

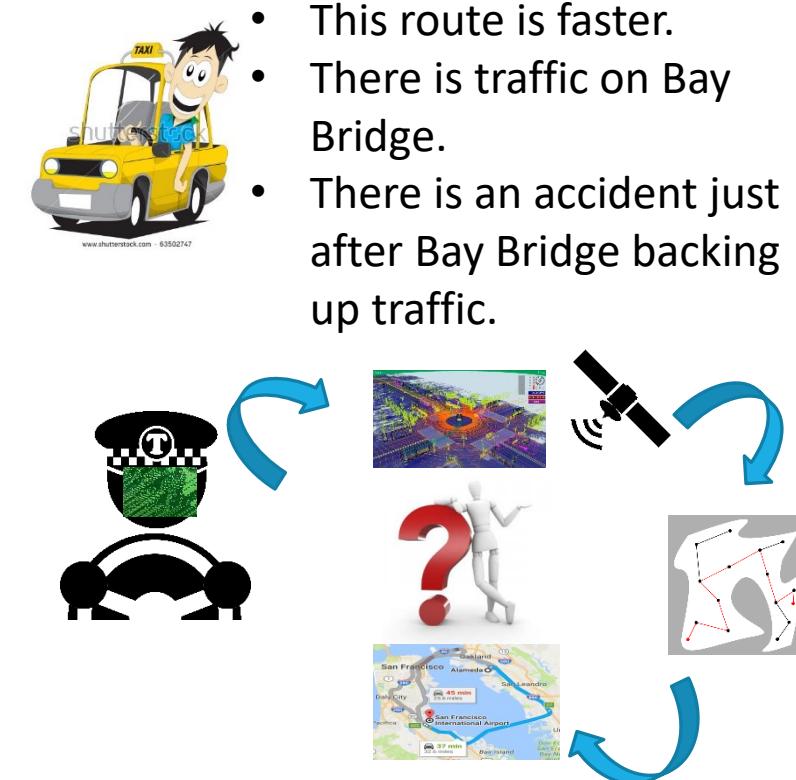
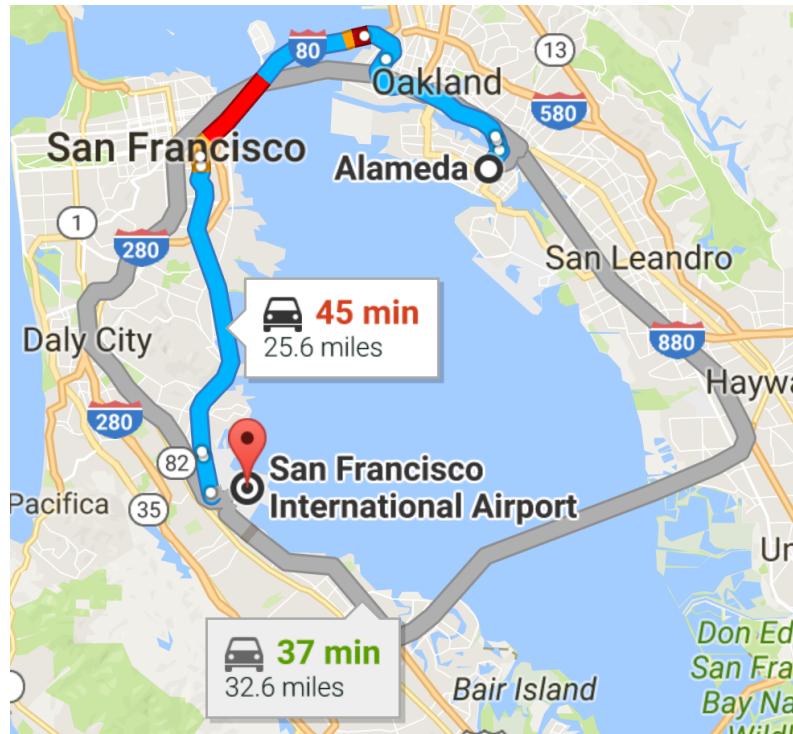
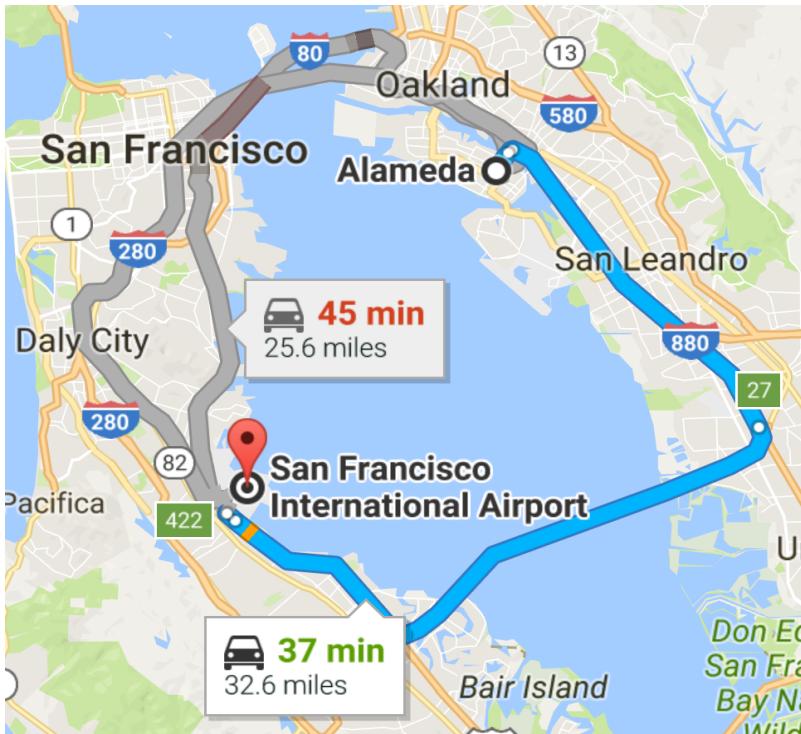
On Learning Sparse Boolean Formulae For Explaining AI Decisions

SUSMIT JHA, SRI INTERNATIONAL

VASUMATHI RAMAN, ZOOX INC.

ALESSANDRO PINTO, TUHIN SAHAI, AND MICHAEL FRANCIS
UNITED TECHNOLOGIES RESEARCH CENTER, BERKELEY

Ubiquitous AI and need for explanation

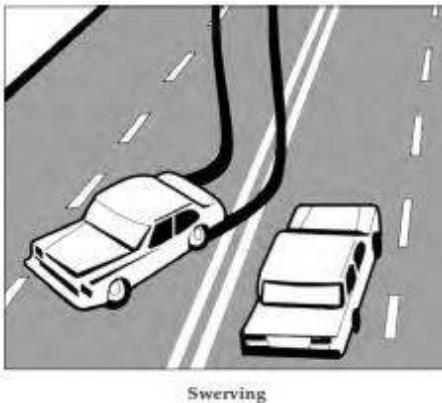


Why did we take the San Mateo bridge instead of the Bay Bridge ?

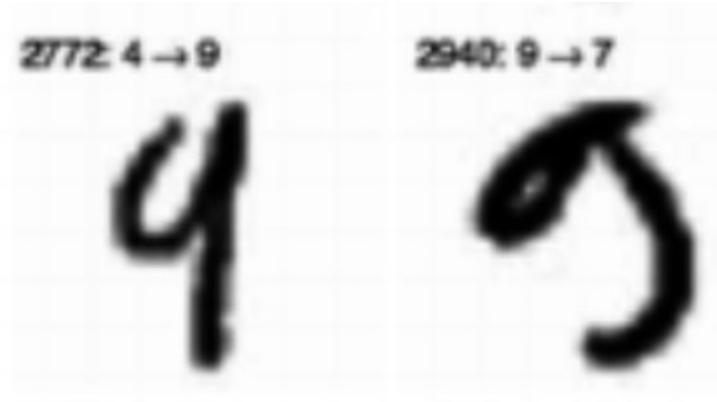
Decision/Recommendation/Classification



Medical Diagnosis



Autonomous Systems Certification

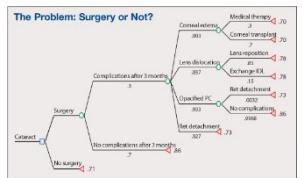


Incorrect ML Recommendations

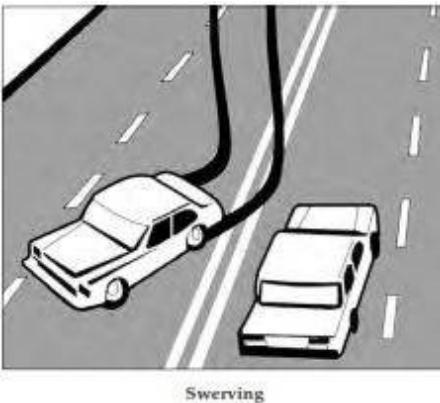
Decision/Recommendation/Classification



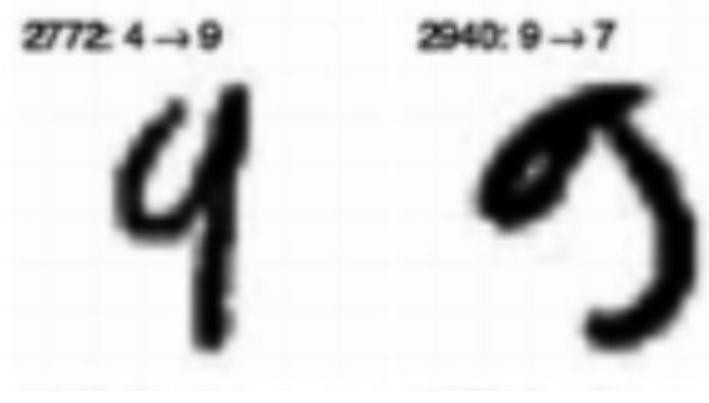
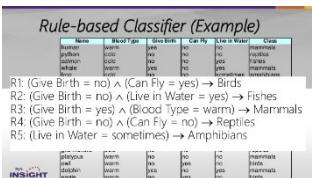
Medical Diagnosis



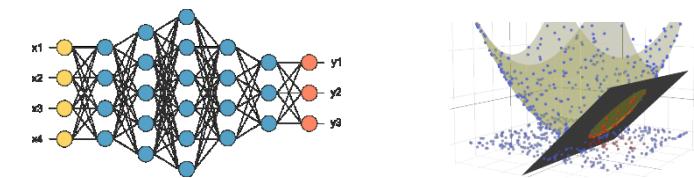
Interpretable but less scalable:
Decision Trees, Propositional Rules



Autonomous Systems Certification



Incorrect ML Recommendations



Scalable but less interpretable :
Neural Networks, Support Vector

Even ‘algorithmic’ decision making: A* Path Planning

A* is an algorithm that:

Uses heuristic to guide search

While ensuring that it will compute a path with
minimum cost

A* computes the function $f(n) = g(n) + h(n)$

$$f(n) = g(n) + h(n)$$

$g(n)$ = “cost from **the starting node** to reach n ”

