

Obedient machine learning

Machine learning models that provably obey constraints

Sebastijan Dumancic

joint work with Kshitij Goyal, Hendrik Blockeel

Is AI trustworthy?

New AI can guess whether you're gay or straight from a photograph

An algorithm deduced the sexuality of people on a dating site with up to 91% accuracy, raising tricky ethical questions

FROM POLITICO PRO

Dutch scandal serves as a warning for Europe over risks of using algorithms

The Dutch tax authority ruined thousands of lives after using an algorithm to spot suspected benefits fraud – and critics say there is little stopping it from happening again.

Uber's self-driving operator charged over fatal crash

Bernard Parker, left, was rated high risk; Dylan

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

ARTIFICIAL INTELLIGENCE | OPINION

AI-Influenced Weapons Need Better Regulation

The weapons are error-prone and could hit the wrong targets

Towards trustworthy systems

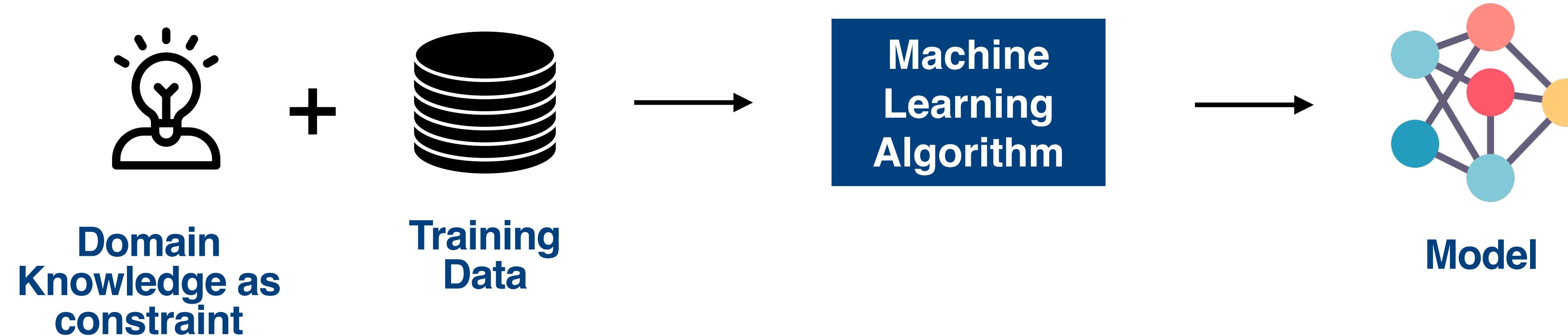
The problem lies in our approach: not everything is in training data...



Towards trustworthy systems

The problem lies in our approach: not everything is in training data...

... many desires for trustworthiness have to be incorporated in some other ways

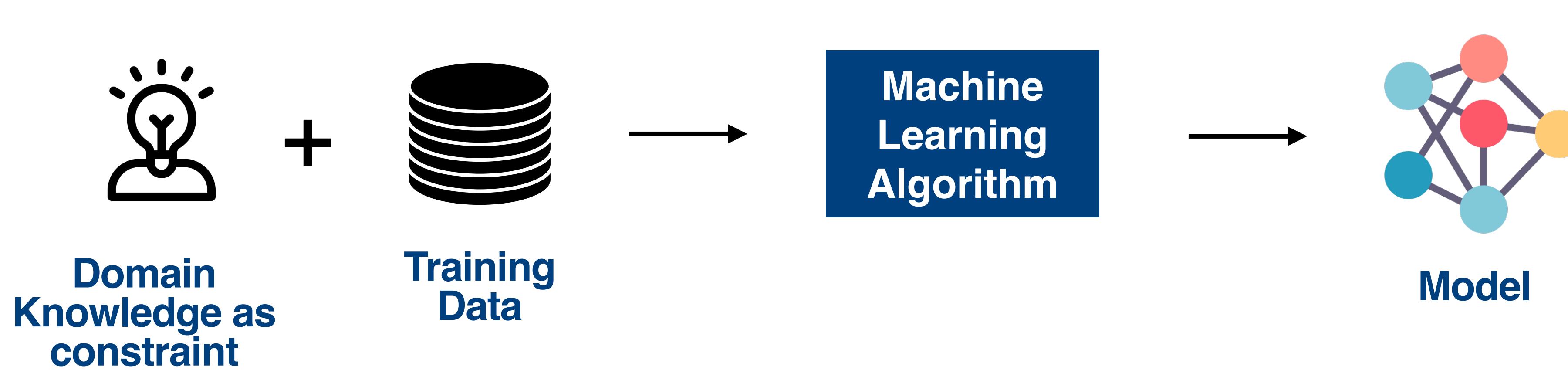


Towards trustworthy systems

The problem lies in our approach: not everything is in training data...

