  $(p \pm \gamma)$ -valid under  $D$ ?
- Use the Hoeffding bound: for any  $\gamma > 0$ ,  
 $\Pr[\frac{1}{m}(\gamma^{(1)} + \gamma^{(2)} + \dots + \gamma^{(m)}) > \mu + \gamma] < e^{-2m\gamma^2}$   
 $\Pr[\frac{1}{m}(\gamma^{(1)} + \gamma^{(2)} + \dots + \gamma^{(m)}) < \mu - \gamma] < e^{-2m\gamma^2}$
- For  $m \geq \frac{1}{2\gamma^2} \ln \frac{2}{\delta}$ , can check that  $e^{-2m\gamma^2} \leq \frac{\delta}{2}$ ; take a union bound of upper and lower bounds to conclude  $p$  is within  $\pm\gamma$  of  $\Pr_D[\varphi(X)]$ .

44

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Comparison: PAC-semantics versus Inductive Logic Programming

### PAC-semantics

- Examples drawn from larger distribution  $D$ 
  - $D$  mostly unseen
- Rules partially capture  $D$
- Rules can be (a little) inconsistent with examples

### Inductive Logic Programming

- Examples define domain
  - Closed-world assumption
- Rules fully capture domain
- Rules must be faithful to defining examples

45

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Occam's Razor [Blumer et al. '87]

- **Theorem.** Let  $\mathcal{H}$  be a set of formulas that can be specified using at most  $B$  bits. Suppose we draw  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$  independently from  $D$  for  $m \geq \frac{1}{2\gamma^2} ((B+1)\ln 2 + \ln \frac{2}{\delta})$ . With probability  $1-\delta$ ,  $\frac{1}{m}([\mathcal{h}(X^{(1)})] + [\mathcal{h}(X^{(2)})] + \dots + [\mathcal{h}(X^{(m)})])$  is within  $\pm\gamma$  of  $\Pr_D[\mathcal{h}(X)]$  for every  $\mathcal{h}$  in  $\mathcal{H}$ .
- We know:  $\frac{1}{m}([\mathcal{h}(X^{(1)})] + [\mathcal{h}(X^{(2)})] + \dots + [\mathcal{h}(X^{(m)})])$  is within  $\pm\gamma$  of  $\Pr_D[\mathcal{h}(X)]$  for any *single*  $\mathcal{h}$  in  $\mathcal{H}$  with probability  $1-\delta/2^{B+1}$ .
- There are  $< 2^{B+1}$  strings of  $B$  bits, so there are fewer than  $2^{B+1}$   $\mathcal{h}$  in  $\mathcal{H}$ . We take a union bound. ■

46

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## (Simplified) Learning to reason using Occam's Razor [Khaddon-Roth '97]

- *Recall:* Deciding if a 3-DNF (ORs of ANDs of at most three literals) is valid is NP-complete.
- There are  $2^{O(n^3)}$  3-DNFs.
- **Theorem.** Suppose we draw  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$  independently from  $D$  for  $m = O(\frac{1}{\gamma^2}(n^3 + \ln \frac{1}{\delta}))$ . Then with probability  $1-\delta$ ,  $\frac{1}{m}([\mathcal{h}(X^{(1)})] + [\mathcal{h}(X^{(2)})] + \dots + [\mathcal{h}(X^{(m)})])$  is within  $\pm\gamma$  of  $\Pr_D[\mathcal{h}(X)]$  for every 3-DNF  $\mathcal{h}$ .
- In particular, if all  $\mathcal{h}(X^{(i)})=1$ , we guarantee  $\mathcal{h}$  is at least  $(1-\gamma)$ -valid with probability  $1-\delta$ .

47

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Learning a chaining rule using Elimination [Valiant '84]

- Suppose  $x_t$  is defined by a *conjunction*, i.e.,  $x_t \Leftrightarrow \ell_1 \wedge \ell_2 \wedge \dots \wedge \ell_k$  where  $\ell_1, \ell_2, \dots, \ell_k$  are literals.
- There are only  $2^{2n}$  conjunctions of literals on  $n$  propositional variables.
- **Occam's Razor:** any rule  $x_t \Leftrightarrow \ell_1 \wedge \ell_2 \wedge \dots \wedge \ell_k$  that is consistent with  $m = O(\frac{1}{\gamma^2}(n + \ln \frac{1}{\delta}))$  examples is  $(1-\gamma)$ -valid with probability  $1-\delta$ .
- We only need to find such a rule.

48

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



### Learning a chaining rule using Elimination [Valiant '84]

- We only need to find  $x_t \Leftrightarrow \ell_1 \wedge \ell_2 \wedge \dots \wedge \ell_k$  consistent with  $m = O(1/\gamma^2(n + \ln^2/\delta))$  examples
- Elimination**: for each  $i$ th example, if  $X_t^{(i)} = 1$ , then any literal  $\ell$  with  $\ell(X^{(i)}) = 0$  cannot be included in the rule. Delete it.
- The rule given by the conjunction of the remaining literals must contain the actual defining conjunction  $\ell_1 \wedge \ell_2 \wedge \dots \wedge \ell_k$ .
- Therefore, we find a conjunction that is consistent with all  $m$  examples, as needed. ■

49

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Learning a chaining rule using Elimination [Valiant '84]

- We find a  $(1-\gamma)$ -**valid** rule using  $O(1/\gamma^2(n + \ln^2/\delta))$  examples by taking the conjunction of all literals for which whenever  $\ell(X^{(i)}) = 0$ ,  $X_t^{(i)} = 0$ .
- As stated, the algorithm runs in time  $O(n/\gamma^2(n + \ln^2/\delta))$  so this is efficient.
  - Note*: actually,  $O(1/\gamma(n + \ln^2/\delta))$  examples will do

50

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### B. Introduction to PAC-Semantics

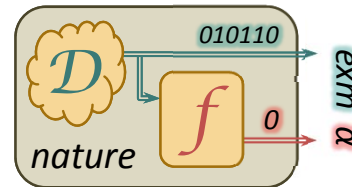
#### Coping with Partial Information

51

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### PAC Learning (with Complete Info)



Given access to  $(exm, \alpha)$ 's drawn during training, w.h.p. and efficiently produce  $h$  s.t. on  $(exm', \alpha')$  drawn during testing,  $h(exm') \neq \alpha'$  w.p. at most  $\epsilon$ .

52

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Learning and Prediction Scenarios

- Construct hypothesis to predict target  $x_1$ .
  - 1101100011011 learning
  - ?011101001010 prediction
- What if we want to predict multiple  $x_j$ 's?
  - 1101100011011 learning
  - ?01?1??0010?0 prediction
- What if we want to learn **autonomously**?
  - 1?011?001?0?? learning
  - ?01?10?001010 prediction

53

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Complete the Missing Value

Marital Status	Age	Smoking Habits	Drinking Habits	Blood Pressure	Cholesterol Level	Exercising Habits
single			3 weekly	130/90		
married parent	49	cigarettes 10 daily			high	
	32		no drinking		normal	jogging daily
divorced			2 daily	145/100		
	68	cigars 3 daily		?		no exercise
single parent		no smoking		125/80	slightly elevated	
	25		rarely		normal	gym 3 weekly

54

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Structure in Missingness?

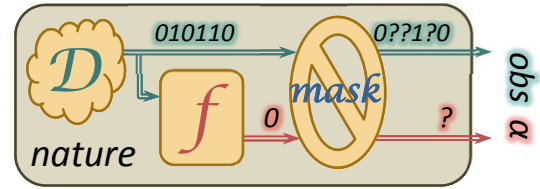
- Q: What is the number of your credit card?
  - The responder may not wish to share it.
  - What if the responder **does not have** one?
- Q: When was the last time you ate apples?
  - The responder may have a poor memory.
- Q: Where you ever convicted for murder?
  - The responder may not wish to answer...
  - ... especially **if the answer** would be “yes”!

55

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Autodidactic Learning (PAC + Missing)



Given access to  $(obs, \alpha)$ 's drawn during training, w.h.p. and efficiently produce  $h$  s.t. on  $(obs', \alpha')$  drawn during testing,  $h(obs') \neq \alpha'$  w.p. at most  $\epsilon$ .

56

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reading Between the Lines

[Michael '09]



57

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Learning the Hidden Reality

- Need learning algorithms able to:
  - identify structure in the hidden example,
  - given access only to partial observations,
  - without knowing how masking works.
- Autodidactic learning [Michael '08,10]:
  - Suffices to learn rules **consistent** with obs.
    - e.g., predict that  $x_1 = 0$  in  $?010??11$
  - W.h.p., these rules make **accurate** predictions.