$h(n)$ = “estimate of the cost of the cheapest
path from n to the **goal node**”

“estimated cost”

“actual cost”

Algorithm 1: A*

Input: $start, goal(n), h(n), expand(n)$

Output: path

```
1 if  $goal(start) = true$  then return  $makePath(start)$ 
2
3  $open \leftarrow start$ 
4  $closed \leftarrow \emptyset$ 
5 while  $open \neq \emptyset$  do
6   sort( $open$ )
7    $n \leftarrow open.pop()$ 
8   kids  $\leftarrow expand(n)$ 
9   forall the  $kid \in kids$  do
10     $kid.f \leftarrow (n.g + 1) + h(kid)$ 
11    if  $goal(kid) = true$  then return  $makePath(kid)$ 
12    if  $kid \cap closed = \emptyset$  then  $open \leftarrow kid$ 
13
14  $closed \leftarrow n$ 
15
16 return  $\emptyset$ 
```

Example: A* Path Planner

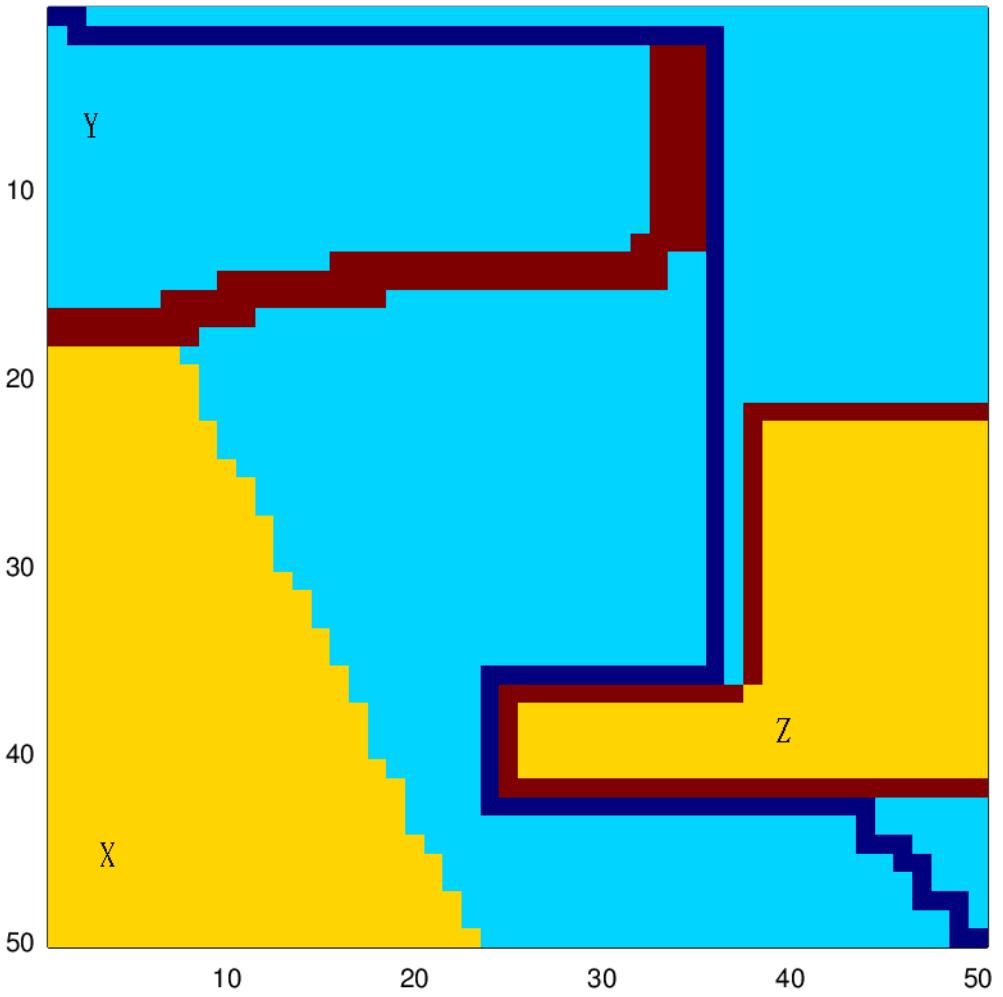


Algorithm 1: A*

Input: start, $goal(n)$, $h(n)$, $expand(n)$
Output: path

```
1 if  $goal(start) = true$  then return  $makePath(start)$ 
2
3  $open \leftarrow start$ 
4  $closed \leftarrow \emptyset$ 
5 while  $open \neq \emptyset$  do
6   sort( $open$ )
7    $n \leftarrow open.pop()$ 
8    $kids \leftarrow expand(n)$ 
9   forall the  $kid \in kids$  do
10     $kid.f \leftarrow (n.g + 1) + h(kid)$ 
11    if  $goal(kid) = true$  then return  $makePath(kid)$ 
12    if  $kid \cap closed = \emptyset$  then  $open \leftarrow kid$ 
13     $closed \leftarrow n$ 
14 return  $\emptyset$ 
```

Example: A* Path Planner



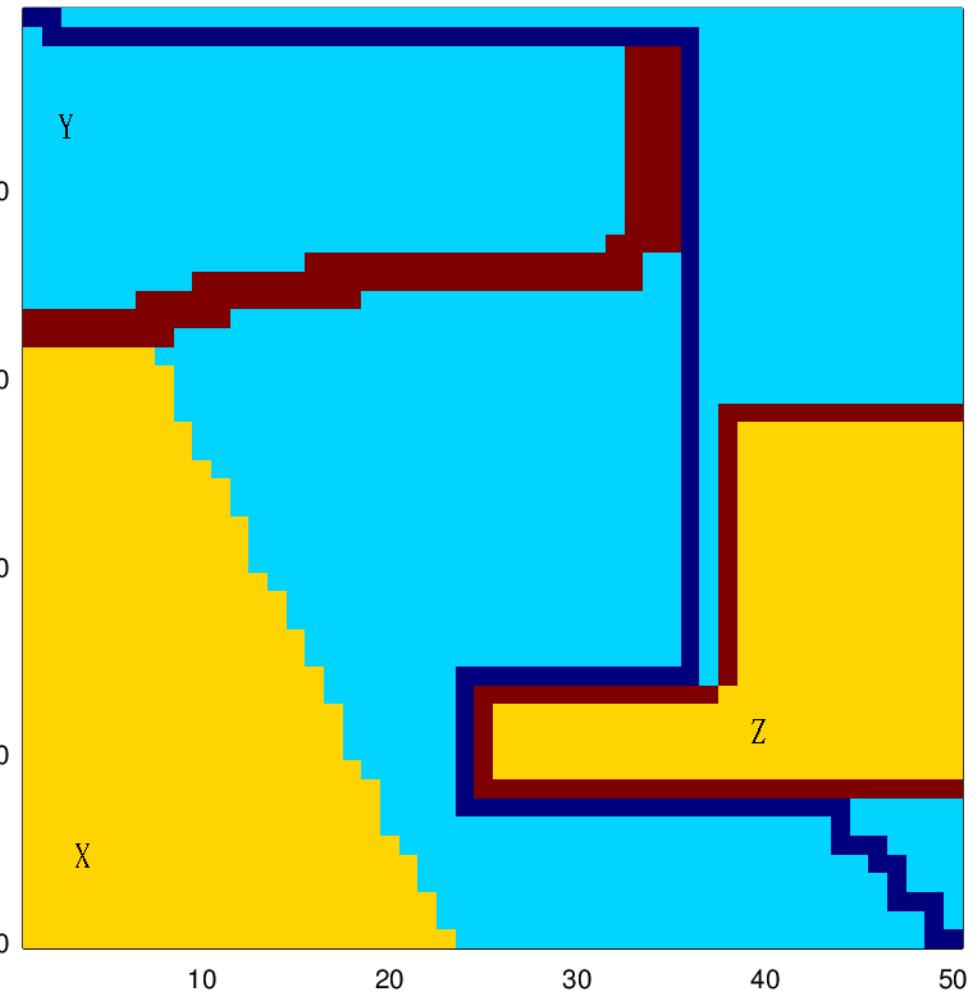
Algorithm 1: A*

Input: start , $\text{goal}(n)$, $h(n)$, $\text{expand}(n)$
Output: path

```
1 if  $\text{goal}(\text{start}) = \text{true}$  then return  $\text{makePath}(\text{start})$ 
2
3  $\text{open} \leftarrow \text{start}$ 
4  $\text{closed} \leftarrow \emptyset$ 
5 while  $\text{open} \neq \emptyset$  do
6   sort( $\text{open}$ )
7    $n \leftarrow \text{open.pop}()$ 
8    $\text{kids} \leftarrow \text{expand}(n)$ 
9   forall the  $\text{kid} \in \text{kids}$  do
10     $\text{kid}.f \leftarrow (n.g + 1) + h(\text{kid})$ 
11    if  $\text{goal}(\text{kid}) = \text{true}$  then return  $\text{makePath}(\text{kid})$ 
12    if  $\text{kid} \cap \text{closed} = \emptyset$  then  $\text{open} \leftarrow \text{kid}$ 
13    $\text{closed} \leftarrow n$ 
14 return  $\emptyset$ 
```