... many desires for trustworthiness have to be incorporated in some other ways

The law requires a certification that the system satisfies constraints



Desiderata

Learning machine learning models that

- guarantee constraint satisfaction
- of a wide range of constraints
- on all possible instances

ML and constraints: A note on constraints

Syntactic Constraints

Constraints on the
model structure

Semantic Constraints

Constraints on the behaviour of
the model predictions

A note on constraints

Syntactic Constraints

Constraints on the
model structure

- **Box Constraints:** upper and lower bound on one of more parameters of the model
- **Depth in decision trees:** to learn smaller trees which are more interpretable

Semantic Constraints

Constraints on the behaviour of
the model predictions

- **Safety Constraint:** If there's a stationary object right in front of a moving autonomous car, the car must make a hard stop in all such scenarios.
- **Multi Target Predictions with Constrained Outputs:** The speech to text translation must follow certain linguistic constraints for all predictions.

What have we tried so far

Loss + $\lambda \times$ penalty

Semantic loss, DL2,
Semantic-Based Regularisation

Constraints are soft and
not guaranteed

Very difficult to encode anything
complex-ish

Specific architecture

MultiplexNet

Constraints are guaranteed

Very difficult to support wide constraints

Semantic layers

If only we had....

Techniques that solve constraints

Constraint satisfaction problems



Constraint satisfaction problems

Decision variables =
which colour is a state?

```
% Colouring Australia using nc colours
int: nc = 3;

var 1..nc: wa;  var 1..nc: nt;  var 1..nc: sa;  var 1..nc: q;
var 1..nc: nsw; var 1..nc: v;   var 1..nc: t;
```



Constraint satisfaction problems

Decision variables =
which colour is a state?

```
% Colouring Australia using nc colours
int: nc = 3;

var 1..nc: wa;  var 1..nc: nt;  var 1..nc: sa;  var 1..nc: q;
var 1..nc: nsw; var 1..nc: v;   var 1..nc: t;

constraint wa != nt;
constraint wa != sa;
constraint nt != sa;
constraint nt != q;
constraint sa != q;
constraint sa != nsw;
constraint sa != v;
constraint q != nsw;
constraint nsw != v;
```

Constraints =
neighbours should have different colours



Constraint satisfaction problems

Decision variables =
which colour is a state?

```
% Colouring Australia using nc colours
int: nc = 3;

var 1..nc: wa;  var 1..nc: nt;  var 1..nc: sa;  var 1..nc: q;
var 1..nc: nsw; var 1..nc: v;   var 1..nc: t;

constraint wa != nt;
constraint wa != sa;
constraint nt != sa;
constraint nt != q;
constraint sa != q;
constraint sa != nsw;
constraint sa != v;
constraint q != nsw;
constraint nsw != v;
solve satisfy;
```