58

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Evaluation of an Individual Rule

example	observation	evaluate rule $(X_2 \vee \neg X_3) \leftrightarrow X_1$		
		inference	consistent	accurate
0010110	0?10??0	1	Yes	Yes
0100100	?10??00	1	No	No
0100100	?10??0?	?	Yes	Abstain
0010110	0??01??	1	Yes	Yes
0100100	0??01??	1	Yes	No

**Theorem:** For each mask there is  $\eta$  s.t. each  $(1-\eta \cdot \epsilon)$ -consistent rule is  $(1-\epsilon)$ -accurate. The “discount” is tight.

59

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## When to Abstain from Predictions

- Choosing** to abstain would trivialize learning.
- Hypotheses are still Boolean functions...
- Abstain only if hypothesis not fully determined.
  - Assume hypothesis is  $h = x_1 \wedge \neg x_3 \wedge x_7$
  - If observation is  $x = \alpha 10?1101$  then  $h(x) = ?$
  - If observation is  $x = \alpha 10?1100$  then  $h(x) = 0$
- Does missing info “kill” our ability to predict?

60

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Activity: Predict the Hidden Animal



61

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Autodidactic Learnability [Michael '08,10]

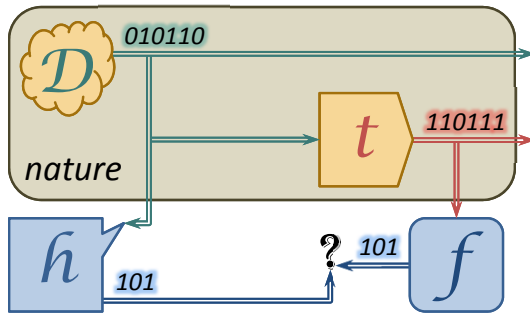
- **Theorem:** All classes of monotone and read-once formulas that are PAC-learnable, are also autodidactically learnable (*arbitrary* masking).
- **Learning algorithm:**
  - Ignore observations where the label is masked.
  - During training map  $a0??1?0$  to  $a0\alpha\alpha1\alpha0$ .
  - PAC-learn from resulting (complete) instances.
- **Theorem:** Parities and monotone term 1-decision lists are *not properly learnable*.

62

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Causal Learnability [Michael '11]



63

Introspective Forecasting

Loizos Michael (OUC)

## B. Introduction to PAC-Semantics

### Dealing with First-Order Expressions

64

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## An Example of Learned Knowledge

- NL query : “members share something”
- logic form :  $member(t_1), something(t_3), share(t_1, t_3)$
- knowledge : (trained on “spyware” web pages)
 

$file(x) \Leftrightarrow$	$threshold(1.0)$	% pos:1852 neg:1831
$\exists v : scan(v,x) \wedge rogue(v)$	$weight(0.962710)$	% pos:16 neg:1
$\exists v : share(v,x)$	$weight(1.627098)$	% pos:11 neg:1
$\exists v : have(x,v) \wedge program(v)$	$weight(0.645691)$	% pos:19 neg:0
$\exists v : open(v,x)$	$weight(1.593269)$	% pos:27 neg:2
- inference :  $file(t_3)$
- NL answer : “something is a file” [Michael '13]

65

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## How To Learn Relational Rules?

- Looking to learn a rule in the correct form:
  - No tokens (they carry no meaning), small arity.
  - $\forall xs [ formula\ over\ IQEs \Leftrightarrow head\_predicate(xs) ]$
- The formula belongs in a concept class that is known to be PAC learnable (cf. [Valiant '00]).
  - Linear thresholds with propositional features.
  - Recall: also learnable from partial observations.

*Relational Scenes  $\rightarrow$  Propositional Examples  $\rightarrow$  Propositional Learning  $\rightarrow$  Relational Hypothesis*

66

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Activity: Learn a Linear Threshold

- Initially assign weight 0.8 to every proposition.
- When a negative example is predicted true: divide by 2 weights of true propositions.
- When a positive example is predicted false: multiply by 2 weights of true propositions.
- Hidden target concept: 

000110, 011101, 101011, 001110, 110010,  
011001, 111111, 101011, 001101, 011101.
- Computed hypothesis:  $\langle 0.2, 0.1, 0.05, 0.8, 0.2 \rangle$

67

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Examples from Relational Scenes

- Instead of Boolean learning examples, we get **relational scenes** giving instances of relations:  $file(t_2), scan(t_7, t_2), rogue(t_7), \neg file(t_5), \neg open(t_5, t_2)$ .
- For a head predicate (e.g.,  $file(x)$ ), **consider all IQEs** that follow a schema (e.g.,  $scan \wedge rogue$ ):
 
$$scan(x, x) \wedge rogue(x), \quad \exists v : scan(x, v) \wedge rogue(v),$$

$$\exists v : scan(x, v) \wedge rogue(x), \quad \exists v : scan(x, v) \wedge rogue(v),$$

$$\exists v : scan(v, x) \wedge rogue(x), \quad \exists v : scan(v, x) \wedge rogue(v),$$

$$\exists u, v : scan(u, v) \wedge rogue(v), \quad \exists u, v : scan(u, v) \wedge rogue(u).$$

68

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Propositionalization of Learning

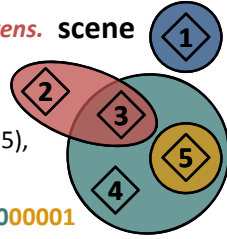
**Full info. Small arity. Few tokens. scene**

$P(1), P(2), \dots, P(5),$   
 $P(1,1), \dots, P(2,3), \dots, P(5,5),$   
 $P(1,1,1), \dots, P(3,5,4), \dots, P(5,5,5),$   
 $P(1), P(2), \dots, P(5).$

1000000000010...00...010...000001

$x=3 \quad -P(x)$  with  $\exists P(x, a, b), \exists P(a, x), \dots$

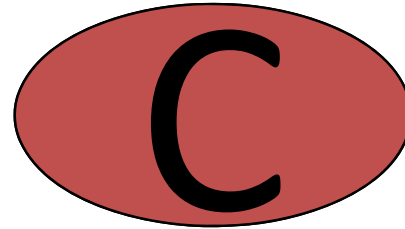
$x=5 \quad +P(x)$  with  $\exists P(a, x, b), \exists P(a, x, b) \wedge P(c, a), \dots$



69

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



70

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## High-Level Tutorial Roadmap

- Why KRR Should Embrace Learning
- Introduction to PAC-Semantics
- Integrating Deduction and Induction**
- Reasoning Non-Monotonically
- Overall Summary and Conclusions

71

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## C. Integrating Deduction and Induction

### Interleaving Learning and Reasoning

72

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## The Declaration of Independence (of KRR from Machine Learning)

*We hold these truths to be self-evident:  
an appropriate knowledge base is given  
reasoning is evaluated against the KB  
chaining of rules is trivially beneficial  
KB rules can be applied in any order  
acquisition of KB can be done a priori*

AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Reasoning ≈ Fill-In Missing Values

Marital Status	Age	Smoking Habits	Drinking Habits	Blood Pressure	Cholesterol Level	Exercising Habits
single			3 weekly	130/90		
married parent	49	cigarettes 10 daily			high	
	32		no drinking		normal	jogging daily
divorced			2 daily	145/100		
	68	cigars 3 daily		?		no exercise
single parent		no smoking		125/80	slightly elevated	
	25		rarely		normal	gym 3 weekly

74      AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## KB Performance Evaluation

- Example knowledge base:
  - $(x_3 \text{ and not } x_5)$  determines the value of  $x_2$
  - $(\text{not } x_4 \text{ or not } x_1)$  determines the value of  $x_6$

example	observation	KB inference	KB evaluation
101011	1?101?	101011	no disagreement no "don't know"
001001	0?100?	011001	disagreement no "don't know"
010101	0??101	0??101	no disagreement "don't know"
010001	110000	110001	disagreement no "don't know"

50% sound  
75% complete

75      AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Multiple Targets and Autonomy

- Predict multiple  $x_j$ 's and learn autonomously.
  - 1?011?001?0??      learning
  - ?01?10?001010      prediction
- Options on how to tackle this task:
  - Learn hypotheses first, then apply in parallel.
  - Learn hypotheses first, then apply by chaining.
  - Simultaneously learn hypotheses and predict.
- Which of these approaches is appropriate?

