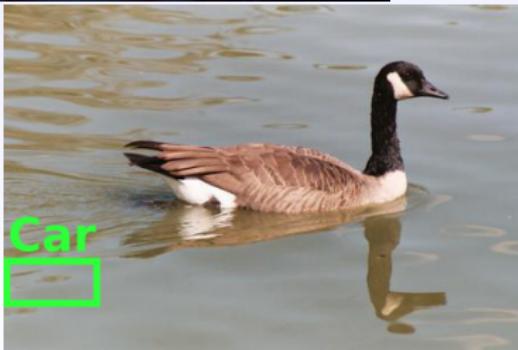
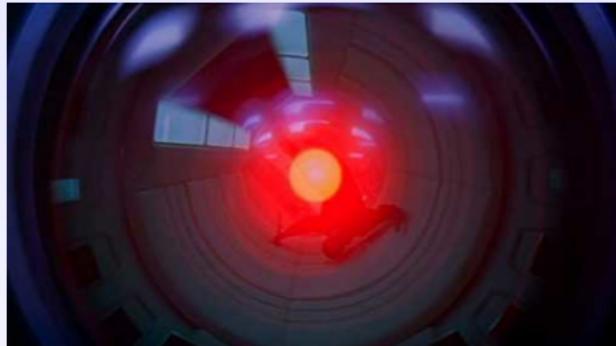


# Learning and reasoning with probabilistic satisfiability modulo theories

Paolo Morettin

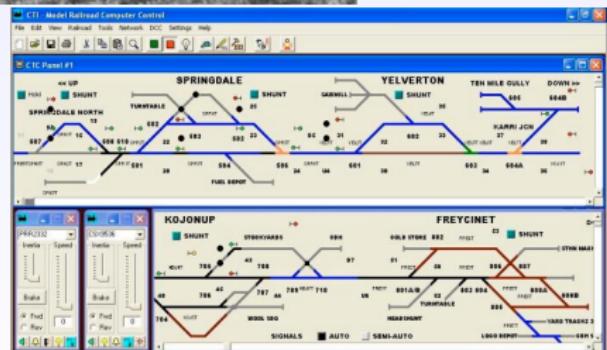


# Concerns on AI



Can we **enforce** and/or **verify**: safety, robustness, ... ?

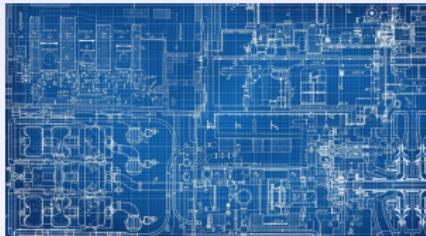
# Concerns on AI other systems?



Not as concerned when we: catch a flight / ride a train / get a CT scan

# Formal verification - the “traditional” approach

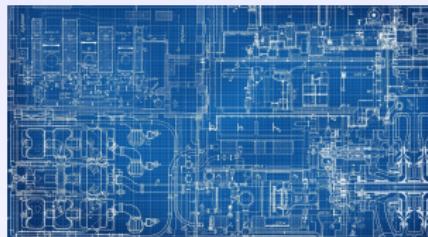
- 1) **Logical modelling** of the system  $S$  and property  $P$ :



*“The elevator door doesn’t open  
if the brakes are off.”*

# Formal verification - the “traditional” approach

1) **Logical modelling** of the system  $S$  and property  $P$ :



$$S = ((\dots \vee \dots) \dots \wedge \dots)$$

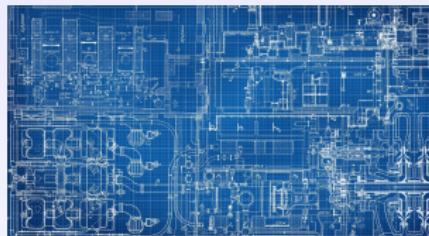
*“The elevator door doesn’t open if the brakes are off.”*



$$P = \neg(DoorOpen \wedge BrakesOff)$$

# Formal verification - the “traditional” approach

1) **Logical modelling** of the system  $S$  and property  $P$ :



*“The elevator door doesn’t open if the brakes are off.”*



$$S = ((\dots \vee \dots) \dots \wedge \dots)$$



$$P = \neg(DoorOpen \wedge BrakesOff)$$

2) Verification of  $P$  is reduced to a **decision problem**:

$$SAT(S \wedge \neg P)?$$



{ YES, NO }

# Formal verification - continuous/discrete models

**Satisfiability Modulo Theories (SMT)** = Logic + Specialized theories

$(A \vee \text{read}(\text{write}(a, i, v), i) = v) \wedge (\text{read}(a, j)) \rightarrow A)$  arrays

$(A \vee (a = b \rightarrow f(a) = g(b))) \wedge ((f(.)) = g(.)) \rightarrow A)$  uninterpreted functions

...

$(A \vee (10x + 13y \leq z + 17/8)) \wedge ((z \geq 1/3) \rightarrow A)$  linear algebra

*"If the elevator speed is greater than k, the alarm is on."*

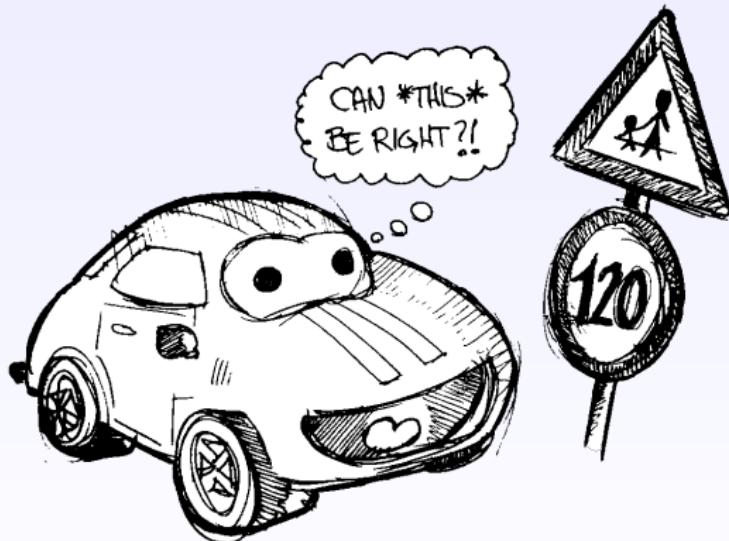
$$P = \neg(\text{speed} > k) \vee \text{AlarmOn}$$

# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?

# When YES/NO are not enough

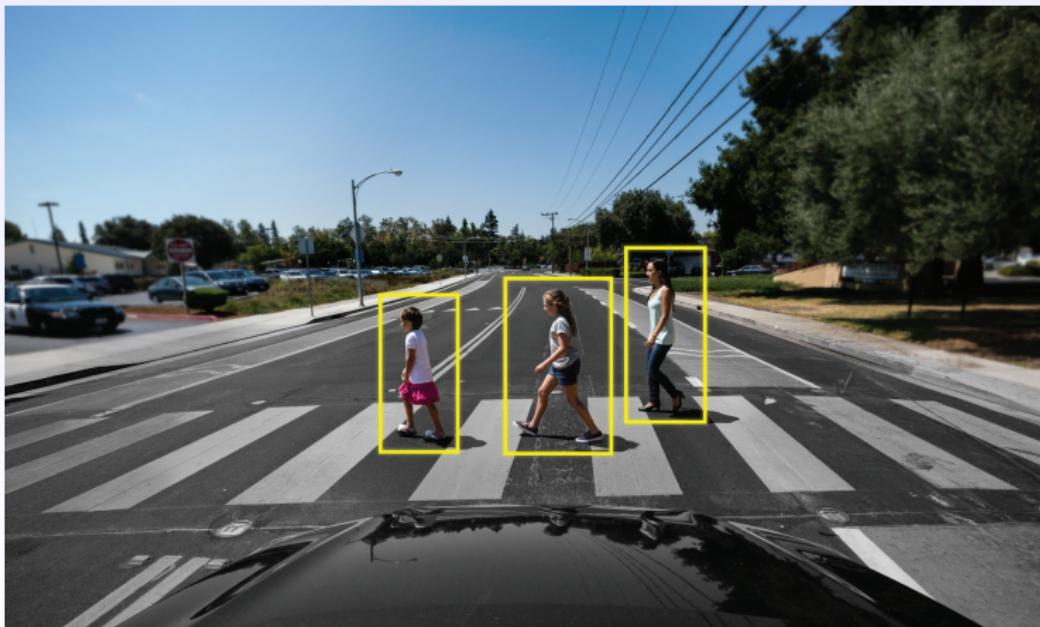
- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?