Constraints =
neighbours should have different colours



Constraint satisfaction

**Boolean
Satisfiability
(SAT)**

$$\neg a \vee (b \wedge c)$$

$a = False$

**Satisfiability
Modulo
Theories**

$$(a > 3) \wedge (a < 3)$$

unsat

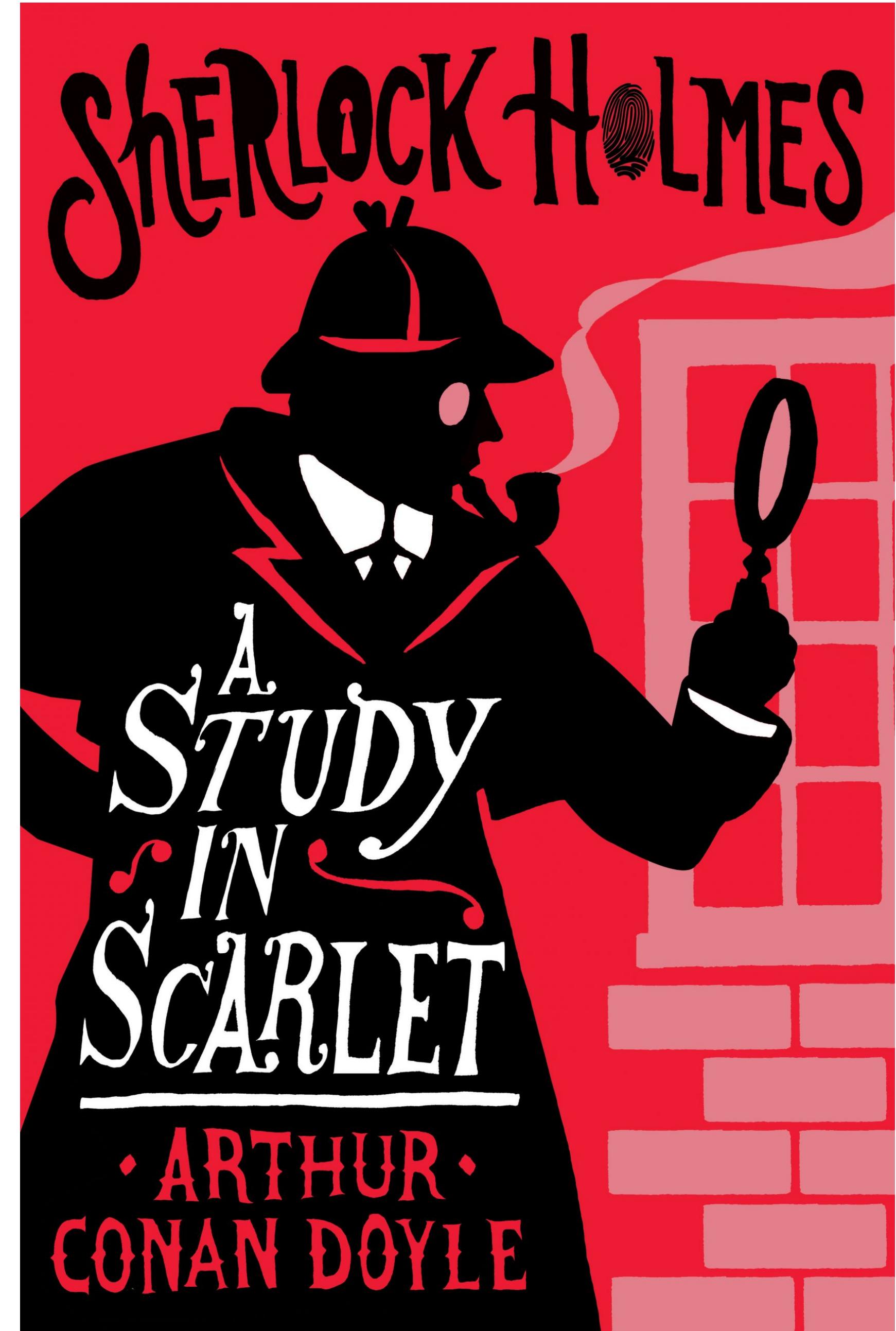
**Maximum
Satisfiability
(MaxSAT/MaxSMT)**

$$\phi_h : \neg a \wedge (b \vee c) \quad \phi_s : a \wedge b$$

$a = False, b = True$

A study in...