76      AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Basic Hypotheses / Predictors

- Given a specified attribute  $x_i$  of the world, and sufficient resources during training,
- produce a PAC **predictor** (=rule)  $P_j$  for  $x_i$ 
  - e.g.,  $x_7 \equiv x_1 \wedge \neg x_3$  is a predictor for  $x_7$
- Predictor  $P_j$  is "classical" and could potentially
  - abstain, when its body is not fully determined
  - predict incorrectly (cf. approximately correct)
  - but generally improve completeness of the DB

77      AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Reasoning Process Semantics

- $exm = \langle -1, 0, 5, 0, 4, -7, 9 \rangle \in S$
- $obs = \langle -1, 0, *, 0, 4, *, 9 \rangle \in S^*$

endogenous and exogenous qualification (hints of NMR?)

Policy  $P$  determines how predict

$(P_1) x_3 \equiv 1x_5 + 1x_7$        $(P_2) x_4 \equiv x_2 \cdot x_6$

e.g.,  $P = \langle \{P_1, P_2, P_3\} \rangle$       e.g.,  $P = \langle \{P_1, P_2\}, \{P_3\} \rangle$

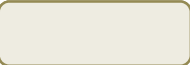
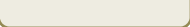
"flat/parallel" policy      "chaining" of predictors

$P(obs) = \langle -1, 0, 5, 0, 4, *, 9 \rangle$        $P(obs) = \langle -1, 0, 5, 0, 4, 3, 9 \rangle$

sound against  $exm$  but **incomplete**      **unsound** against  $exm$  but complete

78      AAAI 2018 Tutorial: Integrating Learning into Reasoning      Juba and Michael

## Activity: Chaining & Completeness

- Example knowledge base:
  - (noon time)  $\Leftrightarrow$  lunch time
  - (lunch time and Bob hungry)  $\Leftrightarrow$  Bob eats
- Statement:
  - "It was not noon time yet, but Bob was hungry."
- Multi-layered application: 
- Single-layered application: 
- Chaining is better than **any** single-layered KB.

79

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Chaining is *Provably* Beneficial

- Definition:** *Chaining collapses* on  $S^*$  against  $S$  if
- for **every** policy  $P$  with some **performance** (soundness + completeness) on  $S^*$  against  $S$ ,
  - there exists a **flat** policy  $P'$  (not a reordering of  $P$ , necessarily) with equally high performance.

**Theorem:** [Michael '14] There exist  $S^*$ ,  $S$  such that chaining does not collapse on  $S^*$  against  $S$ .

*Proof:* Chaining can **simulate non-monotonic reasoning**, which is beyond individual predictors.

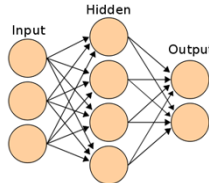
80

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Analogy to Neural Networks

There are multilayered NNs that compute more complex functions than those by **any** NN without hidden layers.



*Trivially, because of larger hypothesis space.*

What if each neuron can compute any function?  
What if neurons can abstain from predictions?

81

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Learn First, Then Predict (& Chain)

- Consider rules  $R_1$  and  $R_2$  obtained to make highly accurate predictions on distribution  $D$ .
- On future examples from  $D$ , apply  $R_1$  **then**  $R_2$ .
- **No!  $R_2$  is applied on distribution  $R_1 \circ D (\neq D)$ . There are no guarantees on that distribution!**
- Are there situations that justify this approach?
  - If information during learning is complete.
  - Could happen, but less realistic assumption.

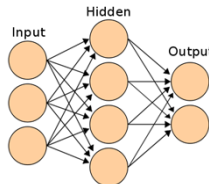
82

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Analogy to Neural Networks

Cannot train each neuron independently and then assign neurons in layers.



*Trivially, because the target to be learned by each specific neuron is not directly observable.*

What if one could train each neuron in isolation?

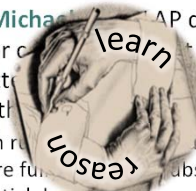
83

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Simultaneously Learn and Predict

- Learn rules for all attributes during training.
  - Apply those rules to enhance the training set.
  - Repeat as necessary to get more KB "layers".
  - **Theorem:** [Michael '14] AP does not reduce soundness or completeness. It may produce KBs with better performance (combined) than **classical** KB.
- Proof:* Train each rule on the distribution, until its predictions are fully accurate. Subsequent training. Doubly-exponential dependence on reasoning depth.



84

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## An Experimental Impasse?

- Scene is an observation. e.g., 1?10010??10
- Goal: Predict what information holds in scene.
- Predictions useful exactly when information is missing from scene. **How to evaluate them?!**
- Evaluate another task with known answers!
  - Use standard approach to solve the task. (1)
  - Make predictions, and give them as input.
  - See whether task performance improves. (2)

85

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Guess the Missing Word [Michael+ '08]

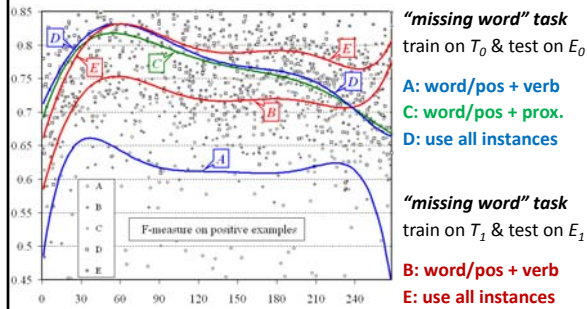
- Training  $T_0$  & testing  $E_0$  from news text corpus.
  - Learn **underlying text structure** [Michael'09].
- Using  $T_0$  learn **relational** thresholds for 268 frequently-used words: *price*, *market*, *stock*, ...
  - Train on  $T_0$  and test on  $E_0$  for each of 268 tasks.
- Using  $T_0$  learn **relational** thresholds for 599 verbs: *buy<sub>sbj,obj</sub>*, *charge<sub>obj</sub>*, *coerce<sub>obj,prd</sub>*, ...
  - **Use verb rules on  $T_0$  to get  $T_1$  / on  $E_0$  to get  $E_1$ .**
  - Train on  $T_1$  and test on  $E_1$  for each of 268 tasks.

86

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Experimental SLAP Results [Michael+ '08]

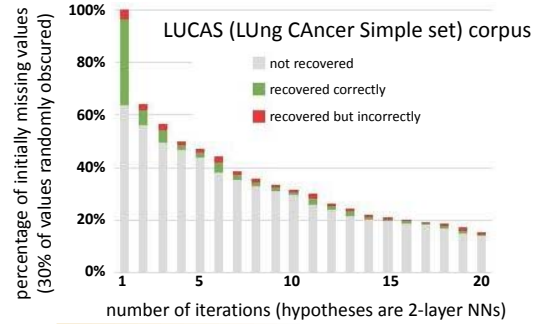


87

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Multiple Iteration SLAP [Skouteli+ '16]

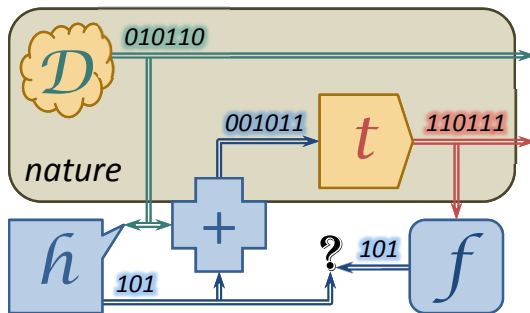


88

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Introspective Forecasting [Michael '15]



89

Introspective Forecasting

Loizos Michael (OUC)

## C. Integrating Deduction and Induction

### Implicit Learning of Testable KBs

90

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### A problem requiring learning and deduction: birds.com Analytics

birds.com™

The (statistical) query: "true with high probability for data source?"

FLY?

The data

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
107	54	Mouse
107	55	Grubs
⋮	⋮	⋮

THE DATA ALONE DOESN'T ANSWER THE QUERY

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### birds.com: Inferring Flying from Food?

- Data:**

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
⋮	⋮	⋮
- Query:** FLY?
- Background knowledge:**
  - PENGUIN  $\vee$  FLY
  - $\neg$ PENGUIN  $\vee$  EAT(FISH)
  - $\neg$ EAT(FISH)  $\vee$   $\neg$ EAT(GRUBS)
  - $\neg$ EAT(FISH)  $\vee$   $\neg$ EAT(MOUSE)
  - $\neg$ EAT(FISH)  $\vee$   $\neg$ EAT(SEED) ...etc...

foods are mutually exclusive...