$$\neg((SpeedLimit = 120) \wedge SchoolCrossing)$$

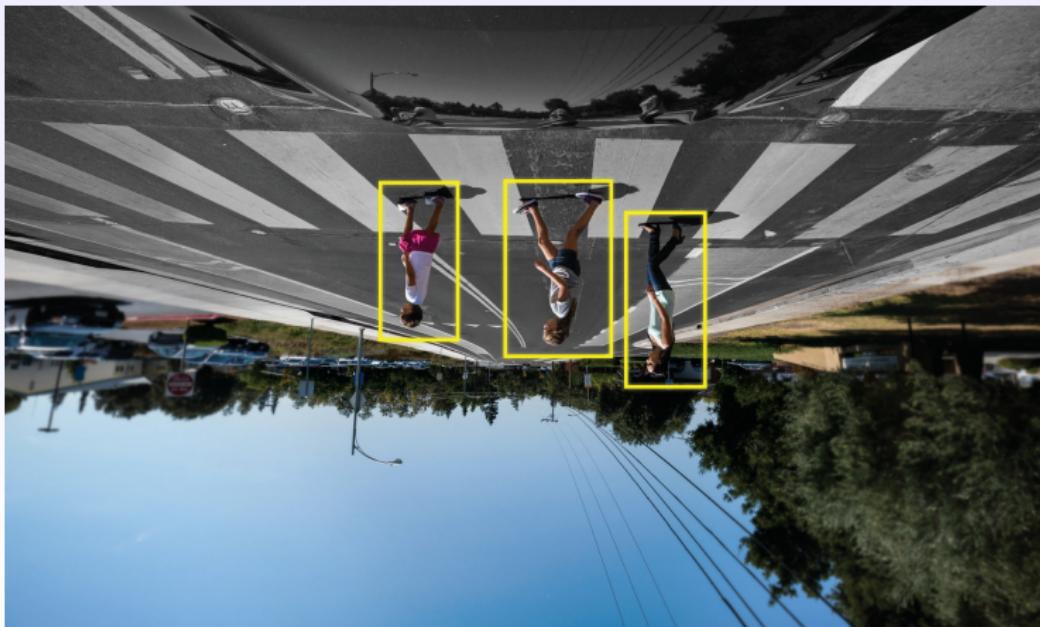
# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?



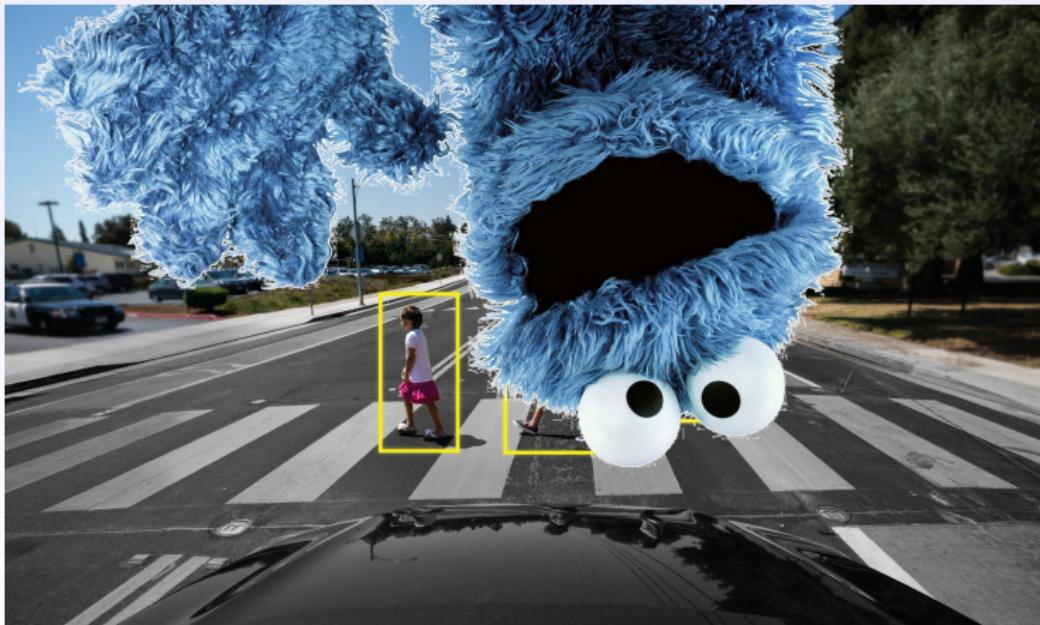
# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?



# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?



# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?
- $Pr(\text{[image]}) = Pr(\text{[image]}) = \epsilon$  ... we don't care about these inputs!

# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?
- $Pr(\text{[image]}) = Pr(\text{[image]}) = \epsilon$  ... we don't care about these inputs!
- We need a probabilistic model of the environment  $P(X)$

# When YES/NO are not enough

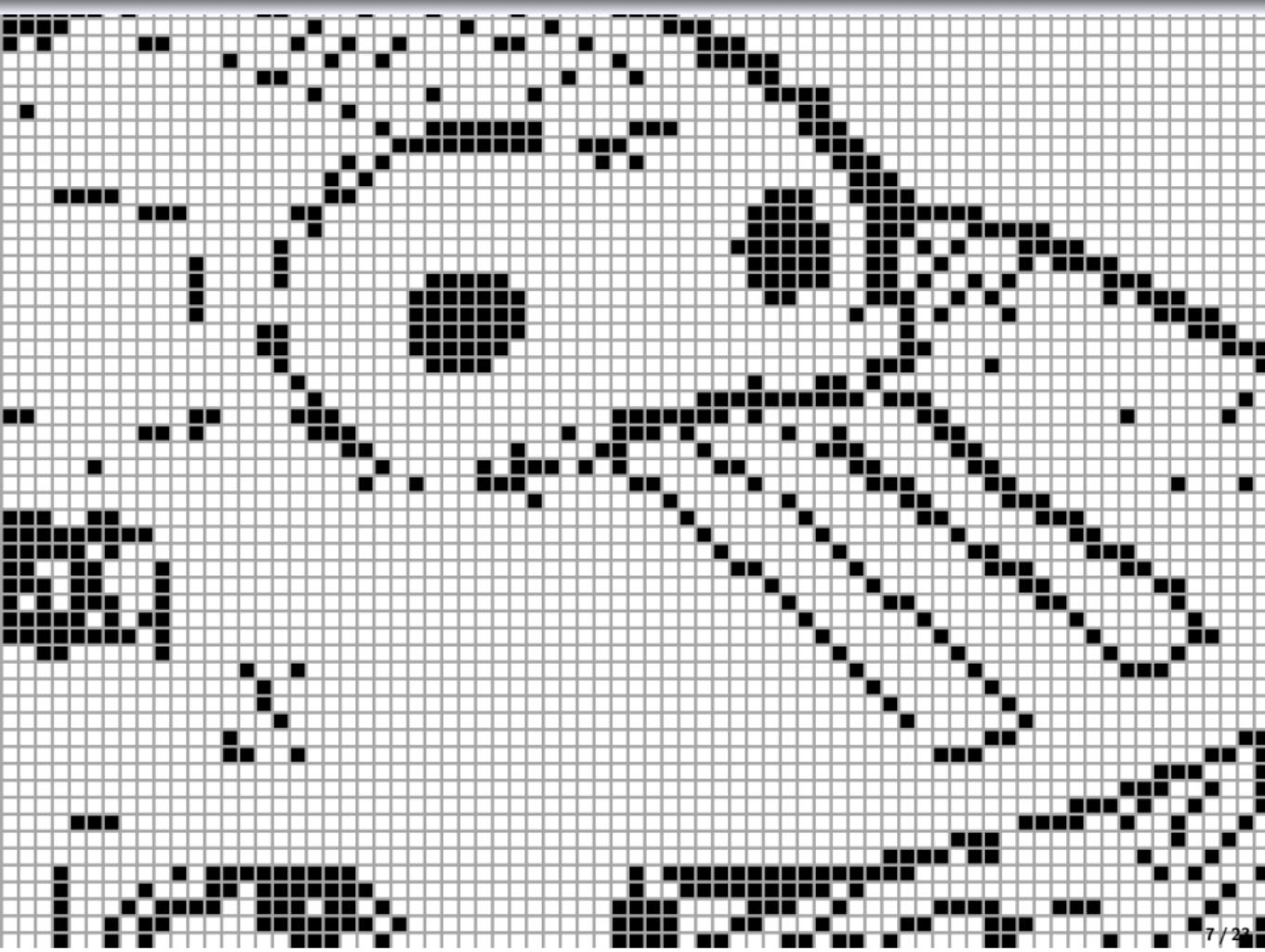
- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?
- $Pr(\text{[image]}) = Pr(\text{[image]}) = \epsilon$  ... we don't care about these inputs!
- We need a probabilistic model of the environment  $P(X)$
- ...and verify whether properties hold with probability  $\geq k$

# When YES/NO are not enough

- we need to *quantify* properties like robustness/fairness?
- ...What if:
  - the system is nondeterministic?
  - the environment/input is high-dimensional and uncertain?