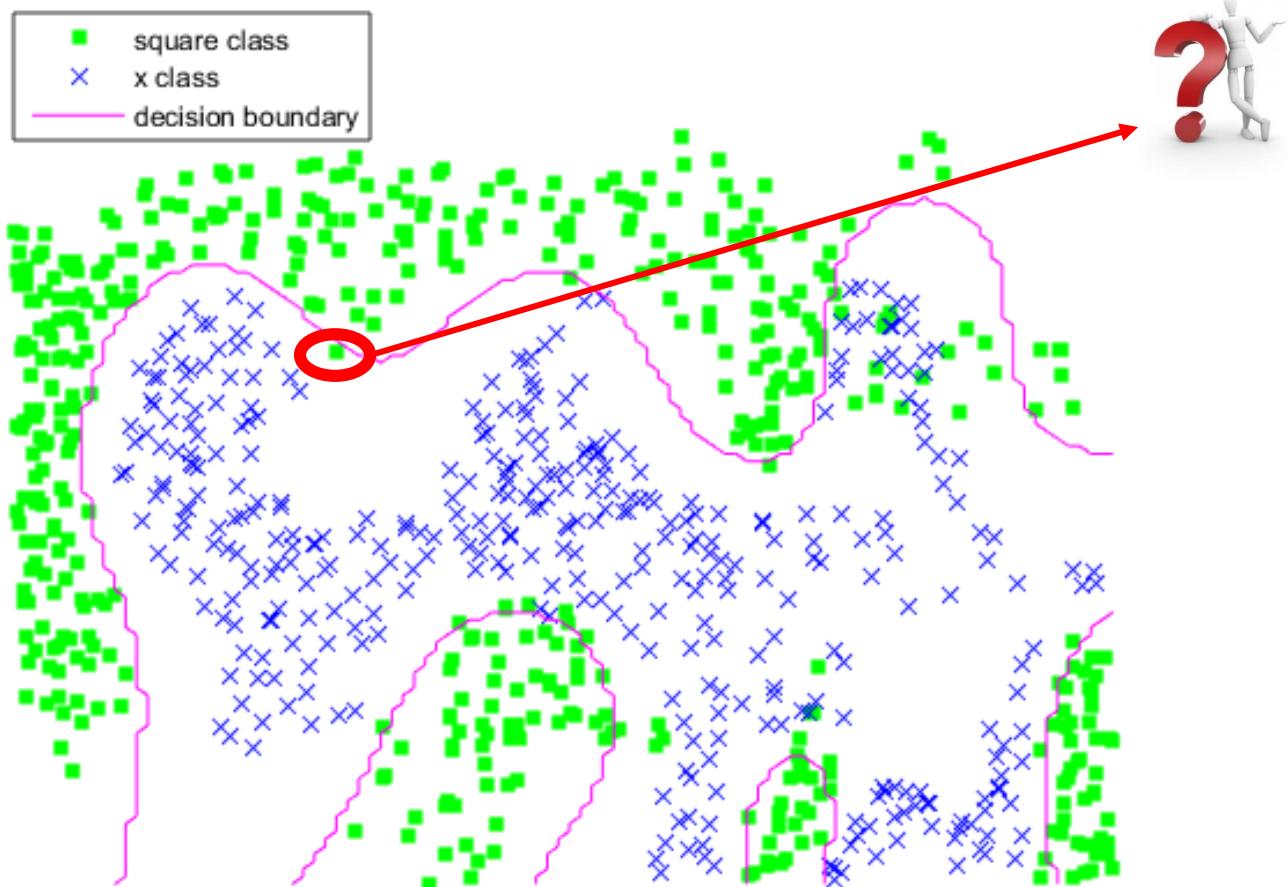
Why didn't we go through X / Z ?

Example: A* Path Planner

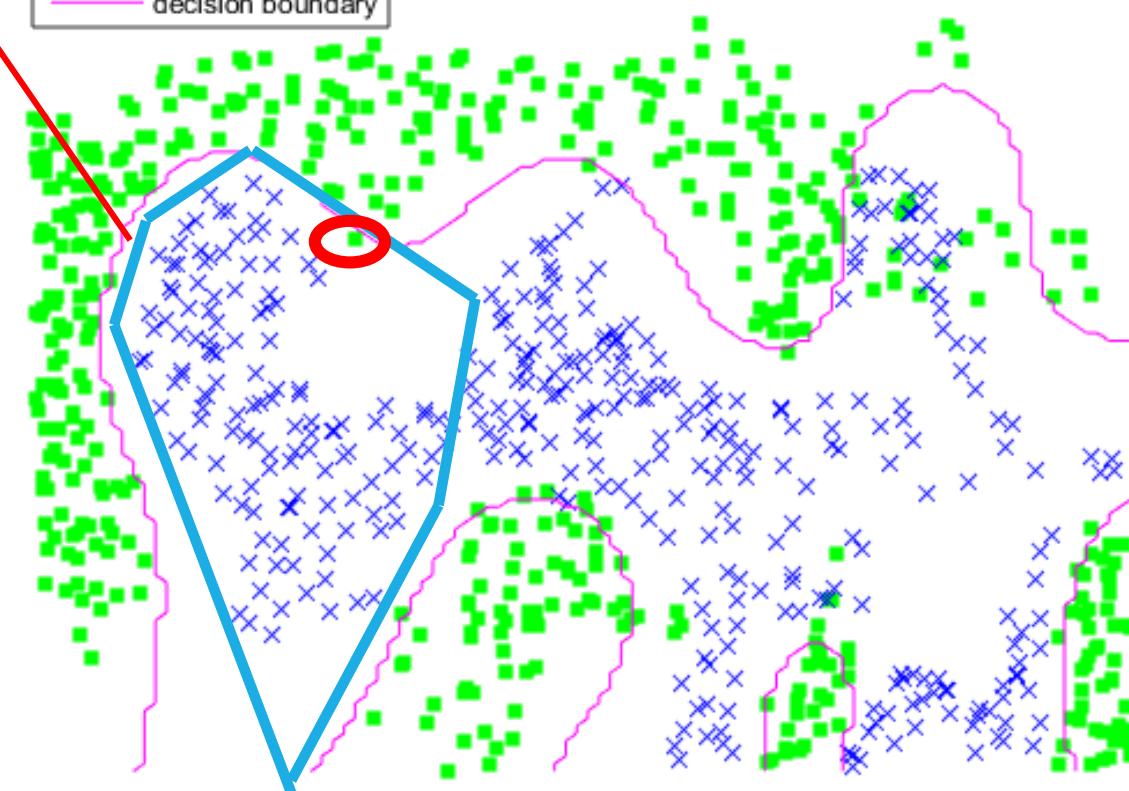
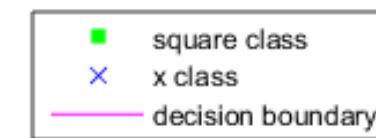
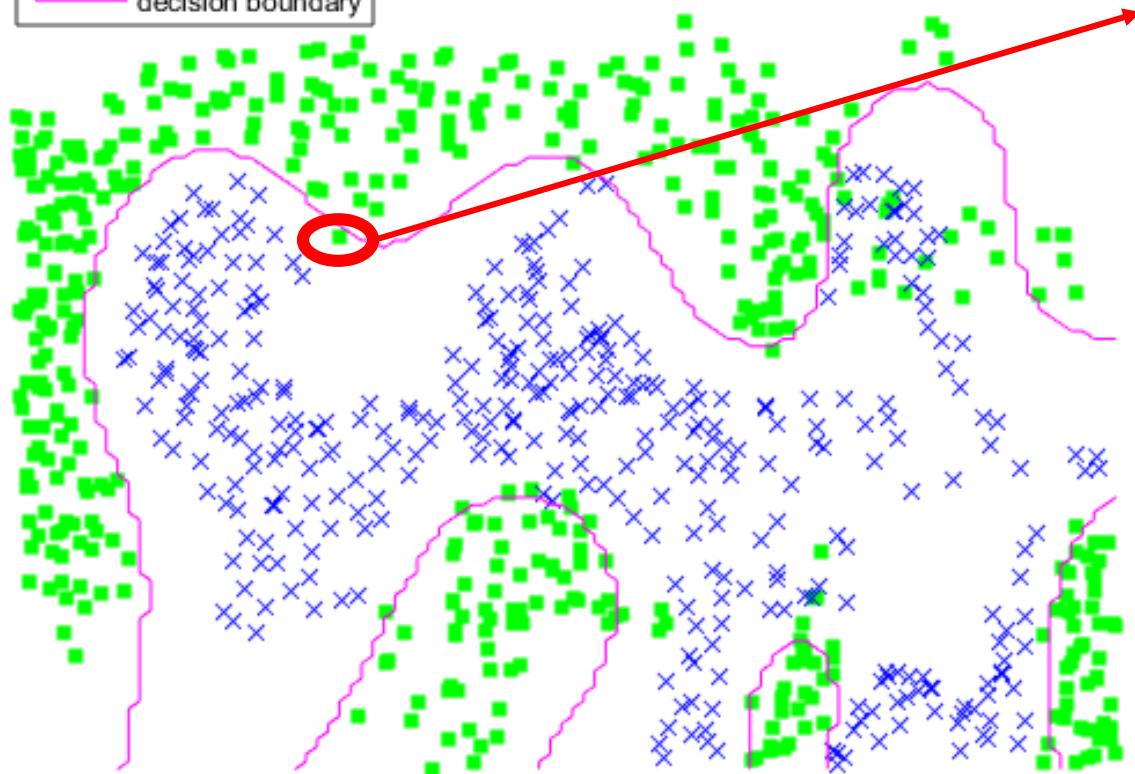
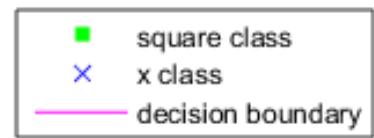


1. Internal details of algorithm and its implementation is unknown to human observer/user. So, explanation must be in terms of a **common vocabulary**.
2. Some explanations need further **deduction and inference that were not performed** by the AI algorithm while making the original decision.
3. Decision making is often accomplished through complex composition of a number of AI algorithms and hence, an explanation process would be practical only if it **did not require detailed modeling of each AI algorithm**.