THE KNOWLEDGE ALONE ALSO DOESN'T ANSWER THE QUERY

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Relevant knowledge we aim to discover from the data

- Relevant property:** the birds eat anything but fish—
  $EAT(GRUBS) \vee EAT(MOUSE) \vee EAT(SEED) \vee \dots$ 
  - Usually satisfied by Data ✓
  - Completing a Proof of Query...

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
⋮	⋮	⋮

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### The Relevant Property completes a Proof of the Query

Negated Query  $\neg$ FLY

FLY

Background Knowledge

Relevant Property

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Relevant knowledge we aim to discover from the data

- Relevant property:** the birds eat anything but fish—
  $EAT(GRUBS) \vee EAT(MOUSE) \vee EAT(SEED) \vee \dots$ 
  - Usually satisfied by Data ✓
  - Completing a Proof of Query ✓
- The **existence** of such knowledge establishes the query (FLY)
  - It isn't important what the actual property says!
  - We are not told what property to search for...

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Approximate Query Answering with Implicit Learning

- Data:**

Day	Bird no.	Food
107	48	Seed
107	49	Grubs
⋮	⋮	⋮
- Background knowledge:**
  - PENGUIN  $\vee$  FLY
  - $\neg$ PENGUIN  $\vee$  EAT(FISH)
  - $\neg$ EAT(FISH)  $\vee$   $\neg$ EAT(GRUBS) ...etc...
- Query:** FLY?
- Verdict:** Does there exist a relevant property e.g.,  $EAT(GRUBS) \vee EAT(MOUSE) \vee EAT(SEED) \vee \dots$ 
  - Usually satisfied by Data
  - Completing Proof of Query

DON'T NEED TO SAY WHAT PROPERTY

AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael



### Proposed Algorithm for Approximate Query Answering

Given a query  $\varphi$ , background knowledge KB, target  $\epsilon$ , and i.i.d. partial examples  $\rho^{(1)}, \dots, \rho^{(m)}$

- For each partial example  $\rho^{(i)}$ ,
  - If  $\varphi |_{\rho^{(i)}}$  provable from KB  $|_{\rho^{(i)}}$ , increment *SUCCESS* (initially, 0)
- If  $SUCCESS/m > 1-\epsilon$ , **ACCEPT**, otherwise **REJECT**

**PARTIAL EVALUATION:** PLUG IN  $\rho^{(i)}$  AND RECURSIVELY REPLACE EACH LOCALLY DETERMINED CONNECTIVE BY ITS VALUE.

$\rho: x=0, y=0$

97 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Proposed Algorithm for Approximate Query Answering

Given a query  $\varphi$ , background knowledge KB, target  $\epsilon$ , and i.i.d. partial examples  $\rho^{(1)}, \dots, \rho^{(m)}$

- For each partial example  $\rho^{(i)}$ ,
  - If  $\varphi |_{\rho^{(i)}}$  provable from KB  $|_{\rho^{(i)}}$ , increment *SUCCESS* (initially, 0)
- If  $SUCCESS/m > 1-\epsilon$ , **ACCEPT**, otherwise **REJECT**

**PARTIAL EVALUATION:** PLUG IN  $\rho^{(i)}$  AND RECURSIVELY REPLACE EACH LOCALLY DETERMINED CONNECTIVE BY ITS VALUE.

98 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Algorithm on birds.com Problem

On most partial examples...

107	48	Seed
107	49	Grubs
107	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
107	54	Mouse
107	55	Grubs
:	:	:

... some background rule simplifies to  $\neg\text{EAT}(\text{FISH})\dots$

**ACCEPT** ... completing a proof of the query

99 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Approximate Queries: The theorem [Juba '13]

**Theorem.** Our algorithm distinguishes:

- The query is only satisfied by  $D$  w.p.  $< 1-\epsilon-\gamma$
- There exists a **"(1- $\epsilon+\gamma$ )-testable"** formula  $\psi$  for which there exists a proof of the query from  $\psi$  and any background knowledge w.p.  $1-\delta$ , given  $1/\gamma^2 \ln 1/\delta$  partial examples from  $M(D)$

100 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Key Definition: Testable formulas

Relevant Properties we can detect [Juba '13]

- Definition.** A formula  $\psi$  is **(1- $\epsilon$ )-testable** under a distribution over partial examples  $M(D)$  if  $\Pr_{M(D)}[\psi |_{\rho}=1] \geq 1-\epsilon$

**PARTIAL EVALUATION:** PLUG IN  $\rho$  AND RECURSIVELY REPLACE EACH LOCALLY DETERMINED CONNECTIVE BY ITS VALUE.

101 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

### Key Definition: Testable formulas

Relevant Properties we can detect [Juba '13]

- Definition.** A formula  $\psi$  is **(1- $\epsilon$ )-testable** under a distribution over partial examples  $M(D)$  if  $\Pr_{M(D)}[\psi |_{\rho}=1] \geq 1-\epsilon$

**birds.com**

Day	Bird no.	Food
97	48	Seed
97	49	Grubs
97	50	Mouse
107	51	Mouse
107	52	Worm
107	53	Fish
:	:	:

$\text{EAT}(\text{GRUBS}) \vee \text{EAT}(\text{MOUSE}) \vee \text{EAT}(\text{SEED}) \vee \dots$

102 AAAI 2018 Tutorial: Integrating Learning into Reasoning Juba and Michael

## Key Definition: Testable formulas

Relevant Properties we can detect [Juba '13]

- **Definition.** A formula  $\psi$  is  $(1-\epsilon)$ -testable under a distribution over partial examples  $M(D)$  if  $\Pr_{M(D)}[\psi|_{\rho}=1] \geq 1-\epsilon$
- Require *more* than truth of Relevant Property
  - Standard cases (clause/linear inequality premises): actually **no** more demanding

103

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Approximate Queries: The theorem

[Juba '13]

- Theorem.** Our algorithm distinguishes:
- The query is only satisfied by  $D$  w.p.  $< 1-\epsilon-\gamma$
  - There exists a  $(1-\epsilon+\gamma)$ -testable formula  $\psi$  for which there exists a proof of the query from  $\psi$  and any background knowledge w.p.  $1-\delta$ , given  $1/\gamma^2 \ln 1/\delta$  partial examples from  $M(D)$

104

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Proposed Algorithm for Approximate Query Answering

Given a query  $\varphi$ , background knowledge KB, target  $\epsilon$ , and i.i.d. partial examples  $\rho^{(1)}, \dots, \rho^{(m)}$

- For each partial example  $\rho^{(i)}$ ,
  - If  $\varphi|_{\rho^{(i)}}$  provable from  $KB|_{\rho^{(i)}}$ , increment *SUCCESS* (initially, 0)
- If  $SUCCESS/m > 1-\epsilon$ , **ACCEPT**, otherwise **REJECT**

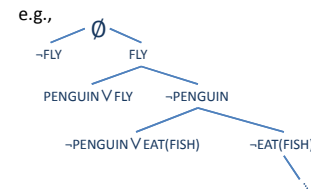
105

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Detecting a Relevant Property... Consider: Tractable Proof Systems

- Bounded-width resolution
- Treelike, bounded clause space resolution



106

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## A Key Property of the Example Tractable Proof Systems

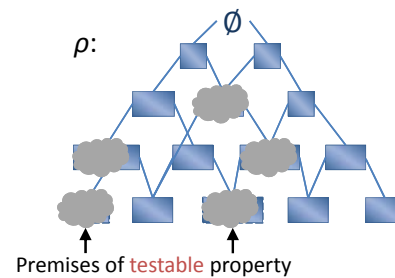
- Bounded-width resolution
  - Treelike, bounded clause space resolution
- ⇒ *Partial evaluation* of proofs of these forms yields proofs of the same form (from a proof of a query  $\varphi$ , we obtain a proof of  $\varphi|_{\rho}$  of the same syntactic form)

107

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Insight: The Testable Premises Drop Out of the Proof!