- $Pr(\text{[image]}) = Pr(\text{[image]}) = \epsilon$  ... we don't care about these inputs!
- We need a probabilistic model of the environment  $P(X)$
- ...and verify whether properties hold with probability  $\geq k$

**probabilistic** formal verification!



# Weighted Model Counting to the rescue!

Weighted sum of the models of a logical formula

$$WMC(\Delta, w) = \sum_{\mu \models \Delta} w(\mu)$$

# Weighted Model Counting to the rescue!

Weighted sum of the models of a logical formula

$$WMC(\Delta, w) = \sum_{\mu \models \Delta} w(\mu)$$

$w$  factorizes over the literals:

$$w(\overbrace{\mu_1}^{A \wedge \neg B}) = k_1 = w(A) \cdot w(\neg B)$$

$$w(\overbrace{\mu_2}^{A \wedge B}) = k_2 = w(A) \cdot w(B)$$

|         |         |        |     |
|---------|---------|--------|-----|
|         |         | $A$    |     |
| $\top$  | $k_1$   | $k_2$  | $B$ |
| $\perp$ | $k_3$   | $k_4$  |     |
|         | $\perp$ | $\top$ |     |
|         |         |        |     |

# Weighted Model Counting to the rescue!

Weighted sum of the models of a logical formula

$$WMC(\Delta, w) = \sum_{\mu \models \Delta} w(\mu)$$

$w$  factorizes over the literals:

$$w(\overbrace{\mu_1}^{A \wedge \neg B}) = k_1 = w(A) \cdot w(\neg B)$$

$$w(\overbrace{\mu_2}^{A \wedge B}) = k_2 = w(A) \cdot w(B)$$

|         |       |        |     |
|---------|-------|--------|-----|
|         |       | $A$    |     |
| $\top$  | $k_1$ | $k_2$  | $B$ |
| $\perp$ | $k_3$ | $k_4$  |     |
| $\perp$ |       | $\top$ |     |
|         |       |        |     |

Marginal inference via WMC

$$\Pr(B|A) = \frac{WMC(A \wedge B, w)}{WMC(A, w)} = \frac{k_2}{k_1 + k_2}$$

**Our model:** a (parametric) PLP  $P$

$\theta_1 :: \text{stress}(X) \leftarrow \text{person}(X).$

$\theta_2 :: \text{influences}(X, Y) \leftarrow \text{person}(X),$   
 $\text{person}(Y).$

$\text{smokes}(X) \leftarrow \text{stress}(X).$

$\text{smokes}(X) \leftarrow \text{friend}(X, Y),$   
 $\text{influences}(Y, X),$   
 $\text{smokes}(Y).$

$\text{person}(a). \text{person}(b). \text{person}(c).$

$\text{friend}(a, b).$

...

$\text{evidence}(\text{friend}(a, b)).$

$\text{evidence}(\text{smokes}(c)).$

...

**Our model:** a (parametric) PLP  $P$

1) Conversion in weighted prop. logic

$$\Delta = \text{person}(a) \wedge \text{person}(b)$$

$$P \rightarrow P_{\text{ground}} \rightarrow \langle \Delta, w \rangle$$

...

$$\wedge (\text{aux}_1 \wedge \text{person}(b) \rightarrow \text{stress}(b)) \vee$$

...

$$w(\text{person}(a)) = 1$$

$$w(\text{person}(b)) = 1$$

...

$$w(\text{aux}_1) = \theta_1$$

$$w(\text{aux}_2) = \theta_2$$

# WMC for safe AI

**Our model:** a (parametric) PLP  $P$

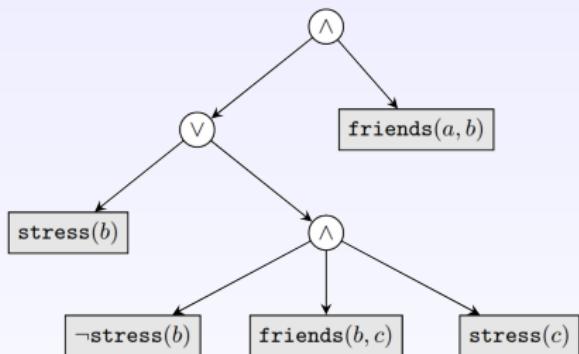
1) Conversion in weighted prop. logic

$$P \rightarrow P_{\text{ground}} \rightarrow \langle \Delta, w \rangle$$

2) Knowledge compilation:

$$\Delta \rightarrow C_\Delta$$

- $WMC(\Delta, w)$  computed in  $\Theta(|C_\Delta|)$
- $C_\Delta$  is differentiable w.r.t.  $w$



**Our model:** a (parametric) PLP  $P$

1) Conversion in weighted prop. logic

$$P \rightarrow P_{\text{ground}} \rightarrow \langle \Delta, w \rangle$$

2) Knowledge compilation:

$$\Delta \rightarrow C_\Delta$$

- $WMC(\Delta, w)$  computed in  $\Theta(|C_\Delta|)$
- $C_\Delta$  is differentiable w.r.t.  $w$

| Query ▾                 | Probability |
|-------------------------|-------------|
| funds(fwo,paolo)        | 0           |
| influences(paolo,pollo) | 0.2         |
| smokes(paolo)           | 0.3         |

Inference (and learning) in a model that  
satisfy constraints by construction!

# WMC for safe AI

WMC-based inference can be expensive (even with neural predicates)..

# WMC for safe AI

WMC-based inference can be expensive (even with neural predicates)..

..what if we can't afford it at inference time?

# WMC for safe AI

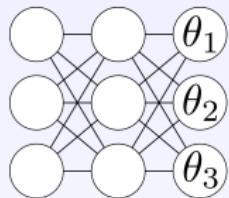
WMC-based inference can be expensive (even with neural predicates)..

..what if we can't afford it at inference time?

Can we still leverage our background knowledge  $\Delta$ ?

**Our model:** a differentiable function

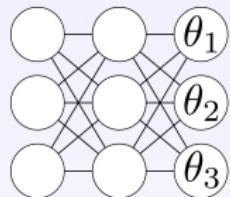
$$f(\mathbf{x}) = \Pr(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \text{ where } \Pr(Y_i = \top|\mathbf{X}) = \theta_i$$



**Our model:** a differentiable function

$$f(\mathbf{x}) = \Pr(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \text{ where } \Pr(Y_i = \top|\mathbf{X}) = \theta_i$$

Regularize  $f$  w.r.t. a constraint  $\Delta$  over  $\mathbf{Y}$ :



$\mathbf{Y} \models \Delta?$

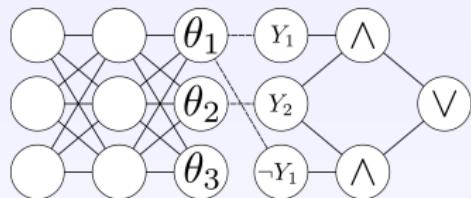
$$\begin{aligned} \text{SL}_\Delta(\Theta) &= -\log \text{WMC}(\Delta, \Theta) \\ &= -\log \sum_{\mu \models \Delta} \prod_{Y_i \in \mu} \theta_i \prod_{\neg Y_j \in \mu} (1 - \theta_j) \end{aligned}$$

...it's a **big** log-polynomial!