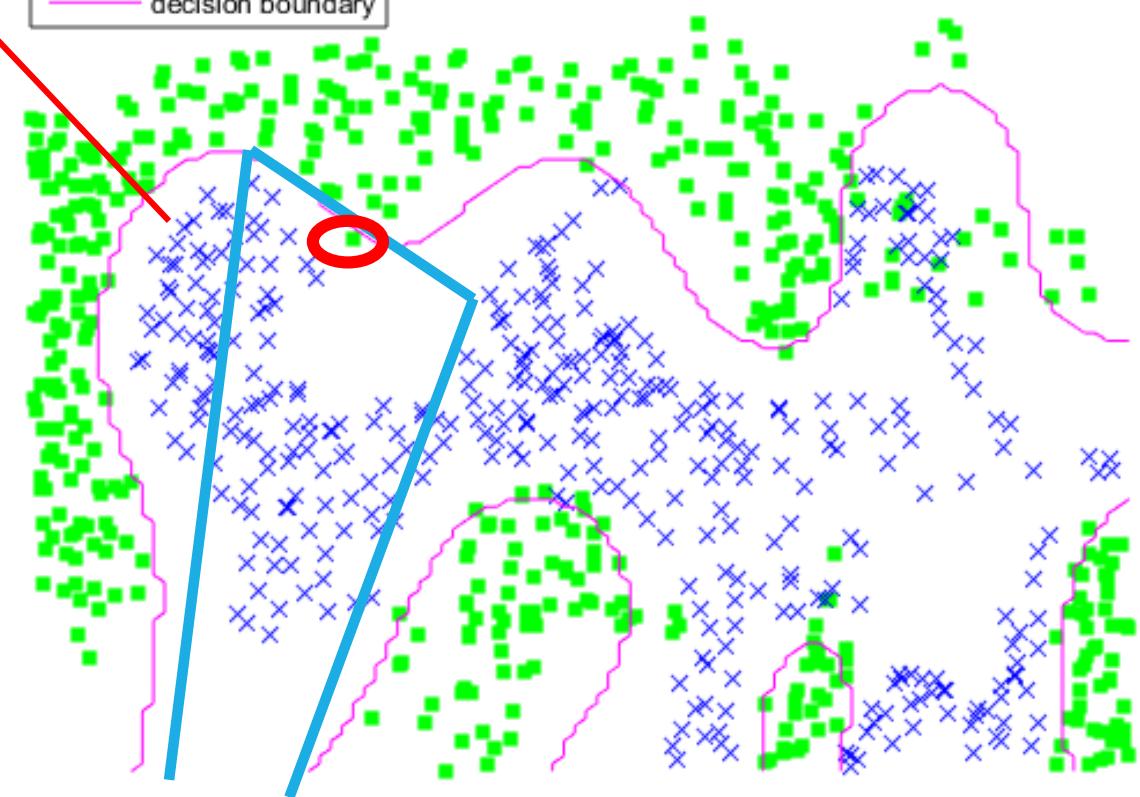
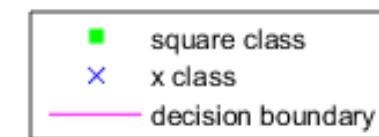
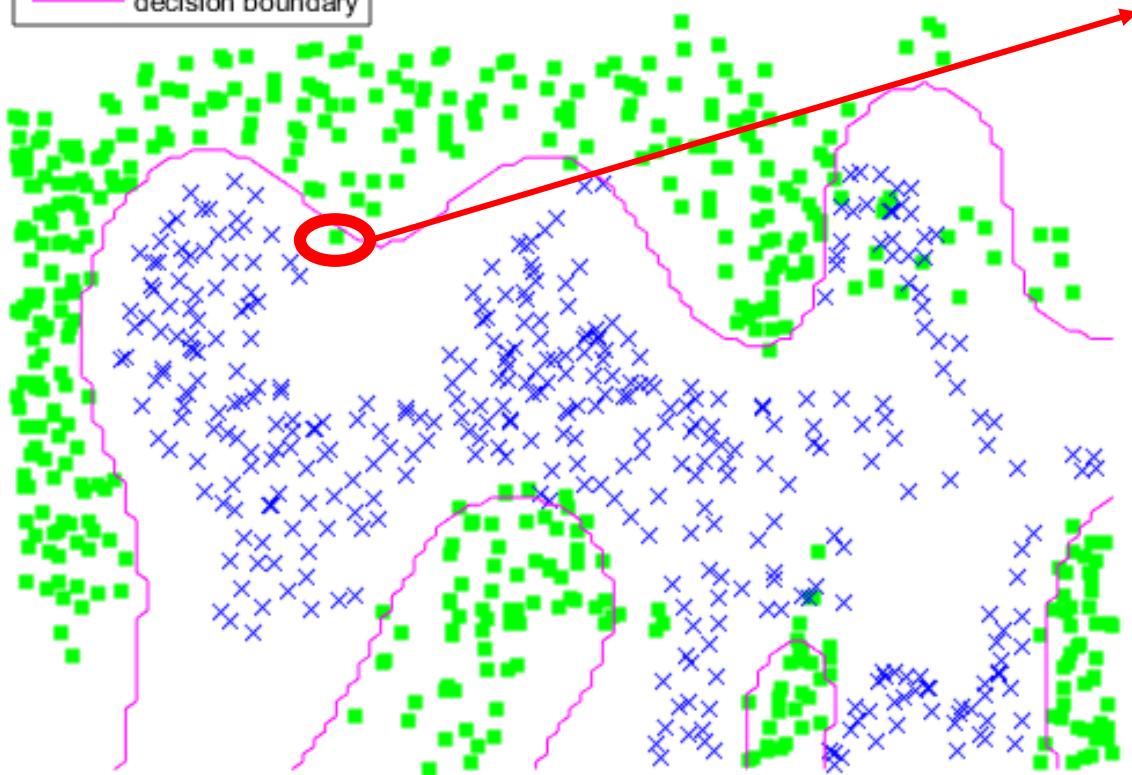
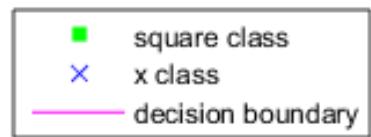
Local Explanations of Complex Models



Local Explanations of Complex Models



Local Explanations of Complex Models



Local Explanations in AI

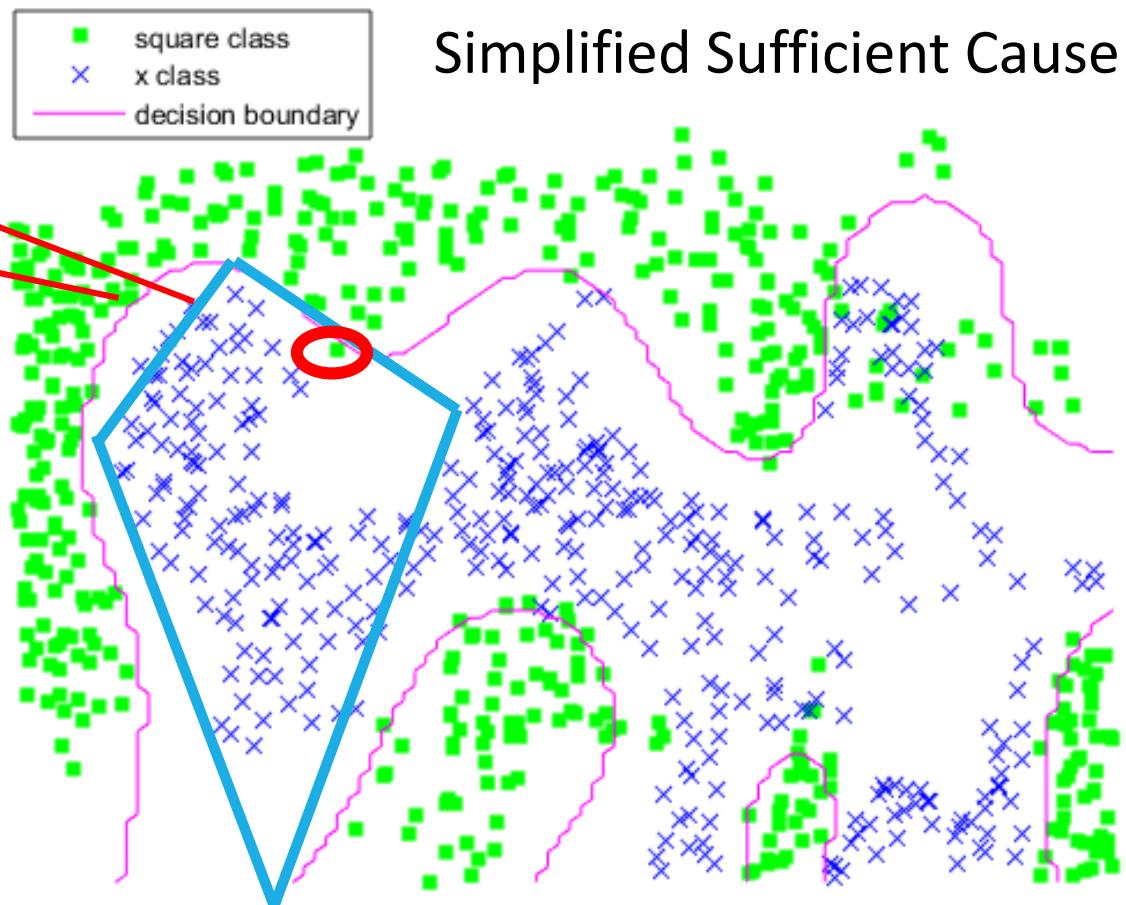
$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

$\mathcal{L}(f, g, \pi_x)$ Measure of how well g approximates f

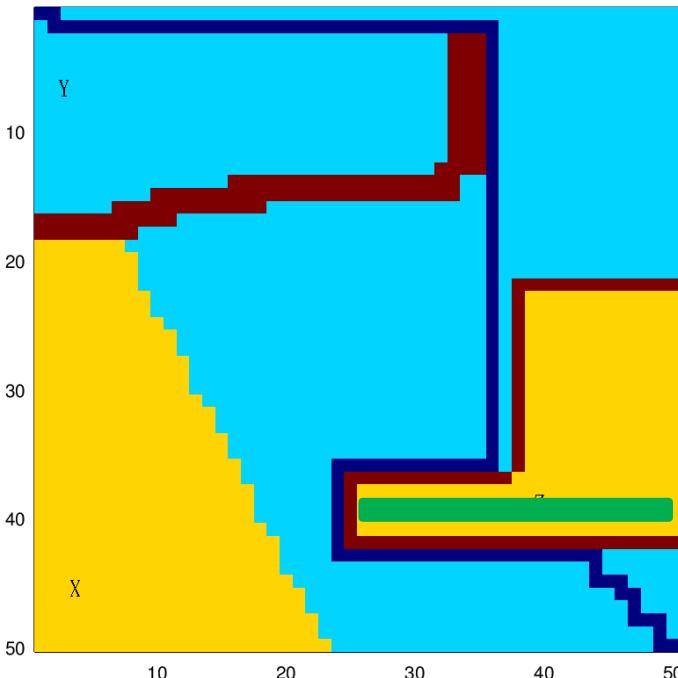
$\Omega(g)$ Measure of complexity of g

Formulation in AI:

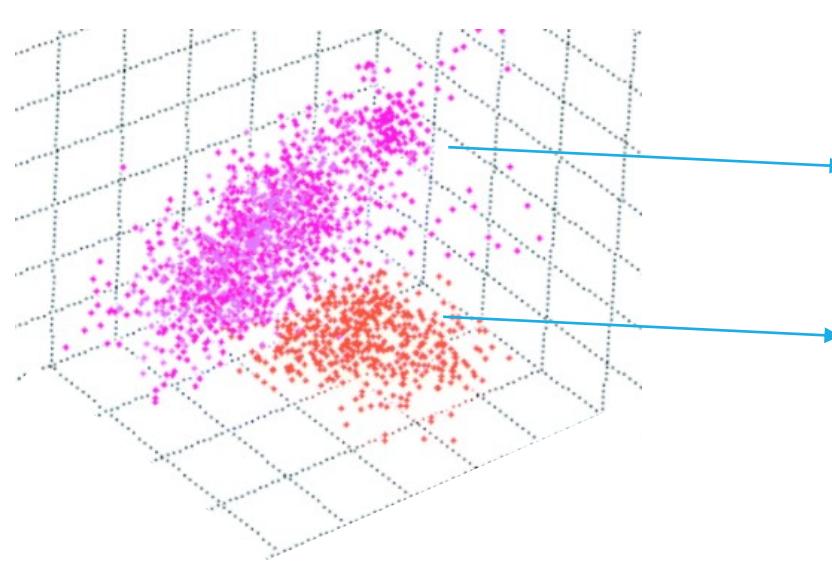
- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?: Explaining the Predictions of Any Classifier." *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016.
- Hayes, Bradley, and Julie A. Shah. "Improving Robot Controller Transparency Through Autonomous Policy Explanation." *Proceedings of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*. ACM, 2017.



Model Agnostic Explanation through Boolean Learning



Why does the path not go through Green?