108

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Recall: Proposed Algorithm Detects These Residual Proofs

Given a query  $\varphi$ , background knowledge KB, target  $\epsilon$ , and i.i.d. partial examples  $\rho^{(1)}, \dots, \rho^{(m)}$

- For each partial example  $\rho^{(i)}$ ,
  - If  $\varphi \mid \rho^{(i)}$  provable from KB  $\mid \rho^{(i)}$ , increment *success* (initially 0)
- If  $\text{SUCCESS} / m > 1 - \epsilon$ , **ACCEPT**, otherwise **REJECT**

**THE THEOREM THEREFORE  
FOLLOWS IMMEDIATELY FROM  
HOEFFDING'S INEQUALITY**

109

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

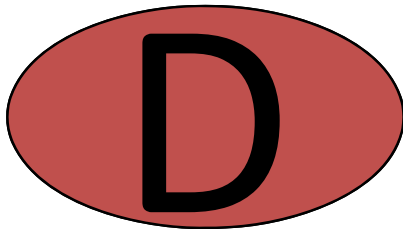
### Pros and cons of Implicit Learning

- **Pro:** utilizes rules with imperfect **validity**
  - Usually intractable to learn explicitly
  - Captures kinds of commonsense reasoning (*next part*)
- **Pro:** reasoning time complexity independent of size of implicit KB
  - Actually, may reduce reasoning complexity in some circumstances [Juba '15]
- **Con:** cannot report rules used to support conclusion

110

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



111

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### High-Level Tutorial Roadmap

- Why KRR Should Embrace Learning
- Introduction to PAC-Semantics
- Integrating Deduction and Induction
- Reasoning Non-Monotonically**
- Overall Summary and Conclusions

112

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### D. Reasoning Non-Monotonically

#### Conditional Probability and NMR

113

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- *We have seen:* deciding queries by counting the frequency with which they are provable
  - Allows us to answer "hard" queries
  - Draws on a potentially large KB of implicitly learned knowledge as needed
- *New twist:* suppose we incorporate a hypothesis by **filtering out examples that do not satisfy the hypothesis**
  - Produces desirable non-monotonic inferences, appropriate for "*commonsense reasoning*"

114

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- **Example:** suppose we are reasoning about “birds” – we filter the set of examples to *only* include  $\text{bird} = 1$ :

bird	fly	has_beak	red	purple	penguin
1	1	1	0	0	0
1	1	1	1	0	0
1	0	1	0	0	1
1	0	1	0	0	1
1	1	1	0	0	0
1	1	1	1	0	0
1	1	1	0	0	0

- We find  $\text{has\_beak}$  is  $(1-\gamma)$ -valid,  $\text{fly}$  is  $(\frac{5}{7}-\gamma)$ -valid, but  $\text{red}$  and  $\text{penguin}$  are at most  $(\frac{2}{7}+\gamma)$ -valid...

115

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- **Example:** Now suppose we consider specifically red birds, filtering to *only* include  $\text{bird} \wedge \text{red} = 1$ :

bird	fly	has_beak	red	purple	penguin
1	1	1	1	0	0
1	1	1	1	0	0
1	1	1	1	0	0
1	1	1	1	0	0
1	1	1	1	0	0
1	1	1	1	0	0
1	1	1	1	0	0

- Now  $\text{has\_beak}$  and  $\text{fly}$  are (still)  $(1-\gamma')$ -valid, and  $\text{penguin}$  is at most  $\gamma'$ -valid (for some  $\gamma' > \gamma$ )...

116

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- **Example:** now suppose we are considering penguins, filtering to only include  $\text{penguin} = 1$ :

bird	fly	has_beak	red	purple	penguin
1	0	1	0	0	1
1	0	1	0	0	1

- We find  $\text{has\_beak}$  is still  $(1-\gamma')$ -valid, but  $\text{fly}$  and  $\text{red}$  are now at most  $\gamma'$ -valid

117

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- **Example:** now considering specifically “penguins with beaks,” we only include  $\text{penguin} \wedge \text{has\_beak} = 1$ :

bird	fly	has_beak	red	purple	penguin
1	0	1	0	0	1
1	0	1	0	0	1

- Again, we find  $\text{fly}$  and  $\text{red}$  are at most  $\gamma'$ -valid – not affected by  $\text{has\_beak}$

118

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity in learning to reason [Roth '95]

- **Example:** if we instead consider “purple penguins,” we only include examples with  $\text{penguin} \wedge \text{purple} = 1$ :

bird	fly	has_beak	red	purple	penguin
1	0	1	0	0	1
1	0	1	0	0	1

- With no examples remaining, we cannot draw *any* conclusions (except perhaps from a given KB)

119

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Nonmonotonicity and conditional probability distributions

- A set of examples filtered to satisfy a formula  $h$  has the conditional probability distribution  $D|[h(X)=1]$  (we “condition on  $h$ ”)
- So: as long as we consider  $(1-\epsilon)$ -validity for some  $\epsilon > 0$  (e.g.,  $\epsilon = 1/3$  could have sufficed), conditioning may have a *non-monotonic* effect
- *Note:* this requires the use of non-negligibly large  $\epsilon$ 
  - Since: we must have enough examples in the conditional distribution to support inferences

120

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Classical issues in commonsense reasoning, in PAC-semantics

- **Qualification problem:** taking  $\epsilon > 0$  permits a rule to fail for any number of unspecified and unmodeled reasons. The commonsense rule we implicitly use is never written out.
- **Elaboration tolerance:** we can simply add a new example to our set of examples, and the next time we answer a query the count will be slightly different. But the implicit KB does not need to be “edited” in any way.

121

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Classical issues in commonsense reasoning, in PAC-semantics

- **Ramification problem:** any further consequences  $r$  to a hypothesis  $h$  are included: for any example  $x$  satisfying  $h$ , we are given that  $x$  also satisfies  $r$ , so  $r$  will be highly **valid** in  $D \mid [h(X)=1]$ . So,  $r$  is included in the implicit KB of  $D \mid [h(X)=1]$  without further consideration.
  - E.g.:  $\neg \text{fly}$  is highly **valid** in  $D \mid [\text{penguin}(X)=1]$ .

122

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Context in reasoning using “preconditions” [Valiant '94,95,00]

- The condition  $h$  in  $D \mid [h(X)=1]$  is sometimes called a “**precondition**”
- Valiant proposed: when answering a query, a **precondition** capturing the current context should be used to filter examples.
  - E.g., given by the units currently firing in the “neuroidal” cognitive model [Valiant '94]
- **Problem:** may be too specific.
- *How should we choose a precondition?*

123

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## D. Reasoning Non-Monotonically

### Preconditions and Abduction

124

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning with preconditions [Juba '16]

- One possible formulation of reasoning with **preconditions**:
 

“Does there exist a **precondition**  $h$  from a class of representations  $H$  such that...

  1.  $h$  supports the query  $\varphi$
  2.  $h$  is common
  3.  $h$  is consistent with the current context  $x^*$ ?”
- If **any precondition** supports the query, we will find one.

125

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning with preconditions: formalization [Juba '16]

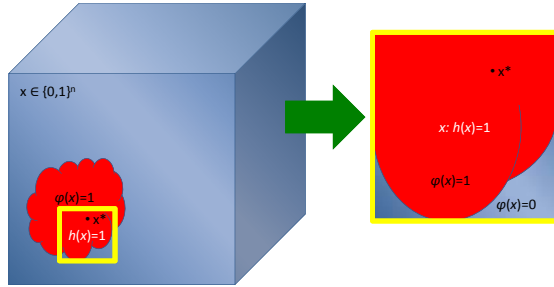
- Fix a class of Boolean representations  $\mathcal{H}$
- **Given** query formula  $\varphi$ ; context assignment  $x^*$ ,  $\epsilon, \delta, \mu \in (0,1)$ ; access to examples from  $D$ ,
- Suppose that there exists a  $h^* \in \mathcal{H}$  such that
  1.  $h^*$  supports  $\varphi$ :  $\Pr_D[\varphi(X) \mid h^*(X)] = 1$
  2.  $h^*$  is common:  $\Pr_D[h^*(X)] \geq \mu$
  3.  $h^*$  is consistent with context  $x^*$ :  $h^*(x^*) = 1$
- **Find** a  $h$  (ideally in  $\mathcal{H}$ ) such that with prob.  $1-\delta$ ,
  1.  $\Pr_D[\varphi(X) \mid h(X)] \geq 1-\epsilon$
  2.  $\Pr_D[h(X)] \geq \mu'$  for some  $\mu'$  (ideally close to  $\mu$ )
  3.  $h(x^*) = 1$

126

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## The precondition task, in pictures...