**Our model:** a differentiable function

$$f(\mathbf{x}) = \Pr(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \text{ where } \Pr(Y_i = \top|\mathbf{X}) = \theta_i$$

Regularize  $f$  w.r.t. a constraint  $\Delta$  over  $\mathbf{Y}$ :



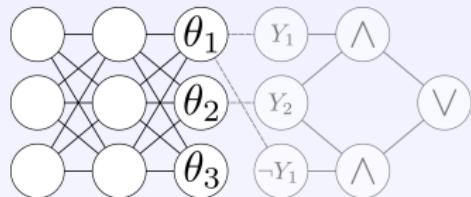
$$\begin{aligned} \text{SL}_\Delta(\Theta) &= -\log \text{WMC}(\Delta, \Theta) \\ &= -\log \sum_{\mu \models \Delta} \prod_{Y_i \in \mu} \theta_i \prod_{\neg Y_j \in \mu} (1 - \theta_j) \end{aligned}$$

but we can use KC!

**Our model:** a differentiable function

$$f(\mathbf{x}) = \Pr(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \text{ where } \Pr(Y_i = \top|\mathbf{X}) = \theta_i$$

Regularize  $f$  w.r.t. a constraint  $\Delta$  over  $\mathbf{Y}$ :

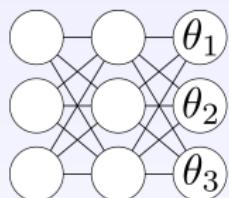


$$\begin{aligned} \text{SL}_\Delta(\Theta) &= -\log \text{WMC}(\Delta, \Theta) \\ &= -\log \sum_{\mu \models \Delta} \prod_{Y_i \in \mu} \theta_i \prod_{\neg Y_j \in \mu} (1 - \theta_j) \end{aligned}$$

Once  $f$  is trained, we can discard  $C_\Delta$

**Our model:** a differentiable function

$$f(\mathbf{x}) = \Pr(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \text{ where } \Pr(Y_i = \top|\mathbf{X}) = \theta_i$$



$$\mathbf{Y} \models \Delta!$$

Regularize  $f$  w.r.t. a constraint  $\Delta$  over  $\mathbf{Y}$ :

$$\begin{aligned} \text{SL}_\Delta(\Theta) &= -\log \text{WMC}(\Delta, \Theta) \\ &= -\log \sum_{\mu \models \Delta} \prod_{Y_i \in \mu} \theta_i \prod_{\neg Y_j \in \mu} (1 - \theta_j) \end{aligned}$$

Once  $f$  is trained, we can discard  $C_\Delta$

Efficient inference that

satisfy constraints in expectation!

What does “*in expectation*” mean?

What does “*in expectation*” mean?

Point-wise evaluations can be misleading..

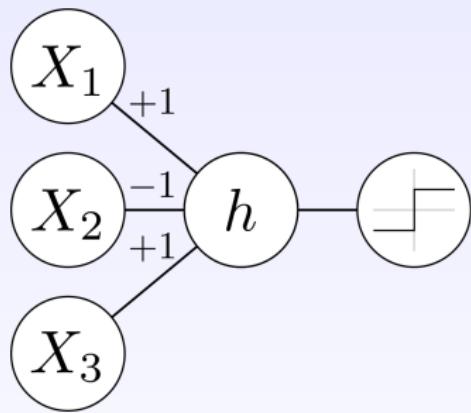
What does "*in expectation*" mean?

Point-wise evaluations can be misleading..

..how can we be sure that  $\Delta$  is satisfied enough in the real world?

**Our model:** a binarized NN

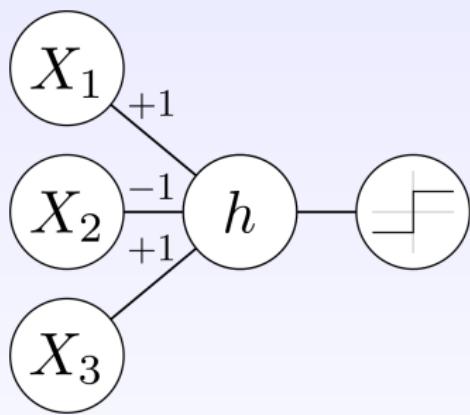
- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations



$$Y = \text{sign}(\langle w, \mathbf{X} \rangle + b)$$

**Our model:** a binarized NN

- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations

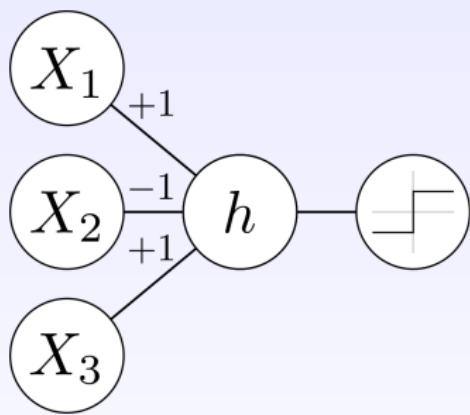


1) Conversion  $\mathcal{N} \rightarrow \Delta$

$$Y = \text{sign}(\langle w, \mathbf{X} \rangle + b)$$

**Our model:** a binarized NN

- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations



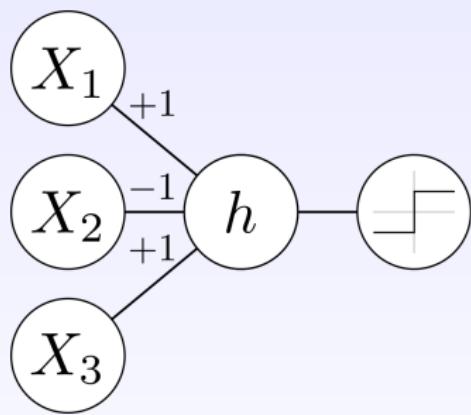
1) Conversion  $\mathcal{N} \rightarrow \Delta$

- *trivial* encoding  $\Delta$  is exponentially large

$$Y = \text{sign}(\langle w, \mathbf{X} \rangle + b)$$

**Our model:** a binarized NN

- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations



1) Conversion  $\mathcal{N} \rightarrow \Delta$

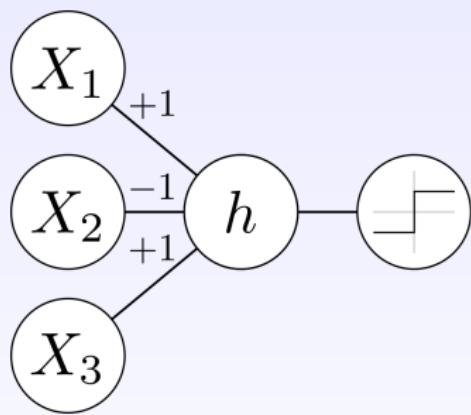
- *trivial* encoding  $\Delta$  is exponentially large
- *equicardinal*  $\Delta'$  is found using MILP

$$\#SAT(\Delta) = \#SAT(\Delta')$$

$$Y = \text{sign} (\langle w, \mathbf{X} \rangle + b)$$

**Our model:** a binarized NN

- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations



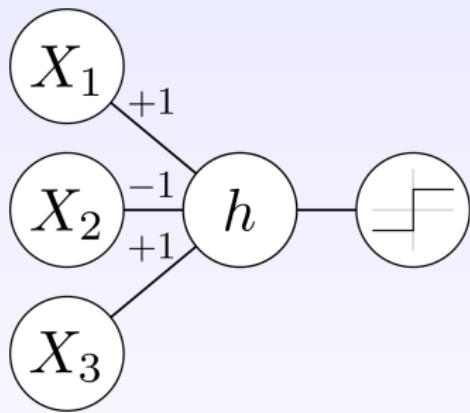
$$Y = \text{sign}(\langle w, \mathbf{X} \rangle + b)$$

1) Conversion  $\mathcal{N} \rightarrow \Delta$

- *trivial* encoding  $\Delta$  is exponentially large
- *equicardinal*  $\Delta'$  is found using MILP

$$\#SAT(\Delta) = \#SAT(\Delta')$$

2) Prior  $P(\mathbf{X})$ ? Reduction WMC  $\rightarrow \#SAT$



$$Y = \text{sign}(\langle w, \mathbf{X} \rangle + b)$$

**Our model:** a binarized NN

- $\mathcal{N} : \mathbb{Z}^N \rightarrow \mathbb{B}^M$
- weights  $w \in \{-1, 1\}$  and step activations