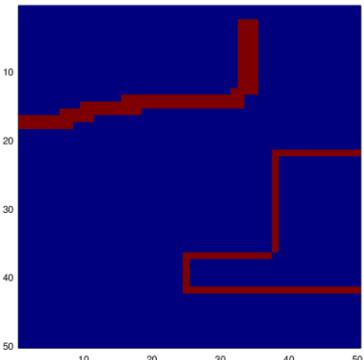
Let each point in k-dimensions (for some k) correspond to a map.

Maps in which optimum path goes via green
Maps in which optimum path does not go via green

Find a Boolean formula ϕ such that

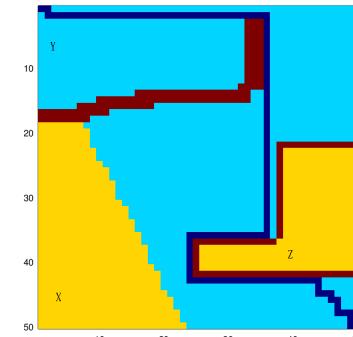
$$\begin{aligned}\phi &\Leftrightarrow \text{Path contain } z \\ \phi &\Rightarrow \text{Path contain } z\end{aligned}$$

Example: Explanations in A*



```
Algorithm 1: A*
Input: start, goal(n), h(n), openSet(n)
Output: null
1 if goal(start) = true then return makePath(start)
2
3 open ← start
4 closed ← ∅
5 while open ≠ ∅ do
6   sort(open)
7   n ← open.pop()
8   kids ← expand(n)
9   forall the kid ∈ kids do
10    kid.f ← (g + 1) + h(kid)
11    if goal(kid) = true then return makePath(kid)
12    if kid ∉ closed then open ← kid
13   closed ← n
14 return ∅
```

A*



$\phi_{explain}$:

Using explanation vocabulary

Ex: Obstacle presence

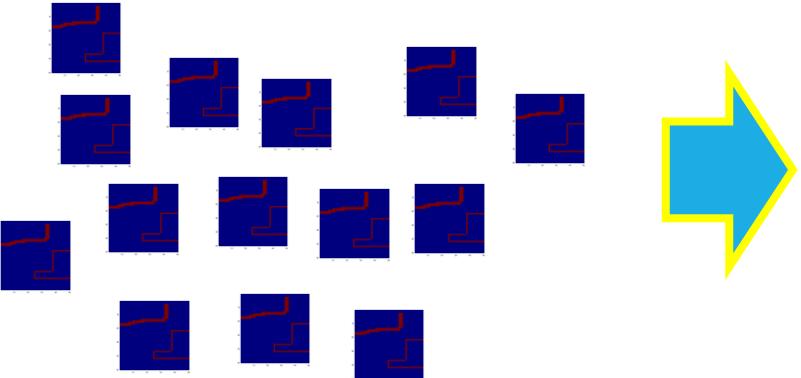
ϕ_{query} :

Some property of the output

Ex: Some cells not selected

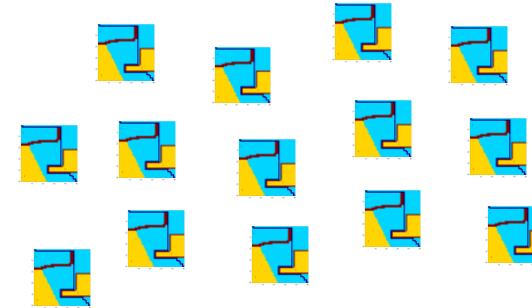
$$\begin{aligned}\phi_{explain} &\Rightarrow \phi_{query} \\ \phi_{explain} &\Leftrightarrow \phi_{query}\end{aligned}$$

Example: Explanations in A*



```
Algorithm 1: A*
Input: start, goal(n), h(n), expand(n)
Output: null
1 if goal(start) = true then return makePath(start)
2
3 open ← start
4 closed ← ∅
5 while open ≠ ∅ do
6   sort(open)
7   n ← open.pop()
8   kids ← expand(n)
9   forall the kid ∈ kids do
10    kid.f ← (n.g + 1) + h(kid)
11    if goal(kid) = true then return makePath(kid)
12    if kid ∉ closed then open ← kid
13   closed ← n
14 return ∅
```

A*



$\phi_{explain}$:

Using explanation vocabulary

Ex: Obstacle presence

ϕ_{query} :

Some property of the output

Ex: Some cells not selected

*Learn Decision Trees For $\phi_{explain}$
Using Labels for ϕ_{query}*

Example: Explanations in A*

50x50 grid has $2^{2^{50 \times 50}}$ possible explanations even if vocabulary only considers presence/absence of obstacles.

Scalability: Usually the feature space or vocabulary is large. For a map, its order of features in the map. For an image, it is order of the image's resolution.

Guarantee: Is the sampled space of maps enough to generate the explanation with some quantifiable probabilistic guarantee?

*Learn Decision Trees For $\phi_{explain}$
Using Labels for ϕ_{query}*

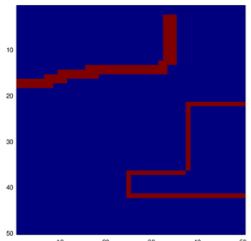
Example: Explanations in A*

Definition (*PAC learnability*) Let \mathcal{C} be a class of boolean functions $f : \{0, 1\}^n \rightarrow 0, 1$. We say that \mathcal{C} is PAC-learnable if there exists an algorithm \mathcal{L} such that

- for every $f \in \mathcal{C}$
- for any probability distribution \mathcal{D}
- for any ϵ (where $0 \leq \epsilon < \frac{1}{2}$)
- for any δ (where $0 \leq \delta < 1$)
- algorithm \mathcal{L} on input ϵ and δ and a set of random examples picked from any probability distribution \mathcal{D} outputs at least with a probability $1 - \delta$, concept h such that $\text{error}(h, f) \leq \epsilon$.

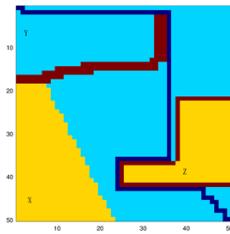
Theoretical Result: Learning Boolean formula even approximately is hard.
3-DNF is not learnable in Probably Approximately Correct framework unless RP = NP.

Two Key Ideas



```
Algorithm 1: A*
Input: start, goal(s), N(i), open(i)
Output: null
1: If goal(s) = true then return solvePath(start)
2: open ← start
3: closed ← ∅
4: while open ≠ ∅ do
5:   sort(open)
6:   n ← open.pop(0)
7:   N ← N(n)
8:   for each child in N do
9:     if child not in closed then
10:      f ← f(n, 1) + f(child)
11:      if f(child) <= f(best) or best = null then
12:        if f(child) <= f(best) then open ← {child}
13:        closed ← n
14: return ∅
```

A*



$\phi_{explain}$:

Using explanation vocabulary
Ex: Obstacle presence

- 1 Vocabulary is large.
iy samples (and what
on) to consider for
explanation ?
Boolean formula with
antees is hard.

ϕ_{query} :

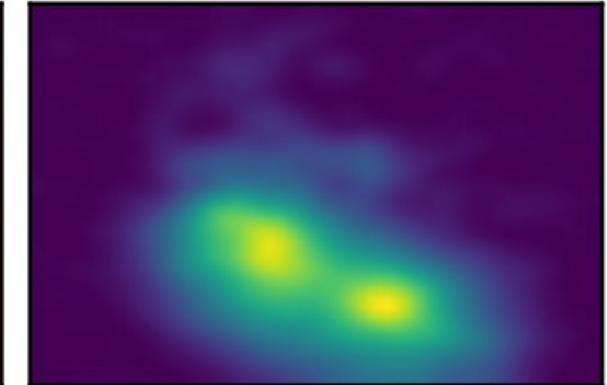
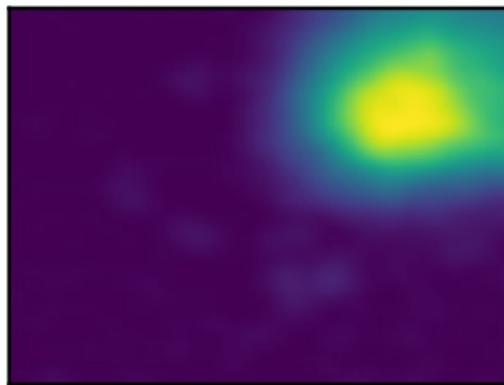
Some property of the output
Ex: Some cells not selected



Active learning Boolean formula $\phi_{explain}$ and not learning from fixed sample.

Explanations are often short and involve only few variables !

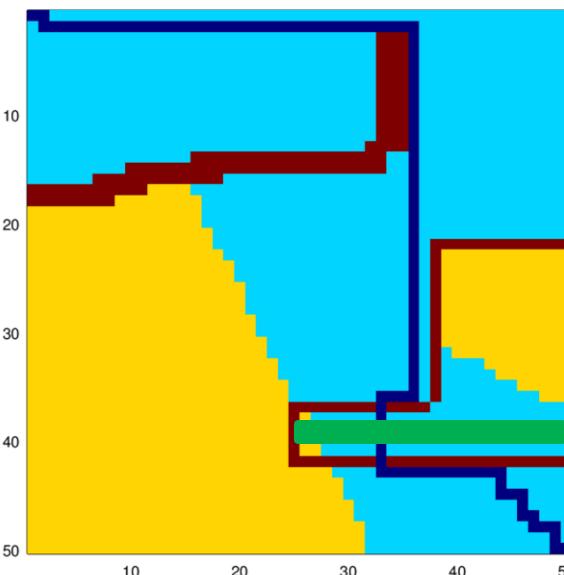
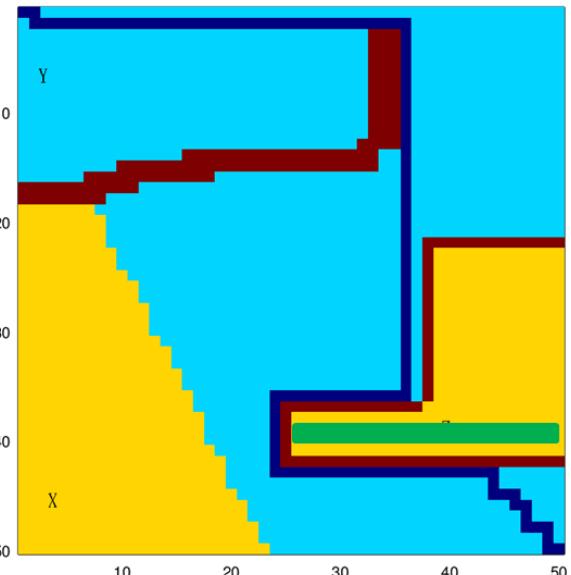
Two Key Ideas



Active learning Boolean formula $\phi_{explain}$ and not learning from fixed sample.

Explanations are often short and involve only few variables !

Two Key Ideas



Involves only two variables.
If we knew which two, we had
only $2^{2^2} = 16$
possible explanations.

How do we find these relevant
variables?