making logic do something it
was not supposed to do



SaDe:
ML+Constraints with satisfiability descent

Desiderata revisited

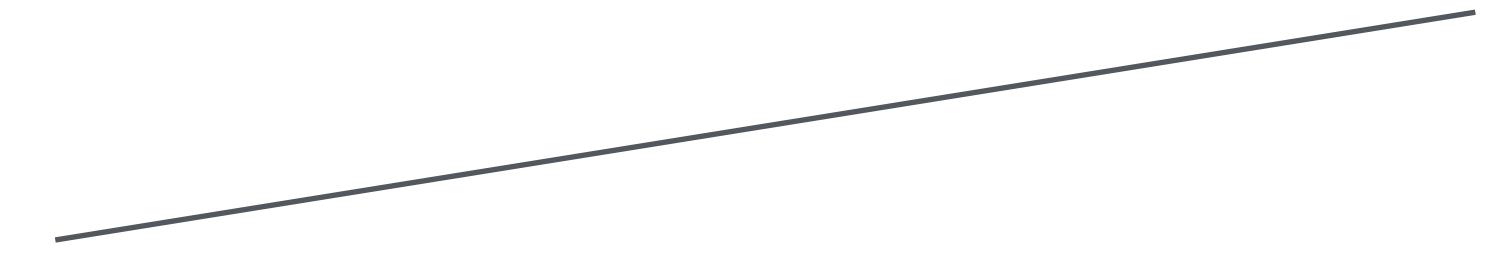
Given

- a dataset D
- (domain) constraints C
- domain of instances I

Desiderata revisited

Given

- a dataset D
- (domain) constraints C
- domain of instances I

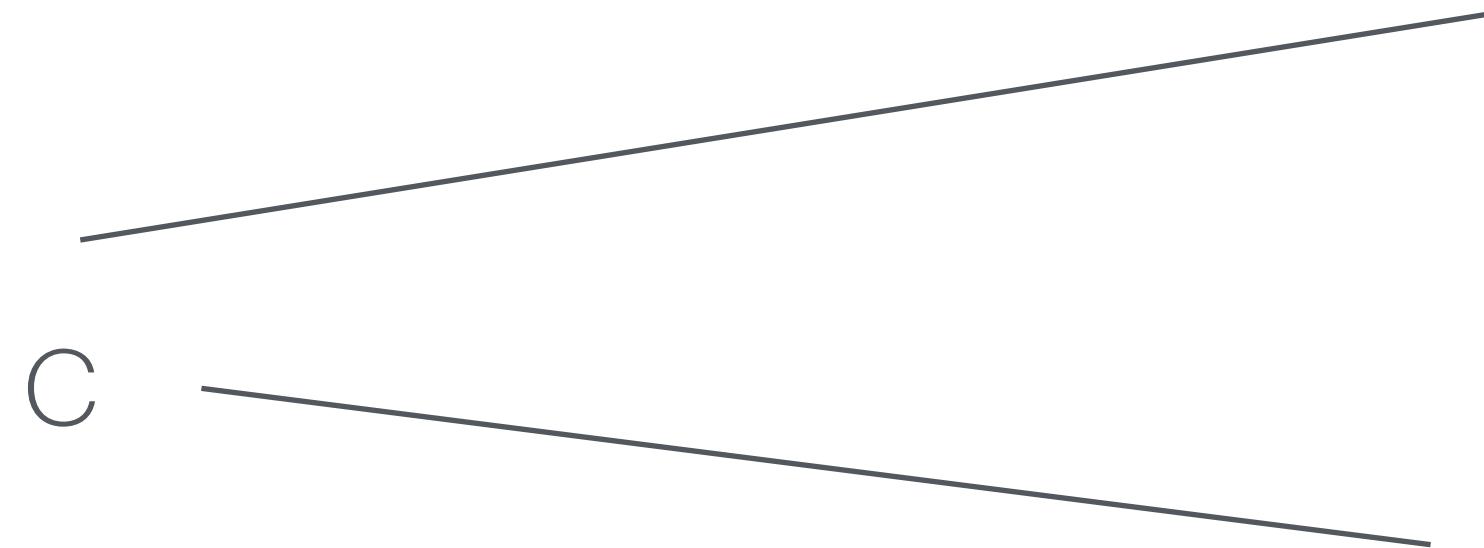


x₁	x₂	y
1	4	0
2	3	0
-3	2	1

Desiderata revisited

Given

- a dataset D
- (domain) constraints C
- domain of instances I



x₁	x₂	y
1	4	0
2	3	0
-3	2	1

$$\forall x_1, x_2 \ f(x_1, x_2) > 0$$

Desiderata revisited

Given

- a dataset D
- (domain) constraints C
- domain of instances I

x₁	x₂	y
1	4	0
2	3	0
-3	2	1

$$\forall x_1, x_2 \quad f(x_1, x_2) > 0$$

$$x_1 \in [-4, 2] \quad x_2 \in [0, 6]$$

Desiderata revisited

Given

- a dataset D
- (domain) constraints C
- domain of instances I

Learn a model that

- satisfies constraints C
on entire instance space I
- maximises predictive performance

A closer look into the learning problem

Learn a model that

- satisfies constraints C
on entire instance space I
- maximises predictive performance