1) Conversion  $\mathcal{N} \rightarrow \Delta$

- *trivial* encoding  $\Delta$  is exponentially large
- *equicardinal*  $\Delta'$  is found using MILP

$$\#SAT(\Delta) = \#SAT(\Delta')$$

2) Prior  $P(\mathbf{X})$ ? Reduction WMC  $\rightarrow \#SAT$

Approximate  $\langle \epsilon, \delta \rangle$ -counting for

**verifying** quantitative properties!

# WMC for safe AI

In discrete settings, WMC can be used for:

learning models that satisfy constraints by construction

learning models that satisfy constraints in expectation

verifying quantitative properties of learned models

# WMC for safe AI

In discrete settings, WMC can be used for:

learning models that satisfy constraints by construction

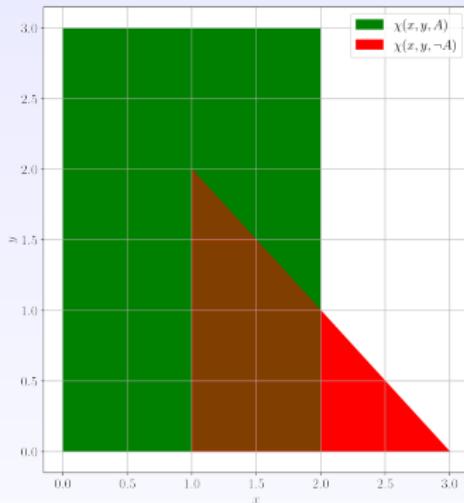
learning models that satisfy constraints in expectation

verifying quantitative properties of learned models

..but what about continuous or hybrid models?

# Generalizing WMC

SMT formulas  
w/ algebraic constraints

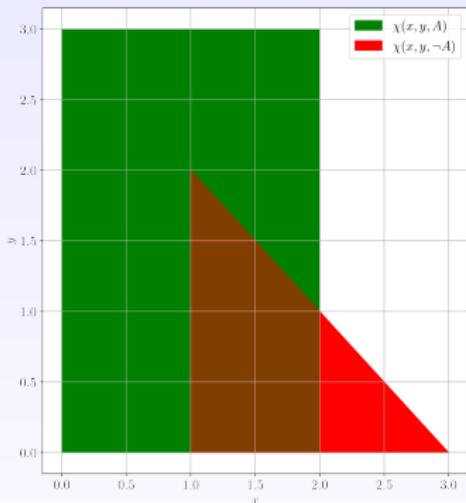


$$\begin{aligned}\chi = & (0 \leq y) \wedge (y \leq 3) \wedge (0 \leq x) \\ \wedge & (A \rightarrow (x \leq 2)) \\ \wedge & (\neg A \rightarrow ((1 \leq x) \wedge (x + y \leq 3)))\end{aligned}$$

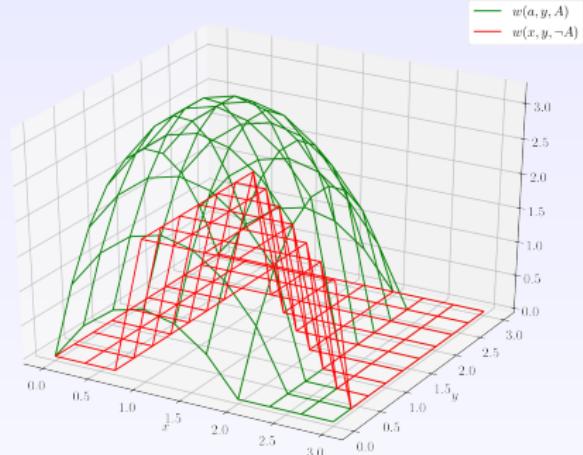
# Generalizing WMC

SMT formulas  
w/ algebraic constraints

+ density functions  
over continuous variables



$$\begin{aligned}\chi = & (0 \leq y) \wedge (y \leq 3) \wedge (0 \leq x) \\ \wedge & (\textcolor{green}{A} \rightarrow (x \leq 2)) \\ \wedge & (\neg A \rightarrow ((1 \leq x) \wedge (x + y \leq 3)))\end{aligned}$$



$$\begin{aligned}w(x, y, A) = & \llbracket A \rrbracket (-x^2 - y^2 + 2x + 3y) \\ & + \llbracket \neg A \rrbracket (-2x - 2y + 6)\end{aligned}$$

# Weighted Model Integration

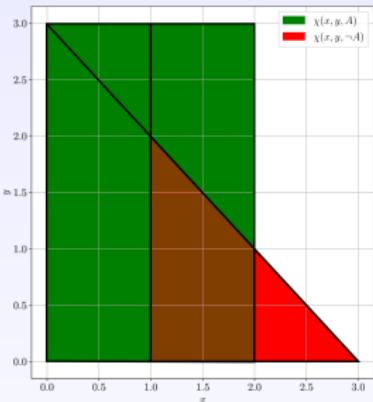
$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu)$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$

# Weighted Model Integration

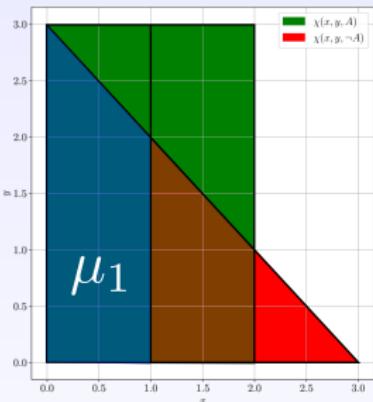
$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



$$WMI(\chi, w) =$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



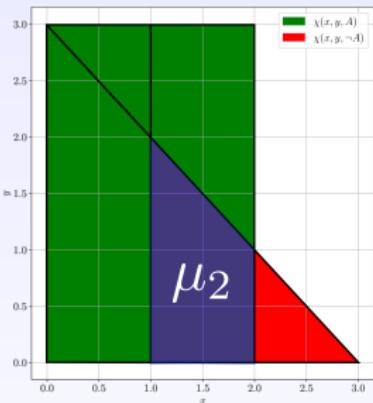
$$\mu_1 = \begin{cases} A \\ \neg(1 \leq x) \\ (x \leq 2) \\ (x + y \leq 3) \\ \dots \end{cases}$$

$$w(\mu_1) = \int_0^1 \int_0^{3-x} -x^2 - y^2 + 2x + 3y \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) +$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



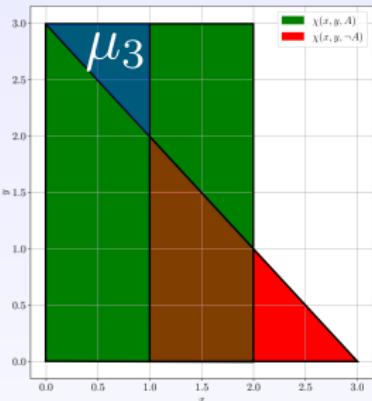
$$\mu_2 = \begin{cases} A & (1 \leq x) \\ & (x \leq 2) \\ & (x + y \leq 3) \\ & \dots \end{cases}$$

$$w(\mu_2) = \int_1^2 \int_0^{3-x} -x^2 - y^2 + 2x + 3y \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) + w(\mu_2) +$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



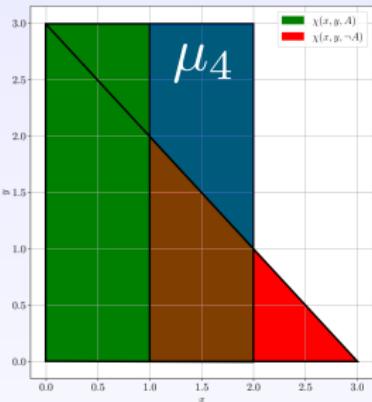
$$\mu_3 = \begin{cases} A \\ \neg(1 \leq x) \\ (x \leq 2) \\ \neg(x + y \leq 3) \\ \dots \end{cases}$$

$$w(\mu_3) = \int_0^1 \int_{3-x}^3 -x^2 - y^2 + 2x + 3y \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) + w(\mu_2) + w(\mu_3) +$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