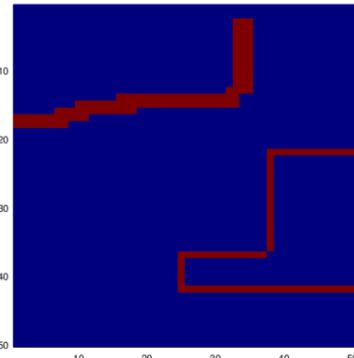


Active learning Boolean formula $\phi_{explain}$ and not learning from fixed sample.

Explanations are often short and involve only few variables !

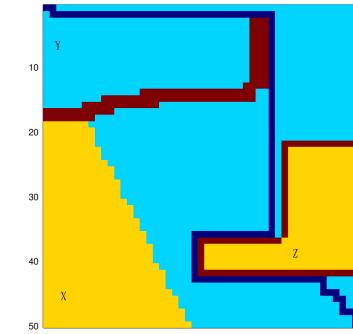
Actively Learning Boolean Formula

Oracle



```
Algorithm 1: A*
Input: start, goal(n), h(n), expand(n)
Output: null
1 if goal(start) = true then return makePath(start)
2
3 open ← start
4 closed ← ∅
5 while open ≠ ∅ do
6   sort(open)
7   n ← open.pop()
8   kids ← expand(n)
9   forall the kid ∈ kids do
10    kid.f ← (n.g + 1) + h(kid)
11    if goal(kid) = true then return makePath(kid)
12    if kid ∉ closed then open ← kid
13 closed ← n
14 return ∅
```

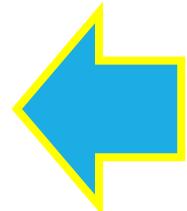
A*



ϕ_{query} :
Some property of the output
Ex: Some cells not selected



Assignments to V
 $m_1 = (0,0,0,1,1,0,1)$
 $m_2 = (0,0,1,1,0,1,0)$



$\phi_{explain}(V)$:
Using explanation vocabulary
Ex: Obstacle presence

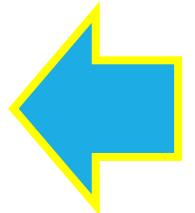
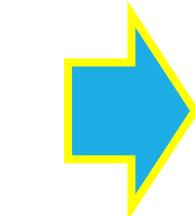
ϕ
Evaluates assignments and returns T,F

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Assignments to V

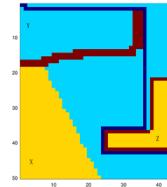
$m1 = (0,0,0,1,1,0,1)$



```
Algorithm 1: A*
Input: start, goal, N, C, openList
Output: sol
1 if openList = {} then return calculateSol(start)
2
3 openList ← {start}
4 closedList ← {}
5 while openList ≠ {} do
6   current ← openList.pop(0)
7   if current = goal then
8     return calculateSol(current)
9   for each child in children do
10     if f - f(g, h) >= 44.05
11     if 2 * g(f) - f(g) >= 0 then return calculateSol()
12     if f(g) - f(g) = 0 then openList.append(child)
13   closedList.append(current)
14 return None
```



A*



ϕ_{query} :
Some property of the output
Ex: Some cells not selected

Oracle

$m1$: True

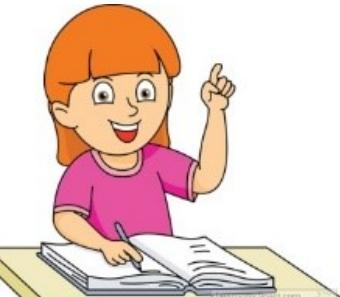
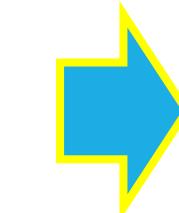
Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

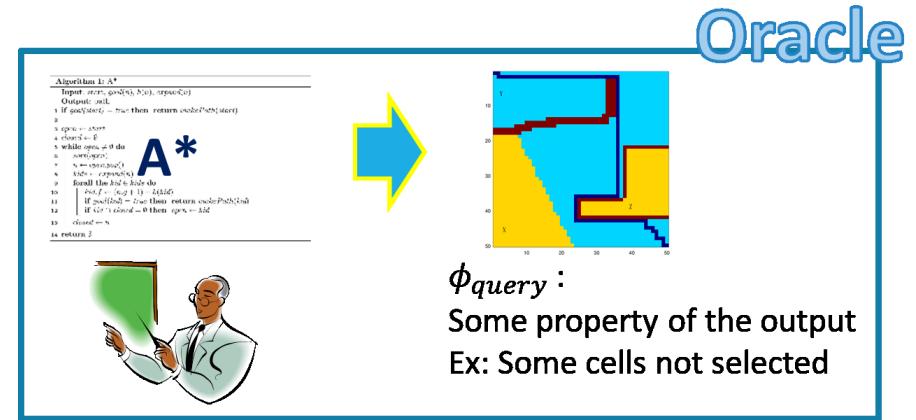
Assignments to V

$m_1 = (0,0,0,1,1,0,1)$

$m_2 = (0,0,1,1,0,1,0)$



Random Sample Till
Oracle differs



m_1 : True, m_2 : False

Actively Learning Relevant Variables

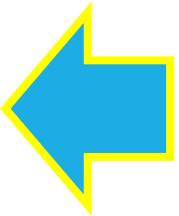
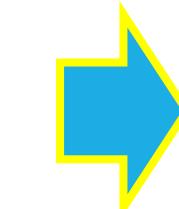
Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Assignments to V

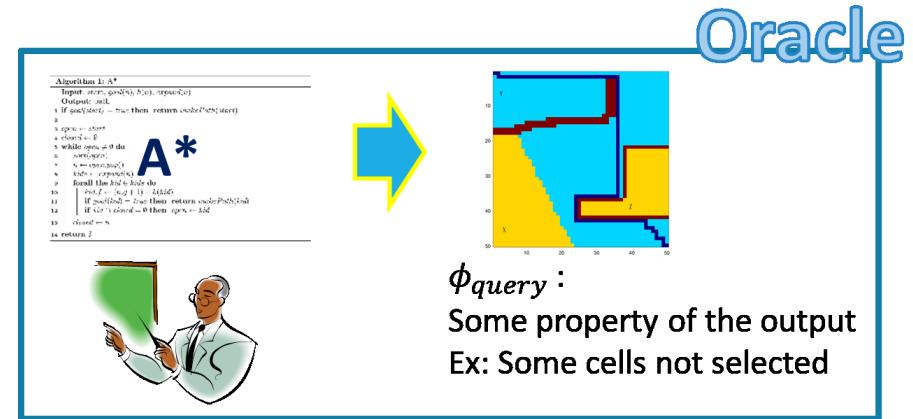
$m1 = (0,0,0,1,1,0,1)$

$m2 = (0,0,1,1,0,1,0)$

$m3 = (0,0,0,1,1,1,0)$



Binary Search Over
Hamming Distance



$m1: \text{True}, m2: \text{False}$

Actively Learning Relevant Variables

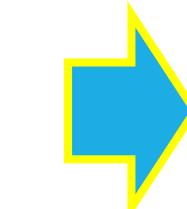
Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Assignments to V

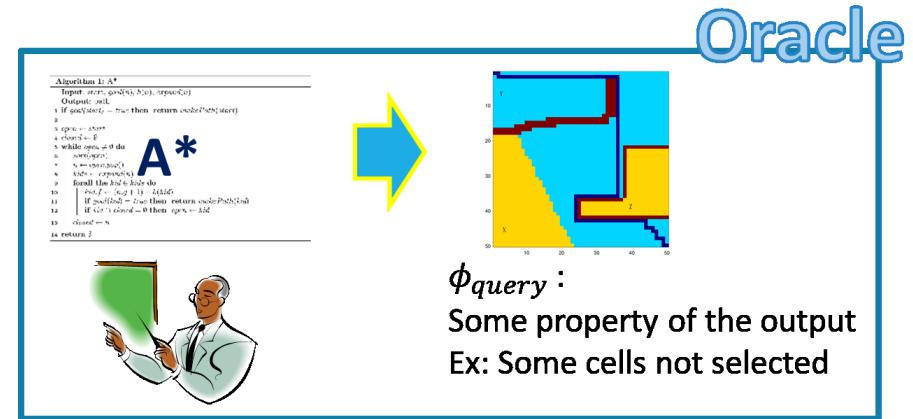
$m1 = (0,0,0,1,1,0,1)$

$m2 = (0,0,1,1,0,1,0)$

$m3 = (0,0,0,1,1,1,0)$



Binary Search Over
Hamming Distance



$m1$: True, $m2$: False
 $m3$: True

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Hamming
Distance = 4

Assignments to V

~~$m_1 = (0, 0, 0, 1, 1, 0, 1)$~~

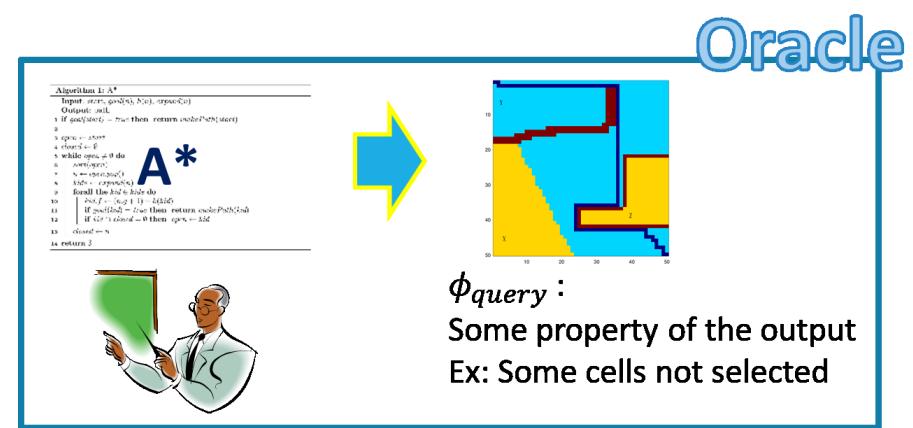
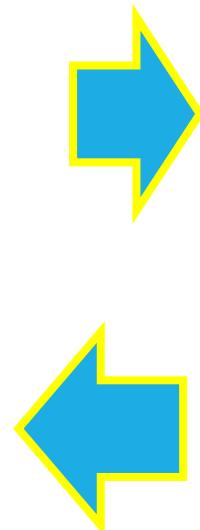
$m_2 = (0, 0, 1, 1, 0, 1, 0)$