127

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning with k-DNF preconditions is possible (complete information)

**Theorem.** If there is a k-DNF  $h^*$  such that

1.  $h^*$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h^*(X)] = 1$
2.  $h^*$  is common:  $\Pr_D[h^*(X)] \geq \mu$
3.  $h^*$  is consistent with context  $x^*$ :  $h^*(x^*) = 1$

then using  $m = O(1/\mu \epsilon (n^k + \log^1/\delta))$  examples, in time  $O(mn^k)$  we can find a k-DNF  $h$  such that with probability  $1-\delta$ ,

1.  $h$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h(X)] \geq 1-\epsilon$
2.  $h$  is common:  $\Pr_D[h(X)] \geq \mu$
3.  $h$  is consistent with context  $x^*$ :  $h(x^*) = 1$

128

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Finding a supporting k-DNF precondition using Elimination

- Start with  $h$  as an OR over all terms of size  $k$
- For each example  $x^{(1)}, \dots, x^{(m)}$ 
  - If  $\varphi(x^{(i)}) = 0$ , delete all terms  $T$  from  $h$  such that  $T(x^{(i)}) = 1$
- If  $h(x^*) = 1$ , return  $h$
- Else return FAIL (no supporting precondition)

Running time is still clearly  $O(mn^k)$

129

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Analysis pt 1: $\Pr_D[h(X)] \geq \mu$ , $h(x^*) = 1$

- We are given that some k-DNF  $h^*$  has
  1.  $h^*$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h^*(X)] = 1$
  2.  $h^*$  is common:  $\Pr_D[h^*(X)] \geq \mu$
  3.  $h^*$  is consistent with context  $x^*$ :  $h^*(x^*) = 1$
- Initially, every term of  $h^*$  is in  $h$
- Terms of  $h^*$  are never true when  $\varphi(x)=0$  by 1.
- $\Rightarrow$  every term of  $h^*$  remains in  $h$
- $\Rightarrow h^*$  implies  $h$ , so  $\Pr_D[h(X)] \geq \Pr_D[h^*(X)] \geq \mu$  and  $h(x^*) = 1$  since  $h^*(x^*) = 1$  by 3.

130

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Analysis pt 2: $\Pr_D[\varphi(X) | h(X)] \geq 1-\epsilon$

- Rewrite conditional probability:  $\Pr_D[[-\varphi(X)] \wedge h(X)] \leq \epsilon \Pr_D[h(X)]$
- We'll show:  $\Pr_D[[-\varphi(X)] \wedge h(X)] \leq \epsilon \mu$  ( $\leq \epsilon \Pr_D[h(X)]$  by part 1)
- Consider any  $h'$  s.t.  $\Pr_D[[-\varphi(X)] \wedge h'(X)] > \epsilon \mu$ 
  - Since each  $X^{(i)}$  is drawn independently from  $D$   $\Pr_D[\text{no } i \text{ has } [-\varphi(X^{(i)})] \wedge h'(X^{(i)})] < (1-\epsilon\mu)^m$
  - A term of  $h'$  is deleted when  $\varphi=0$  and  $h'=1$
  - So,  $h'$  is only possibly output w.p.  $< (1-\epsilon\mu)^m$

131

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Analysis pt 2, cont'd:

### $\Pr_D[\varphi(X) | h(X)] \geq 1-\epsilon$

- We'll show:  $\Pr_D[[-\varphi(X)] \wedge h(X)] \leq \epsilon \mu$
- Consider any  $h'$  s.t.  $\Pr_D[[-\varphi(X)] \wedge h'(X)] > \epsilon \mu$ 
  - $h'$  is only possibly output w.p.  $< (1-\epsilon\mu)^m$
- There are only  $2^{O(n^k)}$  possible k-DNF  $h'$
- Since  $1-z \leq e^{-z}$ ,  $m = O(1/\mu \epsilon (n^k + \log^1/\delta))$  ex's suffice to guarantee that each such  $h'$  is only possibly to output w.p.  $< \delta/2^{O(n^k)}$
- $\Rightarrow$  w.p.  $> 1-\delta$ ,  $h$  has  $\Pr_D[[-\varphi(X)] \wedge h(X)] \leq \epsilon \mu$ . ■

132

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Reasoning with k-DNF preconditions is possible (complete information)

**Theorem.** If there is a k-DNF  $h^*$  such that

1.  $h^*$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h^*(X)] = 1$
2.  $h^*$  is common:  $\Pr_D[h^*(X)] \geq \mu$
3.  $h^*$  is consistent with context  $x^*$ :  $h^*(x^*) = 1$

then using  $m = O(\frac{1}{\mu\epsilon} (n^k + \log \frac{1}{\delta}))$  examples, in time  $O(mn^k)$  we can find a k-DNF  $h$  such that with probability  $1-\delta$ ,

1.  $h$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h(X)] \geq 1-\epsilon$
2.  $h$  is common:  $\Pr_D[h(X)] \geq \mu$
3.  $h$  is consistent with context  $x^*$ :  $h(x^*) = 1$

133

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Extension: Finding preconditions tolerating $\epsilon > 0$

• **Given only** that some  $h^*$  achieves

1.  $h^*$  supports  $\varphi$ :  $\Pr_D[\varphi(X) | h^*(X)] \geq 1-\epsilon$
2.  $h^*$  is common:  $\Pr_D[h^*(X)] \geq \mu$
3.  $h^*$  is consistent with context  $x^*$ :  $h^*(x^*) = 1$

**Find** an  $h$  such that for some other  $\mu'$  &  $\epsilon'$ ,

1.  $\Pr_D[\varphi(X) | h(X)] \geq 1-\epsilon'$
2.  $\Pr_D[h(X)] \geq \mu'$  for some  $\mu'$  (ideally close to  $\mu$ )
3.  $h(x^*) = 1$

• **Extension of algorithm for k-DNF achieves**  $\mu' = \mu$ ,  $\epsilon' = O(n^k \epsilon)$  (only delete T making  $\epsilon \mu$  errors)

134

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Abductive reasoning: making plausible guesses

- This formulation of **precondition** search is a form of **abductive reasoning**
  - **Given** a conclusion  $c$ , find a "plausible"  $h$  that implies/leads to/...  $c$
  - Proposing a precondition supporting the query
- Two varieties of "plausibility" in common use
  - *Syntactic*: a small  $h$  from which  $c$  follows
  - *Bayesian*: a  $h$  which has large posterior probability when given  $c$  ie., " $\Pr[h \text{ actual rule used} | c \text{ true}] > \dots$ "

135

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Why might we want a new model?

- Existing models only **tractable** in *simple cases*
  - E.g. Horn rules ( $a \wedge b \wedge c \Rightarrow d$  ...no negations), "nice" (conjugate) priors
- The *choice* of formulation, prior distribution, etc. **really matters**
  - And, they are difficult to specify by hand

136

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## New model: abductive reasoning from random examples [Juba '16]

- **Task**: for a conclusion  $c$ , find a  $h$  such that
  1. *Plausibility*:  $\Pr_D[h(X)] \geq \mu$  (for some given  $\mu$ )
  2.  $h$  almost entails  $c$ :  $\Pr_D[c(X) | h(X)] \geq 1-\epsilon$
- **Note**: D now captures **both** the entailment relation **and** the measure of "plausibility"
- **Distinction from earlier precondition search**: no "context" assignment  $x^*$

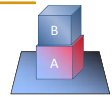
137

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Example: identifying a subgoal

- **Consider**: blocks world. For  $t=1,2,\dots,T$ 
  - Propositional state vars. ("fluents")  $ON_t(A,B)$ ,  $ON_t(A, \text{TABLE})$ ,  $ON_t(C,A)$ , etc.
  - Actions also encoded by propositional vars.  $PUT_t(B,A)$ ,  $PUT_t(C, \text{TABLE})$ , etc.
- **Given many examples of interaction...**
- **Our goal**  $c$ :  $ON_T(A, \text{TABLE}) \wedge ON_T(B,A) \wedge ON_T(C,B)$
- A perhaps plausibly good "subgoal"  $h$ :  $[ON_{T-1}(B,A) \wedge PUT_T(C,B)] \vee [PUT_{T-1}(B,A) \wedge PUT_T(C,B)]$



138

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Formally: *abductive reasoning from random examples* for a class $\mathcal{H}$