A closer look into the learning problem

Learn a model that

- satisfies constraints C on entire instance space I
- maximises predictive performance

Point 1

Impossible to ensure with gradient descent and only from data

A closer look into the learning problem

Learn a model that

- satisfies constraints C on entire instance space I
- maximises predictive performance

Point 1

Impossible to ensure with gradient descent and only from data

Point 2

Domain constraints need to be satisfied

Machine learning as MaxSMT

x₁	x₂	y
1	4	0
2	3	0
-3	2	1

$$\forall x_1, x_2 \quad f(x_1, x_2) > 0$$

$$x_1 \in [-4, 2] \quad x_2 \in [0, 6]$$

Machine learning as MaxSMT

x₁	x₂	y
1	4	0
2	3	0
-3	2	1

Translate model decisions
into soft constraints



Classification

$$\begin{aligned}f_{\mathbf{w}}([1,4]) &< 0 \\f_{\mathbf{w}}([2,3]) &< 0 \\f_{\mathbf{w}}([-3,2]) &> 0\end{aligned}$$

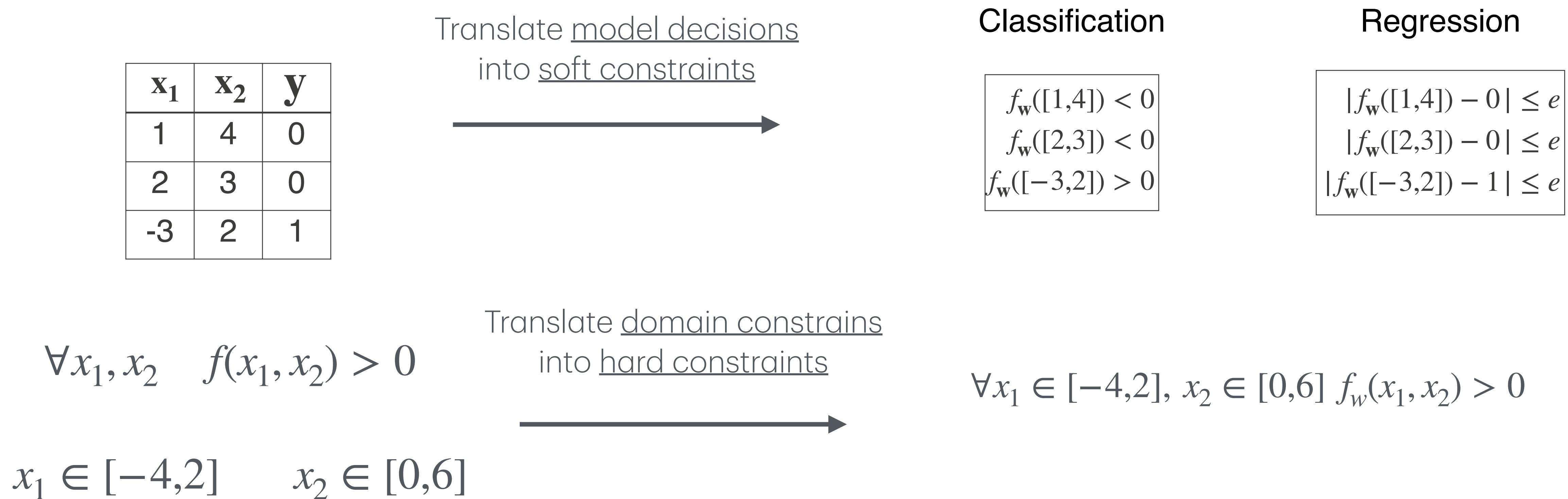
Regression

$$\begin{aligned}|f_{\mathbf{w}}([1,4]) - 0| &\leq e \\|f_{\mathbf{w}}([2,3]) - 0| &\leq e \\|f_{\mathbf{w}}([-3,2]) - 1| &\leq e\end{aligned}$$