$m_3 = (0, 0, 0, 1, 1, 1, 0)$

Hamming
Distance = 2



Binary Search Over
Hamming Distance



~~m_1 : True~~, m_2 : False
 m_3 : True

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Hamming
Distance = 2

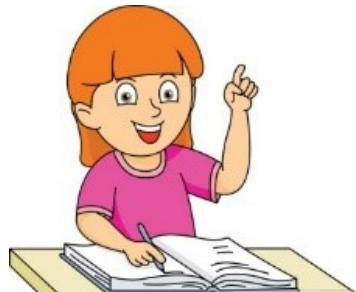
Assignments to V

$m_2 = (0,0,1,1,0,1,0)$

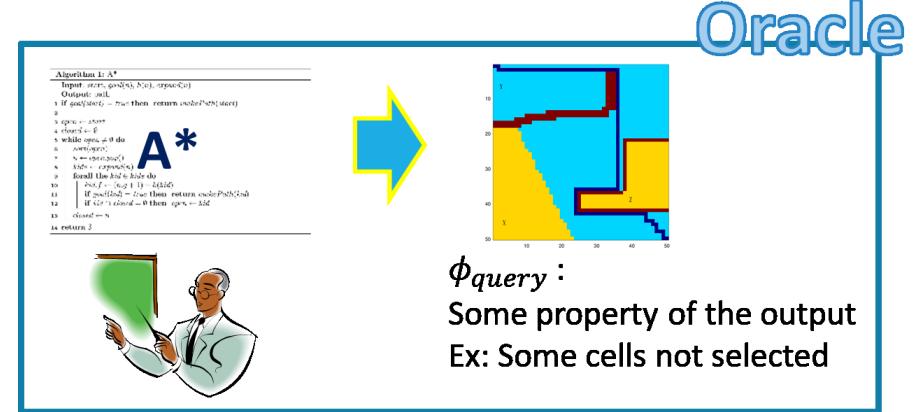
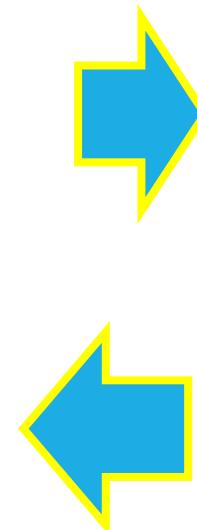
$m_3 = (0,0,0,1,1,1,0)$

$m_4 = (0,0,1,1,1,1,0)$

Hamming
Distance = 1



Binary Search Over
Hamming Distance



m_2 : False, m_3 : True
 m_4 : True

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Hamming
Distance = 2

Assignments to V

$m_2 = (0,0,1,1,0,1,0)$

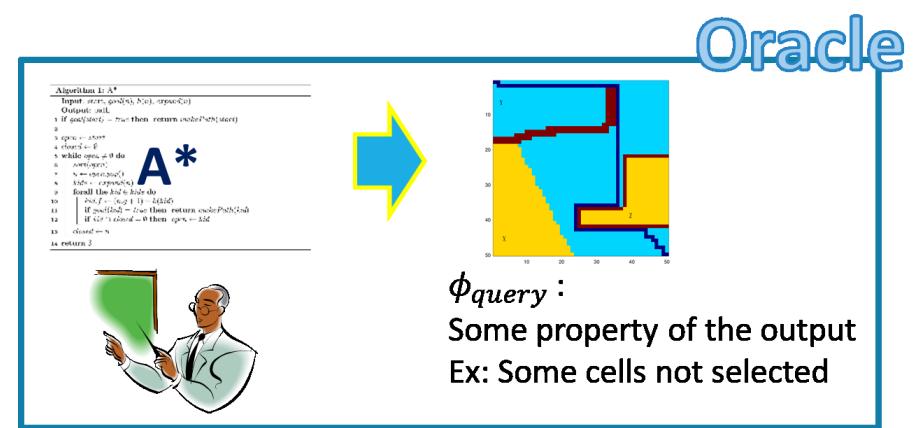
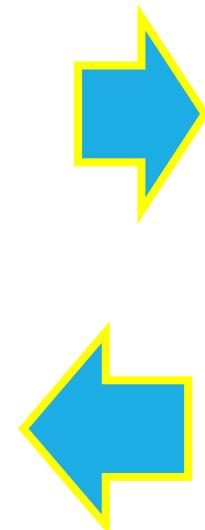
~~$m_3 = (0,0,0,1,1,1,0)$~~

$m_4 = (0,0,1,1,1,1,0)$

Hamming
Distance = 1



Binary Search Over
Hamming Distance



m_2 : False, ~~m_3 : True~~
 m_4 : True

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Hamming
Distance = 1

Assignments to V

$m2 = (0,0,1,1,0,1,0)$
 $m4 = (0,0,1,1,1,1,0)$



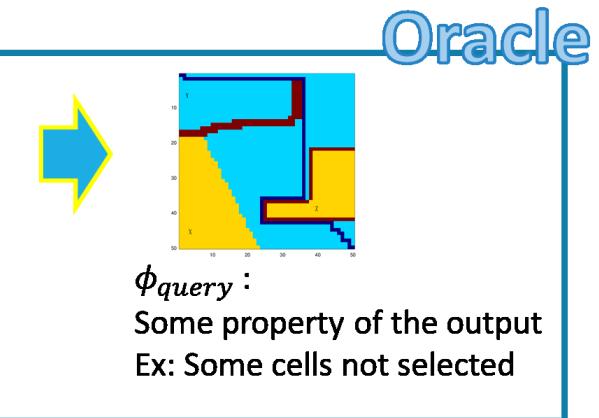
Fifth variable v_5 is relevant !!



Binary Search Over
Hamming Distance

Algorithm 1: A*

```
Input: startNode, N, goalNode, openSet, closedSet
1 If startNode == goalNode then return closedPath(start)
2 If openSet == {} then return closedPath(goal)
3 openSet += {start}
4 closedSet += {}
5 while openSet != {} do
6   current = openSet.pop(0)
7   if current == goalNode then
8     return closedPath(goal)
9   for all neighbors in neighbors(current):
10     if fValue(neighbors) < fValue(openSet[0]):
11       openSet.remove(openSet[0])
12       openSet += [neighbors]
13       closedSet += {current}
14 return goal
```



$m2$: False, $m4$: True

Actively Learning Relevant Variables

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

For each assignment
to relevant variables

$$2^{|U|}$$



Random Sample Till
Oracle differs

$$\ln(1/(1 - \kappa))$$



Binary Search Over
Hamming Distance

$$\ln(|V|)$$

Relevant variables of $\phi_{explain}$ found with confidence κ in

$$2^{|U|} \ln\left(\frac{|V|}{1 - \kappa}\right)$$

Actively Learning Boolean Formula

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Used distinguishing input based approach from ICSE'10



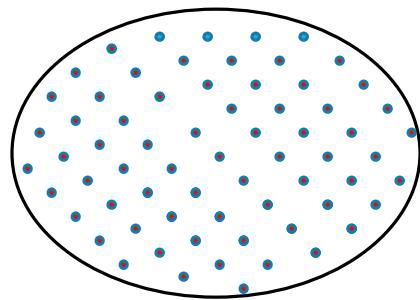
Build Truth Table for
the relevant variables U

Worst Case: $2^{|U|}$

$\phi_{explain}$ found with confidence κ in

$$2^{|U|} \left(1 + \ln\left(\frac{|V|}{1 - \kappa}\right)\right)$$

Distinguishing Example

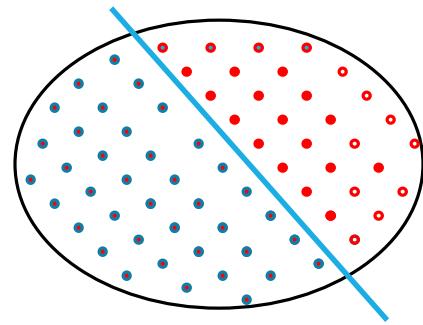


Space of all possible $2^{|U|}$ formulae.
Each dot represents semantically
unique Boolean formula

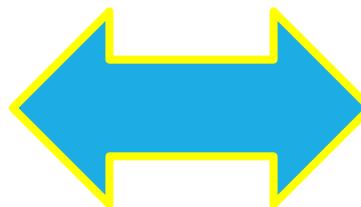
Distinguishing Example



Oracle O: Label e_1 as true iff it satisfies target ϕ



some of the $2^{2^{|U|}}$ formulae are not consistent with example set $\{e_1\}$



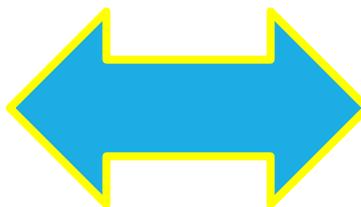
Randomly select - assignment to d Boolean variables - e_1



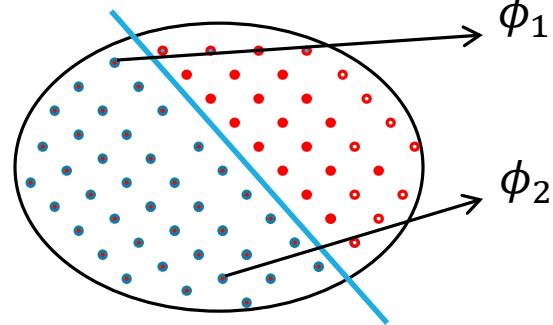
Distinguishing Example



Oracle O: Label e_1 as true iff it satisfies target ϕ



Randomly select - assignment to d Boolean variables - e_1



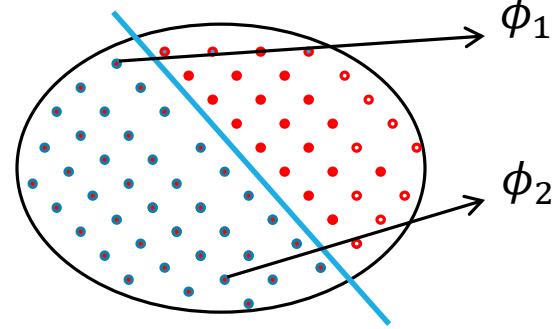
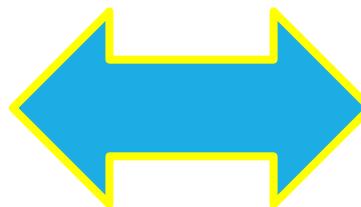
2 Candidate Formulae
Consistent with Example Set

some of the $2^{|U|}$ formulae are not consistent with example set $\{e_1\}$

Distinguishing Example



Oracle O: Label e_2 as true iff it satisfies target ϕ



some of the $2^{2^{|U|}}$ formulae are not consistent with example set
 $\{e_1\}$

Randomly select - assignment to d Boolean variables - e_1

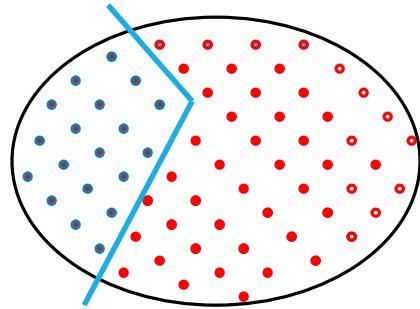


Find distinguishing example
 $e_2 \vDash \phi_1 \oplus \phi_2$

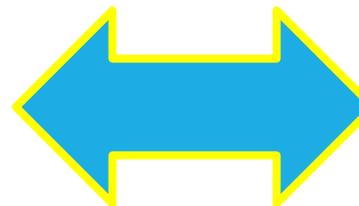
Distinguishing Example



Oracle O: Label e_2 as true iff it satisfies target ϕ



some of the $2^{|U|}$ formulae are not consistent with example set
 $\{e_1, e_2\}$