- Fix a class of Boolean representations  $\mathcal{H}$
- Given** Boolean formula  $c$ ;  $\epsilon, \delta, \mu \in (0,1)$ ; independent examples  $x^{(1)}, \dots, x^{(m)} \in D$ ,
- Suppose that there exists a  $h^* \in \mathcal{H}$  such that
  - Plausibility:  $\Pr_D[h^*(X)] \geq \mu$
  - $h^*$  entails  $c$ :  $\Pr_D[c(X) | h^*(X)] = 1$
- Find**  $h$  (ideally in  $\mathcal{H}$ ) such that with prob.  $1-\delta$ ,
  - Plausibility:  $\Pr_D[h(X)] \geq \mu'$  for some  $\mu'(\mu, n, \epsilon, \delta)$
  - $h$  almost entails  $c$ :  $\Pr_D[c(X) | h(X)] \geq 1-\epsilon$

139

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Abducing k-DNFs is also possible (*complete information*)

**Theorem.** If there is a k-DNF  $h^*$  such that

- Plausibility:  $\Pr_D[h^*(X)] \geq \mu$
- $h^*$  entails  $c$ :  $\Pr_D[c(X) | h^*(X)] = 1$

then using  $m = O(1/\mu\epsilon (n^k + \log^2 1/\delta))$  examples, in time  $O(mn^k)$  we can find a k-DNF  $h$  such that with probability  $1-\delta$ ,

- Plausibility:  $\Pr_D[h(X)] \geq \mu$
- $h$  almost entails  $c$ :  $\Pr_D[c(X) | h(X)] \geq 1-\epsilon$

140

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Elimination Algorithm also solves k-DNF abduction

- Start with  $h$  as an OR over all terms of size  $k$
- For each example  $x^{(1)}, \dots, x^{(m)}$ 
  - If  $c(x^{(i)}) = 0$ , delete all terms  $T$  from  $h$  such that  $T(x^{(i)}) = 1$
- Return  $h$

JUST OMIT THE FINAL TEST FOR CONSISTENCY WITH  $x^*$ . THE ANALYSIS IS ALMOST IDENTICAL.

141

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Extension: *Abductive Reasoning tolerating $\epsilon > 0$*

- Given only** that some  $h^*$  achieves
  - Plausibility:  $\Pr_D[h^*(X)] \geq \mu$
  - $h^*$  almost entails  $c$ :  $\Pr_D[c(X) | h^*(X)] \geq 1-\epsilon'$
- Find** an  $h$  such that for some other  $\mu'$  &  $\epsilon'$ ,
  - Plausibility:  $\Pr_D[h(X)] \geq \mu'$
  - $h$  almost entails  $c$ :  $\Pr_D[c(X) | h(X)] \geq 1-\epsilon'$
- Improved algorithm for k-DNF achieves  $\mu' = (1-\gamma)\mu$ ,  $\epsilon' = \tilde{O}(n^{k/2}\epsilon)$  [Zhang-Mathew-Juba'17]
  - Cf. only obtained  $\epsilon' = O(n^k\epsilon)$  for *preconditions*...

142

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### But, what about abducing conjunctions?? [Juba '16]

**Theorem.** Suppose that a polynomial-time algorithm exists for learning abduction from random examples for conjunctions with  $\mu' = C((1-\gamma)\mu/n)^d$  for some  $C, d$ . Then there is a polynomial-time algorithm for PAC-learning DNF.

- So what?**
- Central open problem in computational learning theory raised in original paper by Valiant (1984)
  - Recent work by Daniely and Shalev-Shwartz (2016) shows that algorithms for PAC-learning DNF would have other surprising consequences.
- In summary* – an algorithm for our problem would constitute an unlikely breakthrough.

143

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

### Extension to abduction with partial examples

- Subtle: *do we condition on...*
  - $h(X)$  true?
  - $h(X)$  provable (under  $\rho$ )?
  - $h|_\rho = 1$  ( $h$  "witnessed" on  $\rho$ )?
- Come see our poster "Learning Abduction Under Partial Observability" (Juba, Li, Miller)
  - Short version: condition on *some term  $T$  of  $h$  provable under  $\rho$*  (i.e.  $T|_\rho$  provable from  $KB|_\rho$ )
  - Can use Elimination; incorporates implicit KB

Generally more desirable...

144

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



## Extent of commonsense reasoning with PAC-semantics

COULD THIS CAPTURE ALL COMMONSENSE REASONING?

- Unlikely!
  - So far, generally doesn't capture "law of inertia," naive theories, ...
- We still seem to require the use of a non-monotonic logic as a foundation
  - At minimum, for "law of inertia"
- **Question:** can simulations provide examples for naive theories?



## Naïve use of preconditions is problematic

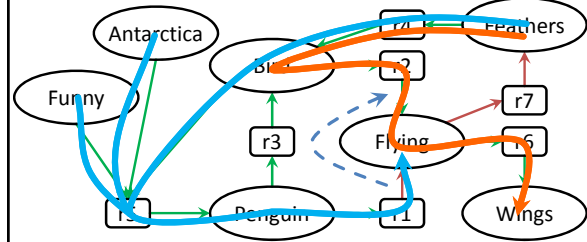
- Our reasoning with preconditions is too credulous
- **Example:** consider the query fly for  $x^*$  with  $penguin(x^*)=1$ ... then  $h^* = bird$ 
  - Is reasonably common ( $Pr_D[bird]$  moderate)
  - Is consistent with  $x^*$  ( $bird(x^*)=1$ )
  - Supports the query ( $Pr[fly|bird]$  high)
- So, we will return a precondition such as  $h = bird$  supporting fly

## D. Reasoning Non-Monotonically

### Learning with Non-Monotonic Logics

## Activity: Find the Arguments

This animal has Feathers, lives in Antarctica, and looks Funny. **Question:** Does it have Wings?



## Why Abandon Equivalences?

- $r_1: penguin \Rightarrow \neg flying, r_2: bird \Rightarrow flying, r_1 > r_2$   
Formula:  $(u \vee bird) \wedge \neg penguin \Leftrightarrow flying$ 
  - Good on "full" scenes. Still abstains on  $\{\neg p, \neg b\}$ .
  - Infers too little... Does not infer  $f$  on  $\{b\}$ . **Bad!**
- $r_1: b \Rightarrow a, r_2: \neg b \Rightarrow a$  Formula:  $true \Leftrightarrow a$ 
  - Infers too much...  $a$  by case analysis on  $\{b\}$ . **Bad!**
- NP-hard reasoning. Still not 1 rule/atom. **Bad!**
- **Thus:** logic-based arguments with preferences.

## How to Learn Arguments?

- Rules with different heads in each argument, therefore, one has to deal with **partial observability**, which then requires **SLAP**.
  - Sufficient to get consistency in predictions.
- Some **relational expressivity** comes for "free".
  - Each rule has an IQE body that is efficiently testable for satisfaction on observations.
- Following **linear thresholds** and **decision lists**.
  - Online, efficient, dealing with priorities, etc.

## Never-Ending Rule Discovery [Michael '16]

1. Get observation *obs*, and **reason with active rules** to get an enhanced observation *obs\**.
2. Find a literal *x* that is observed in *obs* but **not inferred** by active rules triggered in *obs\**. Add *body*  $\Rightarrow$  *x*, for random *body* satisfied by *obs\**.
3. Increase / decrease weight of rules triggered in *obs\** that **concur with / oppose** *obs*.
4. Newly active rules are **weaker than** existing.
5. Newly inactive rules, have **no preferences**.

151

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## What Does it End Up Learning?

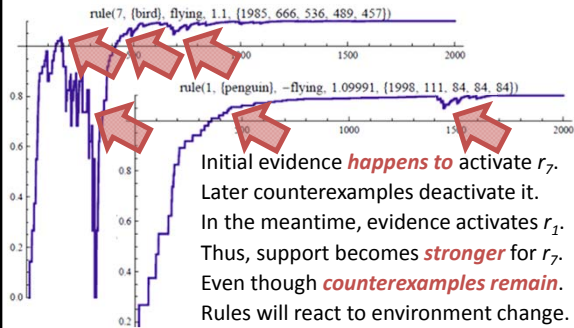
- Hide some attributes from states drawn from:
  - {*-bird*, *-penguin*, *-plane*, *-flying*} w.p. 5/16
  - {*-bird*, *-penguin*, *-plane*, *flying*} w.p. 5/16
  - {*-bird*, *-penguin*, *plane*, *flying*} w.p. 2/16
  - {*bird*, *penguin*, *-plane*, *-flying*} w.p. 1/16
  - {*bird*, *-penguin*, *-plane*, *flying*} w.p. 3/16
- “Intended” learned rules with head (*-*)*flying*:
  - penguin*  $\Rightarrow$  *-flying*, *bird*  $\Rightarrow$  *flying*, *plane*  $\Rightarrow$  *flying*
- But also “picks up”: *plane*  $\Rightarrow$  *-penguin* (mutually exclusive), *-bird*  $\Rightarrow$  *-penguin* (contrapositive), *penguin*  $\Rightarrow$  *bird*, *-flying*  $\Rightarrow$  *-bird* (explaining away).