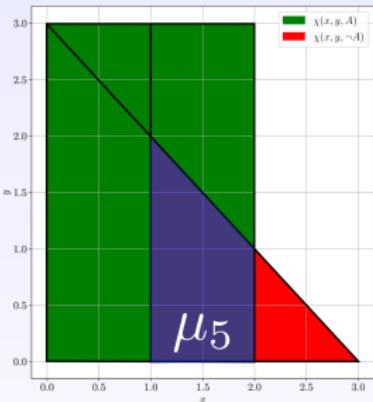
$$\mu_4 = \begin{cases} A & (1 \leq x) \\ & (x \leq 2) \\ & \neg(x + y \leq 3) \\ & \dots \end{cases}$$

$$w(\mu_4) = \int_1^2 \int_{3-x}^3 -x^2 - y^2 + 2x + 3y \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) + w(\mu_2) + w(\mu_3) + w(\mu_4) +$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



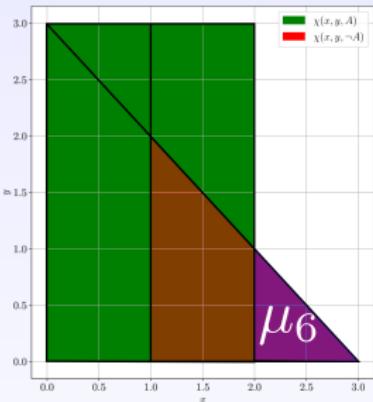
$$\mu_5 = \begin{cases} \neg A \\ (1 \leq x) \\ (x \leq 2) \\ \neg(x + y \leq 3) \\ \dots \end{cases}$$

$$w(\mu_5) = \int_1^2 \int_0^{3-x} -2x - 2y + 6 \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) + w(\mu_2) + w(\mu_3) + w(\mu_4) + w(\mu_5) +$$

# Weighted Model Integration

$$WMI(\Delta, w) = \sum_{\mu \models \Delta} w(\mu) = \sum_{\mu \models \Delta} \int_{\mu} f(\mathbf{x}) d\mathbf{x}$$



$$\mu_6 = \begin{cases} \neg A \\ (1 \leq x) \\ \neg(x \leq 2) \\ \neg(x + y \leq 3) \\ \dots \end{cases}$$

$$w(\mu_6) = \int_2^3 \int_0^{3-x} -2x - 2y + 6 \, dy \, dx$$

$$WMI(\chi, w) = w(\mu_1) + w(\mu_2) + w(\mu_3) + w(\mu_4) + w(\mu_5) + w(\mu_6)$$

# Weighted Model Integration

Two subtasks:

# Weighted Model Integration

Two subtasks:

- 1) **Enumerating** convex integration regions  $\mu$  ( $\lambda$ -SMT)  
→ as hard as #SAT (#P-complete)

# Weighted Model Integration

Two subtasks:

- 1) **Enumerating** convex integration regions  $\mu$  ( $\lambda$ -SMT)

→ as hard as #SAT (#P-complete)

- 2) **Continuous integration**

→ as hard as computing the volume of a polytope (#P-hard)

# Weighted Model Integration

Two subtasks:

1) **Enumerating** convex integration regions  $\mu$  ( $\lambda$ -SMT)

→ as hard as #SAT (#P-complete)

2) **Continuous integration**

→ as hard as computing the volume of a polytope (#P-hard)

The **combination** of 1) and 2) makes it very tricky

→ successful ideas in WMC do not always apply

# Weighted Model Integration

Two subtasks:

1) **Enumerating** convex integration regions  $\mu$  ( $\lambda$ -SMT)

→ as hard as #SAT (#P-complete)

2) **Continuous integration**

→ as hard as computing the volume of a polytope (#P-hard)

The **combination** of 1) and 2) makes it very tricky

→ successful ideas in WMC do not always apply

Why even bother?

# Weighted Model Integration

Two subtasks:

1) **Enumerating** convex integration regions  $\mu$  ( $\lambda$ -SMT)

→ as hard as #SAT (#P-complete)

2) **Continuous integration**

→ as hard as computing the volume of a polytope (#P-hard)

The **combination** of 1) and 2) makes it very tricky

→ successful ideas in WMC do not always apply

Why even bother?

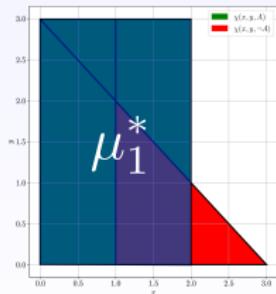
- This is a **very general** framework!
- Many inference algorithms are **specific** polytime WMI “in disguise”

# Computing WMI efficiently

$$\begin{aligned} WMI(\chi, w) = & \int_0^1 \int_0^{3-x} f_A(x, y) dy dx + \int_1^2 \int_0^{3-x} f_A(x, y) dy dx \\ & + \int_0^1 \int_{3-x}^3 f_A(x, y) dy dx + \int_1^2 \int_{3-x}^3 f_A(x, y) dy dx \\ & + \int_1^2 \int_0^{3-x} f_{\neg A}(x, y) dy dx + \int_2^3 \int_0^{3-x} f_{\neg A}(x, y) dy dx \end{aligned}$$

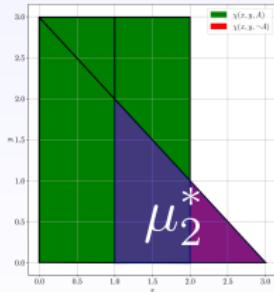
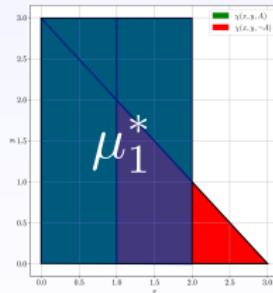
# Computing WMI efficiently

$$\begin{aligned} WMI(\chi, w) &= \int_0^1 \int_0^{3-x} f_A(x, y) dy dx + \int_1^2 \int_0^{3-x} f_A(x, y) dy dx \\ &\quad + \int_0^1 \int_{3-x}^3 f_A(x, y) dy dx + \int_1^2 \int_{3-x}^3 f_A(x, y) dy dx \\ &\quad + \int_1^2 \int_0^{3-x} f_{\neg A}(x, y) dy dx + \int_2^3 \int_0^{3-x} f_{\neg A}(x, y) dy dx \\ &= \overbrace{\int_0^2 \int_0^3 f_A(x, y) dy dx}^{w(\mu_1^*)} \end{aligned}$$



# Computing WMI efficiently

$$\begin{aligned} WMI(\chi, w) &= \int_0^1 \int_0^{3-x} f_A(x, y) dy dx + \int_1^2 \int_0^{3-x} f_A(x, y) dy dx \\ &\quad + \int_0^1 \int_{3-x}^3 f_A(x, y) dy dx + \int_1^2 \int_{3-x}^3 f_A(x, y) dy dx \\ &\quad + \int_1^2 \int_0^{3-x} f_{\neg A}(x, y) dy dx + \int_2^3 \int_0^{3-x} f_{\neg A}(x, y) dy dx \\ &= \overbrace{\int_0^2 \int_0^3 f_A(x, y) dy dx}^{w(\mu_1^*)} + \overbrace{\int_1^3 \int_0^{3-x} f_{\neg A}(x, y) dy dx}^{w(\mu_2^*)} \end{aligned}$$



# Computing WMI efficiently

## How?

Knowledge compilation

SMT oracles

Tractable subclasses

Monte Carlo estimates

# Computing WMI efficiently

## How?

Knowledge compilation

SMT oracles

Tractable subclasses

Monte Carlo estimates

| Algorithm   | Combinatorial Enumeration |               | Integration |      |                   | Expressiveness  |             |
|-------------|---------------------------|---------------|-------------|------|-------------------|-----------------|-------------|
|             | Exact                     | Method        | Exact       | Sym. | Method            | Parametric form | Assumptions |
| WMI-CC      | ✓                         | DPLL-SMT      | ✓           |      | LattE + CC        | Mul. polynomial | UC          |
| WMI-PA      | ✓                         | DPLL-SMT      | ✓           |      | LattE             | Mul. polynomial | -           |
| PRAISE      | ✓                         | DPLL-PIMT     | ✓           | ✓    | PIMT              | Mul. polynomial | -           |
| SVE         | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| BR          | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| F-XSDD      | ✓                         | KC-XSDD       | ✓           | ✓    | XADD / PSI        | Mul. polynomial | -           |
| WMI-SDD     | ✓                         | KC-XSDD       | ✓           |      | Scipy / LattE     | Mul. polynomial | -           |
| Symbo       | ✓                         | KC-XSDD       | ✓           | ✓    | PSI               | Mul. Gaussian   | UC          |
|             | ✓                         | KC-XSDD       |             |      | MC                | Mul. Gaussian   | -           |
| SMI         | ✓                         | AND/OR search | ✓           | ✓    | univ. integration | Biv. monomials  | BC, CNF, A  |
| MP-WMI      | ✓                         | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF, A  |
|             |                           | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF     |
| AprxWMI-CNF |                           | hashing + SMT | ✓           |      | LattE             | Mul. polynomial | CNF         |
| AprxWMI-DNF |                           | FPRAS         | ✓           |      | LattE             | Mul. polynomial | DNF         |