Randomly select - assignment to d Boolean variables - e_1

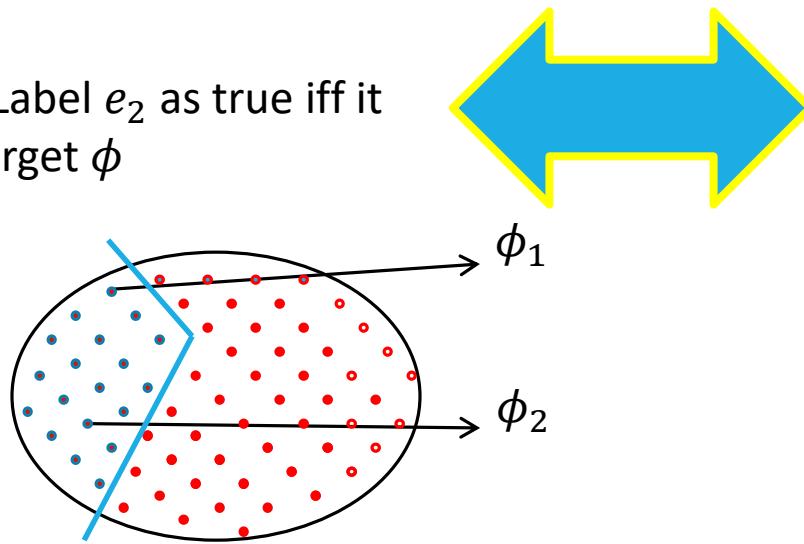


Find distinguishing example
 $e_2 \models \phi_1 \oplus \phi_2$

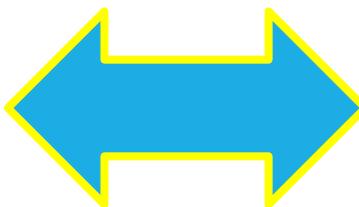
Distinguishing Example



Oracle O: Label e_2 as true iff it satisfies target ϕ



some of the $2^{2^{|U|}}$ formulae are not consistent with example set
 $\{e_1, e_2\}$



Randomly select - assignment to d Boolean variables - e_1

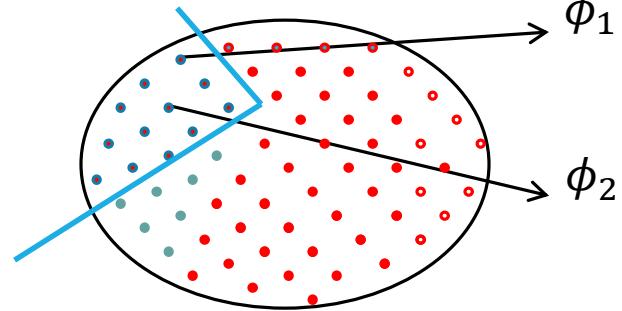
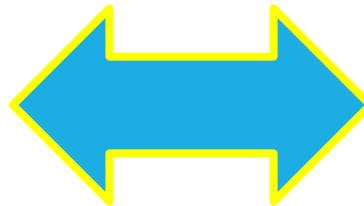


Find distinguishing example
 $e_3 \models \phi_1 \oplus \phi_2$

Distinguishing Example



Oracle O: Label e_2 as true iff it satisfies target ϕ



some of the $2^{|U|}$ formulae are not consistent with example set

$$\{e_1, e_2, e_3\}$$

Randomly select - assignment to d Boolean variables - e_1

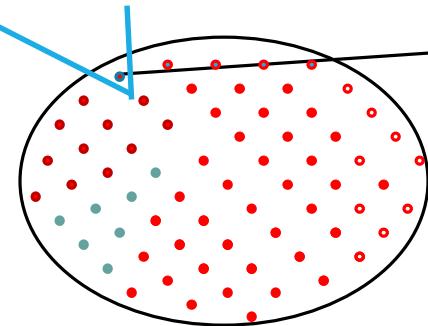


Find distinguishing example
 $e_3 \models \phi_1 \oplus \phi_2$

Distinguishing Example

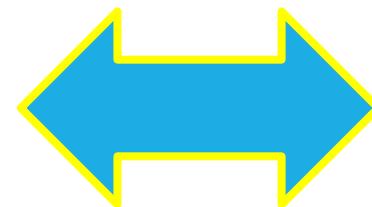


Oracle O: Label e_2 as true iff it satisfies target ϕ



some of the $2^{|U|}$ formulae are not consistent with example set

$$\{e_1, e_2, e_3, \dots, e_k\}$$



Randomly select - assignment to d Boolean variables - e_1



Find distinguishing example
 $e_k \vDash \phi_1 \oplus \phi_2$

Terminate when we are left with a single consistent ϕ w.r.t all the examples.

Actively Learning Boolean Formula

Find U such that $\phi_{explain}(V) \equiv \phi_{explain}(U)$ where $|U| \ll |V|$

Correctness Guarantee: If there exists an explanation using variables in V , then this approach is guaranteed to find it (learn correct Boolean formula).

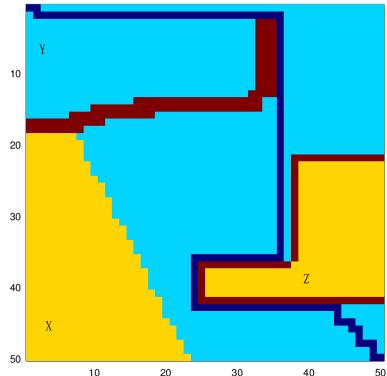
If there does not exist an explanation using V , the approach will either detect it or find a wrong explanation (learn incorrect Boolean formula).

In practice, since we depend only logarithmically on $|V|$, we use large vocabulary.

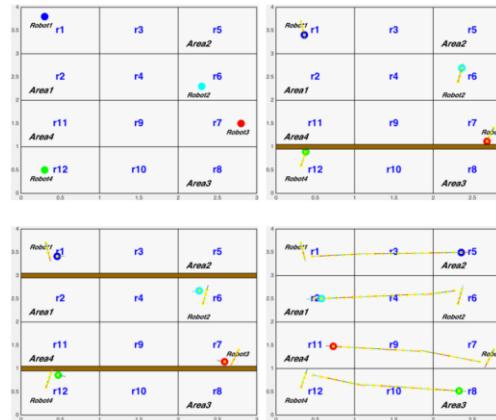


$$2^{|U|} (1 + \ln\left(\frac{|V|}{1 - \kappa}\right))$$

Experiments



Explaining A* Planning
 $|V| = 2500$ 10^153
 $|U| \leq 4$
Runtime < 3 minutes



Reactive Exploration Strategy
 $|V| = 96$ 10^28
 $|U| \leq 2$
Runtime < 5 seconds

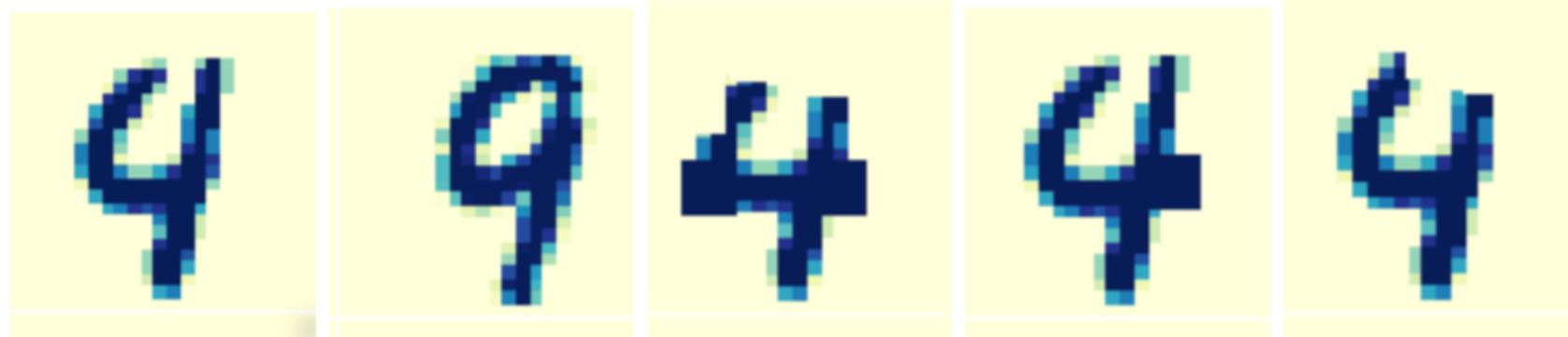


Image Classification
 $|V| = 784$
 $|U| \leq 27$
Runtime < 5 minutes

Challenges

Improving selection of inputs to run to reduce number of introspections needed.

- Exploit structural assumption on nature of explanation
- Exploit probability distribution over inputs

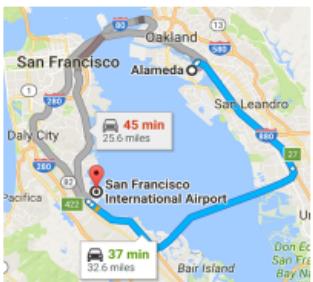
Exploitation of partial knowledge about AI algorithms to enable efficient sampling.

Extension to randomized decision making algorithms which may make different decisions for the same input.

Automated/Iterative Selection of Vocabulary

Thanks!!

Ubiquitous AI and need for explaination



Why did we take the San Mateo bridge instead of the Bay Bridge ?

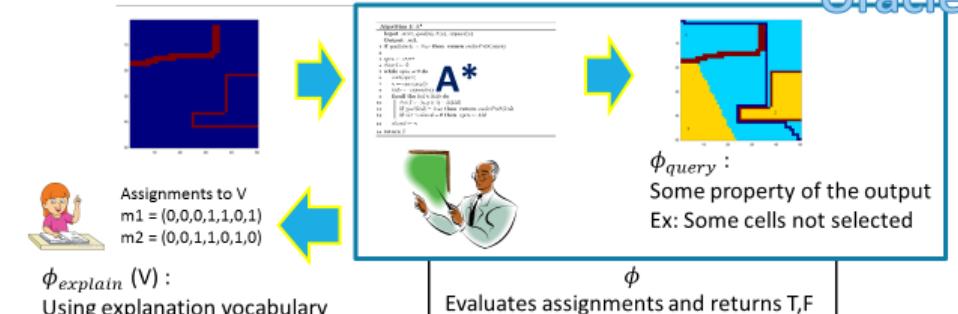
5/15/2017

2

$\phi_{explain}$ found with confidence κ in
 $O(2^{|U|} \ln(|V|/(1 - \kappa)))$

Actively Learning Boolean Formula

Oracle



Thanks!

It is really confusing!!!