152

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Demo of NERD's Online Behavior



153

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Probably Approximately Correct?

- Equivalences are PAC learnable, but suffer from Goldilocks effect: infer too much / little.
  - Are logic-based arguments PAC learnable?
- + Learn with unknown atoms. Learn priorities.  
 - Nested if-then-else's unlearnable from scenes.  
 + Non-adversarial environments. Equiv  $\neq$  Args.  
 - Learning requires reasoning, restricts depth.  
 + Psychological evidence on bounded reasoning.  
 + **On edge of learnable. Ev+...re.**

154

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

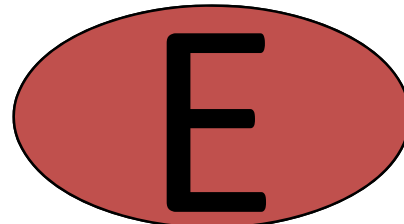
## Machine Coaching [Michael '17]

- Integrate user-provided rules **while** learning.
  - Like online learning, but instead of providing the correct label, “question” part of the argument that leads to the wrong prediction.
  - Possible to give PAC-semantics to the process!
- **Theorem:** Arguments (*ASPIC+ type: grounded semantics, axiomatic premises, defeasible rules, rebutting attacks, last link preferences*) are PAC learnable via **machine coaching** alone.

155

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



156

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## High-Level Tutorial Roadmap

- A. Why KRR Should Embrace Learning
- B. Introduction to PAC-Semantics
- C. Integrating Deduction and Induction
- D. Reasoning Non-Monotonically
- E. Overall Summary and Conclusions

157

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## E. Overall Summary and Conclusions

### Recap and Open Problems

158

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Semantics for Learned Knowledge

- PAC-semantics is a suitable choice:
  - models the world as a probability distribution
  - treats knowledge as high-probability properties
- Offers simple solutions to classic KR problems:
  - non-monotonicity from conditioning of rules
  - exogenous qualification from validity defect
  - elaboration tolerance through implicit learning
  - ramifications incorporated in “implicit KB”
  - natural formulation of abductive inference
  - explicit solutions through argument learning

159

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## KRR and Learning Integration

- Traditional view of reasoning and learning as independent processes must be abandoned.
- When combined under PAC-semantics they:
  - soundly achieve greater completeness
  - circumvent computational barriers in learning
  - enable fast and compact access to “implicit KB”
  - may sometimes reduce reasoning complexity
- Also, seamlessly supports user intervention during learning through machine coaching.

160

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Open Problems with KRR Flavor

- What kind of first-order expressions can we learn from partial examples in reasoning?
  - IQEs via propositionalization. Classical barriers (e.g., Haussler) apply to integrated problem?
- What kind of first-order integrated learning and reasoning is possible?
  - May want first-order expressions to refer to a limited domain. Seems closely related to selection of “preconditions” (see point in next slide) but for limiting the domain of quantifiers.

161

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## Open Problems with KRR Flavor

- When, and how broadly, is the complexity of reasoning reduced in the integrated problem?
  - More natural condition than in [Juba '15]?
  - More broadly: which fragments are tractable?
- Incorporating naive theories into the model?
  - E.g., treat naive simulations as populating data set for implicit KB — how well does this work?
- Preconditions for commonsense reasoning?
  - Perhaps select an “unchallenged” precondition; suggests argument semantics for precondition.

162

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## E. Overall Summary and Conclusions

### Bibliography of Related Work

163

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## References

- Adams. (1975), 'The Logic of Conditionals: An Application of Probability to Deductive Logic', Springer.
- Blumer, Ehrenfeucht, Haussler, Warmuth. (1987), 'Occam's Razor', Information Processing Letters, 24:377–380.
- Daniely, Shalev-Shwartz. (2016), 'Complexity Theoretic Limitations on Learning DNF's', Proc. COLT, volume 49 of JMLR W&CP.
- Diakidoy, Kakas, Michael, Miller. (2014), 'Story Comprehension through Argumentation', Proc. COMMA, p. 31–42.
- Diakidoy, Kakas, Michael, Miller. (2015), 'STAR: A System of Argumentation for Story Comprehension and Beyond', Proc. Commonsense, p. 64–70.
- Dimopoulos, Michael, Athienitou. (2009), 'Ceteris Paribus Preference Elicitation with Predictive Guarantees', Proc. IJCAI, p. 1890–1895.
- Dung. (1995), 'On the Acceptability of Arguments and its Fundamental Role in Nonmonotonic Reasoning, Logic Programming and n-Person Games', Artificial Intelligence 77:321–357.

164

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## References

- Haussler. (1989), 'Learning Conjunctive Concepts in Structural Domains', Machine Learning 4(1):7–40.
- Juba. (2013), 'Implicit Learning of Common Sense for Reasoning', Proc. IJCAI, p. 939–946.
- Juba. (2015), 'Restricted Distribution Automatizability in PAC-Semantics', Proc. ITCS, p. 93–102.
- Juba. (2016), 'Learning Abductive Reasoning Using Random Examples', Proc. AAAI, p. 999–1007.
- Kakas, Michael. (2016), 'Cognitive Systems: Argument and Cognition', IEEE Intelligent Informatics Bulletin 17(1):14–20.
- Kakas, Michael, Miller. (2008), 'Fred meets Tweety', Proc. ECAI, p. 747–748.
- Khardon, Roth. (1997), 'Learning to Reason', J. ACM, 44(5):697–725.
- Michael. (2008), 'Autodidactic Learning and Reasoning', PhD thesis, Harvard.
- Michael. (2009), 'Reading Between the Lines', Proc. IJCAI, p. 1525–1530.
- Michael. (2010), 'Partial Observability and Learnability', Artificial Intelligence 174(11):639–669.

165

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael

## References

- Michael. (2011), 'Causal Learnability', Proc. IJCAI, p. 1014–1020.
- Michael. (2013), 'Machines with WebSense', Proc. Commonsense.
- Michael. (2014), 'Simultaneous Learning and Prediction', Proc. KR, p. 348–357.
- Michael. (2015a), 'The Disembodied Predictor Stance', Pattern Recognition Letters 64(C):21–29.
- Michael. (2015b), 'Introspective Forecasting', Proc. IJCAI, p. 3714–3720.
- Michael. (2016), 'Cognitive Reasoning and Learning Mechanisms', Proc. AIC, p. 2–23.
- Michael. (2017), 'The Advice Taker 2.0', Proc. Commonsense.
- Michael, Papageorgiou. (2013), 'An Empirical Investigation of Ceteris Paribus Learnability', Proc. IJCAI, p. 1537–1543.
- Michael, Valiant. (2008), 'A First Experimental Demonstration of Massive Knowledge Infusion', Proc. KR, p. 378–388.
- Modgil, Prakken. (2014), 'The ASPIC+ Framework for Structured Argumentation: A Tutorial', Argument & Computation 5(1):31–62.

166

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael


## References

- Nilsson. (1986), 'Probabilistic Logic', Artificial intelligence, 28(1):71–87.
- Roth. (1995), 'Learning to Reason: The Non-monotonic Case', Proc. IJCAI, vol. 2, p. 1178–1184.
- Skouteli, Michael. (2016), 'Empirical Investigation of Learning-Based Imputation Policies', Proc. GCAI, p. 161–173.
- Toni. (2014), 'A Tutorial on Assumption-based Argumentation', Argument & Computation 5(1):89–117.
- Valiant. (1984), 'A Theory of the Learnable', Communications of the ACM, 18(11):1134–1142.
- Valiant. (1994), 'Circuits of the Mind', Oxford University Press.
- Valiant. (1995), 'Rationality', Proc. COLT, p. 3–14.
- Valiant. (2000), 'Robust Logics', Artificial Intelligence, 117:231–253.
- Valiant. (2000b), 'A Neuroidal Architecture for Cognitive Computation', J. ACM, 47(5):854–882.
- Zhang, Mathew, Juba. (2017), 'An Improved Algorithm for Learning to Perform Exception-Tolerant Abduction', Proc. AAAI, p. 1257–1265.

167

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael



end

168

AAAI 2018 Tutorial: Integrating Learning into Reasoning

Juba and Michael