# Computing WMI efficiently

## How?

Knowledge compilation

SMT oracles

Tractable subclasses

Monte Carlo estimates

| Algorithm        | Combinatorial Enumeration |               | Integration |      |                   | Expressiveness  |             |
|------------------|---------------------------|---------------|-------------|------|-------------------|-----------------|-------------|
|                  | Exact                     | Method        | Exact       | Sym. | Method            | Parametric form | Assumptions |
| WMI-CC           | ✓                         | DPLL-SMT      | ✓           |      | LattE + CC        | Mul. polynomial | UC          |
| WMI-PA           | ✓                         | DPLL-SMT      | ✓           |      | LattE             | Mul. polynomial | -           |
| PRAISE           | ✓                         | DPLL-PIMT     | ✓           | ✓    | PIMT              | Mul. polynomial | -           |
| SVE              | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| BR               | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| F-XSDD           | ✓                         | KC-XSDD       | ✓           | ✓    | XADD / PSI        | Mul. polynomial | -           |
| WMI-SDD          | ✓                         | KC-XSDD       | ✓           |      | Scipy / LattE     | Mul. polynomial | -           |
| Symbo<br>Sampo   | ✓                         | KC-XSDD       | ✓           | ✓    | PSI               | Mul. Gaussian   | UC          |
|                  | ✓                         | KC-XSDD       |             |      | MC                | Mul. Gaussian   | -           |
| SMI              | ✓                         | AND/OR search | ✓           | ✓    | univ. integration | Biv. monomials  | BC, CNF, A  |
| MP-WMI<br>ReCoIn | ✓                         | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF, A  |
|                  |                           | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF     |
| AprxWMI-CNF      |                           | hashing + SMT | ✓           |      | LattE             | Mul. polynomial | CNF         |
| AprxWMI-DNF      |                           | FPRAS         | ✓           |      | LattE             | Mul. polynomial | DNF         |

Research focused on theory and algorithms

# Computing WMI efficiently

## How?

- Knowledge compilation
- SMT oracles
- Tractable subclasses
- Monte Carlo estimates

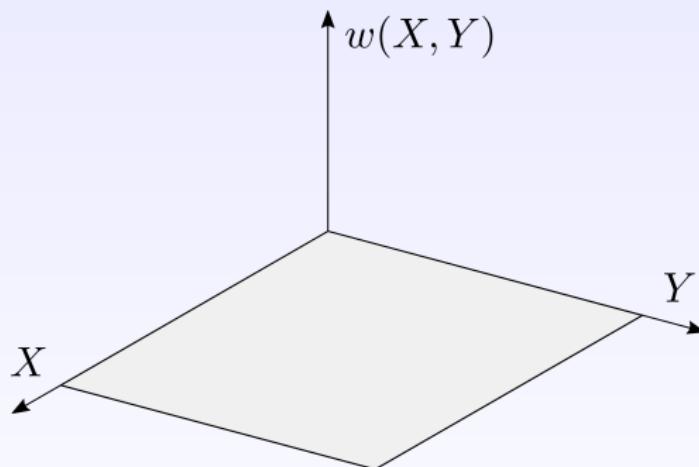
} guided by applications in safe AI

| Algorithm        | Combinatorial Enumeration |               | Integration |      |                   | Expressiveness  |             |
|------------------|---------------------------|---------------|-------------|------|-------------------|-----------------|-------------|
|                  | Exact                     | Method        | Exact       | Sym. | Method            | Parametric form | Assumptions |
| WMI-CC           | ✓                         | DPLL-SMT      | ✓           |      | LattE + CC        | Mul. polynomial | UC          |
| WMI-PA           | ✓                         | DPLL-SMT      | ✓           |      | LattE             | Mul. polynomial | -           |
| PRAISE           | ✓                         | DPLL-PIMT     | ✓           | ✓    | PIMT              | Mul. polynomial | -           |
| SVE              | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| BR               | ✓                         | KC-XADD       | ✓           | ✓    | XADD              | Mul. polynomial | -           |
| F-XSDD           | ✓                         | KC-XSDD       | ✓           | ✓    | XADD / PSI        | Mul. polynomial | -           |
| WMI-SDD          | ✓                         | KC-XSDD       | ✓           |      | Scipy / LattE     | Mul. polynomial | -           |
| Symbo<br>Sampo   | ✓                         | KC-XSDD       | ✓           | ✓    | PSI               | Mul. Gaussian   | UC          |
|                  | ✓                         | KC-XSDD       |             |      | MC                | Mul. Gaussian   | -           |
| SMI              | ✓                         | AND/OR search | ✓           | ✓    | univ. integration | Biv. monomials  | BC, CNF, A  |
| MP-WMI<br>ReCoIn | ✓                         | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF, A  |
|                  |                           | MP            | ✓           | ✓    | Sympy             | Biv. polynomial | BC, CNF     |
| AprxWMI-CNF      |                           | hashing + SMT | ✓           |      | LattE             | Mul. polynomial | CNF         |
| AprxWMI-DNF      |                           | FPRAS         | ✓           |      | LattE             | Mul. polynomial | DNF         |

Research focused on theory and algorithms

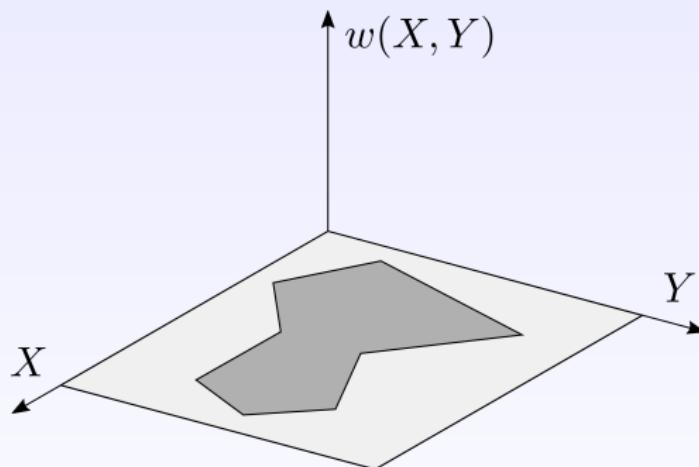
learning models that satisfy constraints by construction?

learning models that satisfy constraints by construction?



## Structure

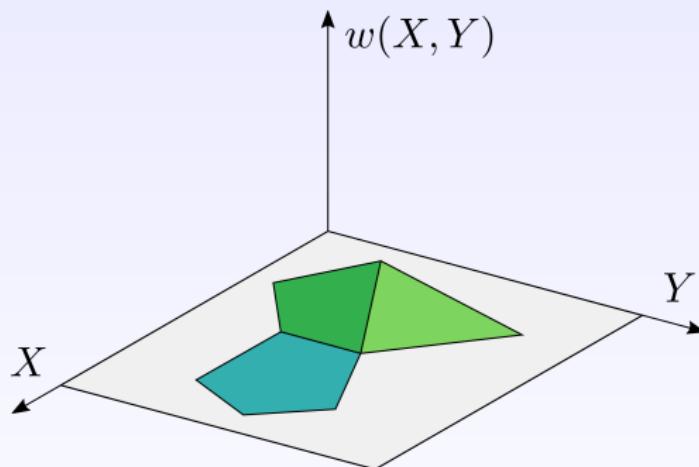
learning models that satisfy constraints by construction?



## Structure

→ support/constraints  $\Delta$  (*if not/partially given*)

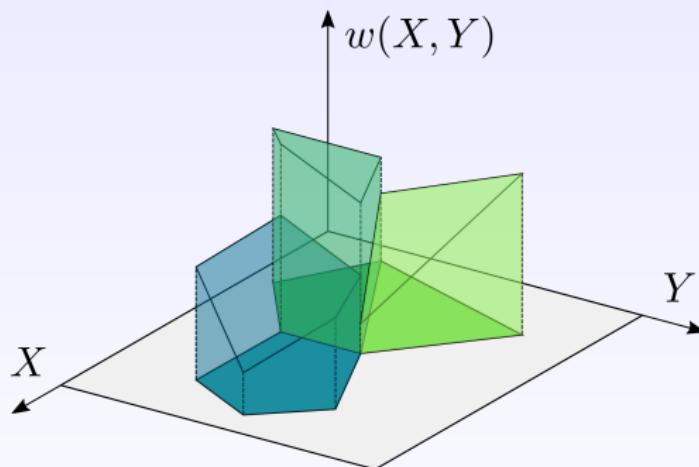
learning models that satisfy constraints by construction?



## Structure

- support/constraints  $\Delta$  (*if not/partially given*)
- $w(\mathbf{X}; \Theta) = [\text{if } \psi \text{ then } w_1(\mathbf{X}; \Theta_1) \text{ else } w_2(\mathbf{X}; \Theta_2)]$   
     $[w_1(\mathbf{X}_1; \Theta_1) \times w_2(\mathbf{X}_2; \Theta_2)]$   
     $[\theta_1 \cdot w_1(\mathbf{X}; \Theta_1) + \theta_2 \cdot w_2(\mathbf{X}; \Theta_2)]$

learning models that satisfy constraints by construction?



## Structure

- support/constraints  $\Delta$  (*if not/partially given*)
- $w(\mathbf{X}; \Theta) = [\text{if } \psi \text{ then } w_1(\mathbf{X}; \Theta_1) \text{ else } w_2(\mathbf{X}; \Theta_2)]$   
 $[w_1(\mathbf{X}_1; \Theta_1) \times w_2(\mathbf{X}_2; \Theta_2)]$   
 $[\theta_1 \cdot w_1(\mathbf{X}; \Theta_1) + \theta_2 \cdot w_2(\mathbf{X}; \Theta_2)]$

## Parameters

- $\operatorname{argmin}_{\Theta} \mathcal{L}(\Theta)$

learning models that satisfy constraints **in expectation?**

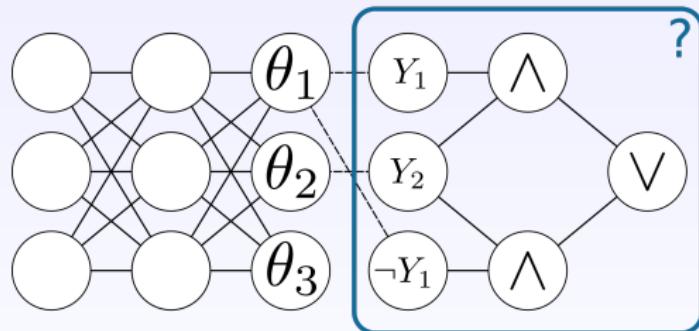
learning models that satisfy constraints **in expectation**?

## Efficient WMI evaluation

Tractable classes?

Knowledge compilation?

Sampling?



learning models that satisfy constraints in expectation?

## Efficient WMI evaluation

Tractable classes?

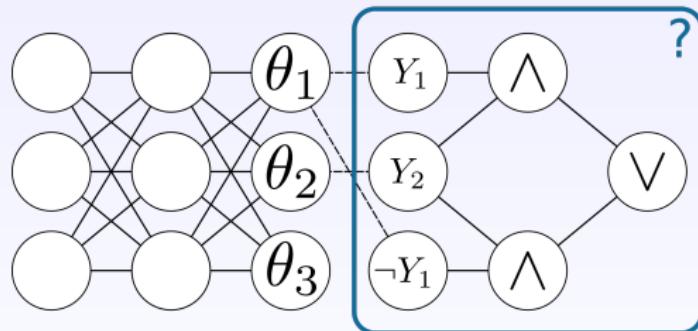
Knowledge compilation?

Sampling?

## Constraints

Logical + linear?

Nonlinear? ( $E = mc^2$ )



learning models that satisfy constraints in expectation?

## Efficient WMI evaluation

Tractable classes?

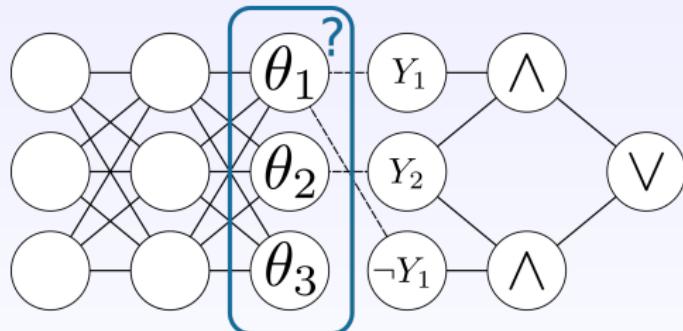
Knowledge compilation?

Sampling?

## Constraints

Logical + linear?

Nonlinear? ( $E = mc^2$ )



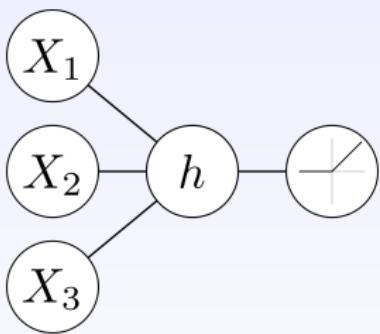
## What models to target

Piecewise polynomials?

Exponential family?

verifying quantitative properties?

## verifying quantitative properties?

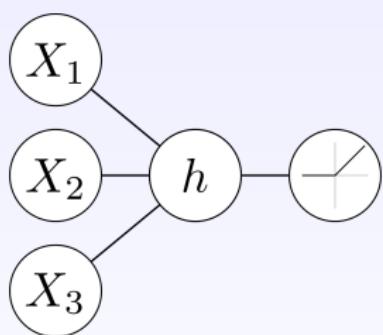


We can encode:

- Decision trees
- Support vector machines
- ReLU networks
- Sum-product networks
- ...

$$\begin{aligned}\Delta &= (h = \langle w, \mathbf{X} \rangle + b) \\ &\wedge ((h \leq 0) \rightarrow (Y = 0)) \\ &\wedge ((h > 0) \rightarrow (Y = h))\end{aligned}$$

## verifying quantitative properties?



We can encode:

- Decision trees
- Support vector machines
- ReLU networks
- Sum-product networks
- ...

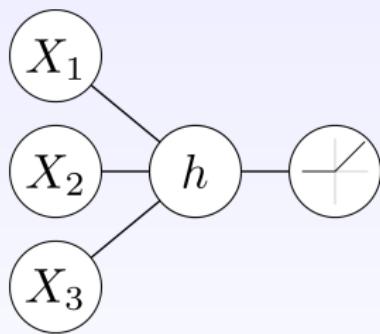
We can verify properties like:

$$\Pr(Hire|Female) = \Pr(Hire|Male)$$

(with arbitrary priors over  $\mathbf{X}$ !)

$$\begin{aligned}\Delta = & (h = \langle w, \mathbf{X} \rangle + b) \\ & \wedge ((h \leq 0) \rightarrow (Y = 0)) \\ & \wedge ((h > 0) \rightarrow (Y = h))\end{aligned}$$

## verifying quantitative properties?



We can encode:

- Decision trees
- Support vector machines
- ReLU networks
- Sum-product networks
- ...

We can verify properties like:

$$\Pr(\text{Hire}|\text{Female}) = \Pr(\text{Hire}|\text{Male})$$

(with arbitrary priors over  $\mathbf{X}$ !)

Can we scale?

Possibly, focusing on specific models/properties!

# The final slide

# The final slide

- WMI is a **very general framework** for constrained inference

# The final slide

- WMI is a **very general framework** for constrained inference

*...and it is not THAT scary!*

# The final slide

- WMI is a **very general framework** for constrained inference

*...and it is not THAT scary!*

- There are **many potential applications** of WMI in safe AI

# The final slide

- WMI is a **very general framework** for constrained inference

*...and it is not THAT scary!*

- There are **many potential applications** of WMI in safe AI

*...also many algorithmic challenges (I haven't talked about)*

# The final slide

- WMI is a **very general framework** for constrained inference

*...and it is not THAT scary!*

- There are **many potential applications** of WMI in safe AI

*...also many algorithmic challenges (I haven't talked about)*

**Thank you!**

*questions / feedback / collaborations are welcome!*