$$\forall x_1, x_2 \quad f(x_1, x_2) > 0$$

$$x_1 \in [-4,2] \quad x_2 \in [0,6]$$

Machine learning as MaxSMT



Machine learning as MaxSMT

Assuming a linear model $f(x_1, x_2) = w_1x_1 + w_2x_2 + w_3$

Machine learning as MaxSMT

Assuming a linear model $f(x_1, x_2) = w_1x_1 + w_2x_2 + w_3$

Decision variables = model parameters

`var w_1 : float`

`var w_2 : float`

`var w_3 : float`

Machine learning as MaxSMT

Assuming a linear model $f(x_1, x_2) = w_1x_1 + w_2x_2 + w_3$

Decision variables = model parameters

`var w_1 : float`

`var w_2 : float`

`var w_3 : float`

Hard constraint = domain constraint

$\forall x_1 \in [-4, 2], x_2 \in [0, 6] \ w_1x_2 + w_2x_2 + w_3 > 0$

Machine learning as MaxSMT

Assuming a linear model $f(x_1, x_2) = w_1x_1 + w_2x_2 + w_3$

Decision variables = model parameters

```
var  $w_1$  : float  
var  $w_2$  : float  
var  $w_3$  : float
```

Hard constraint = domain constraint

$$\forall x_1 \in [-4, 2], x_2 \in [0, 6] \ w_1x_2 + w_2x_2 + w_3 > 0$$

Soft constraints = model predictions

$$4w_1 + 1w_2 + w_3 < 0$$
$$2w_1 + 3w_2 + w_3 < 0$$
$$-3w_1 + 2w_2 + w_3 > 0$$

Representation of domain constraints

Goal: Enforce constraint(s) on the predictions
for all possible instances in the input domain

Representation of domain constraints

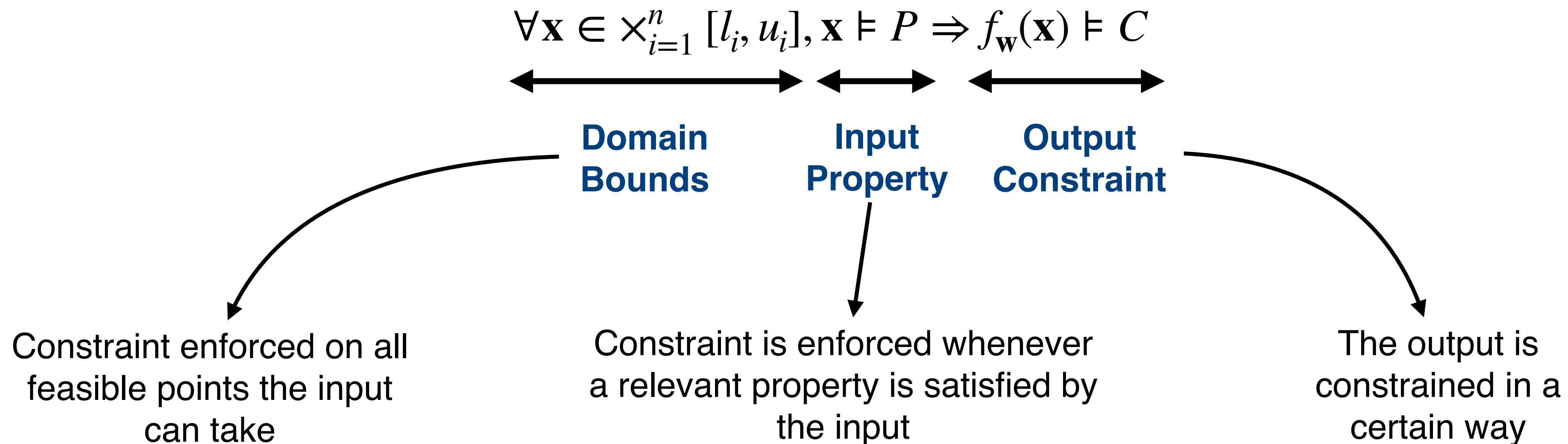
Goal: Enforce constraint(s) on the predictions
for all possible instances in the input domain

→ **Universal Quantifiers (\forall)**

Representation of domain constraints

Goal: Enforce constraint(s) on the predictions for all possible instances in the input domain

→ Universal Quantifiers (\forall)



Representation of domain constraints

Goal: Enforce constraint(s) on the predictions for all possible instances in the input domain

→ **Universal Quantifiers (\forall)**

Solving for such a constraint with constraint solvers potentially allows to find a model ($f_w(\cdot)$) that satisfies the constraint for all predictions

$$\forall \mathbf{x} \in \times_{i=1}^n [l_i, u_i], \mathbf{x} \models P \Rightarrow f_w(\mathbf{x}) \models C$$

Domain
Bounds

Input
Property

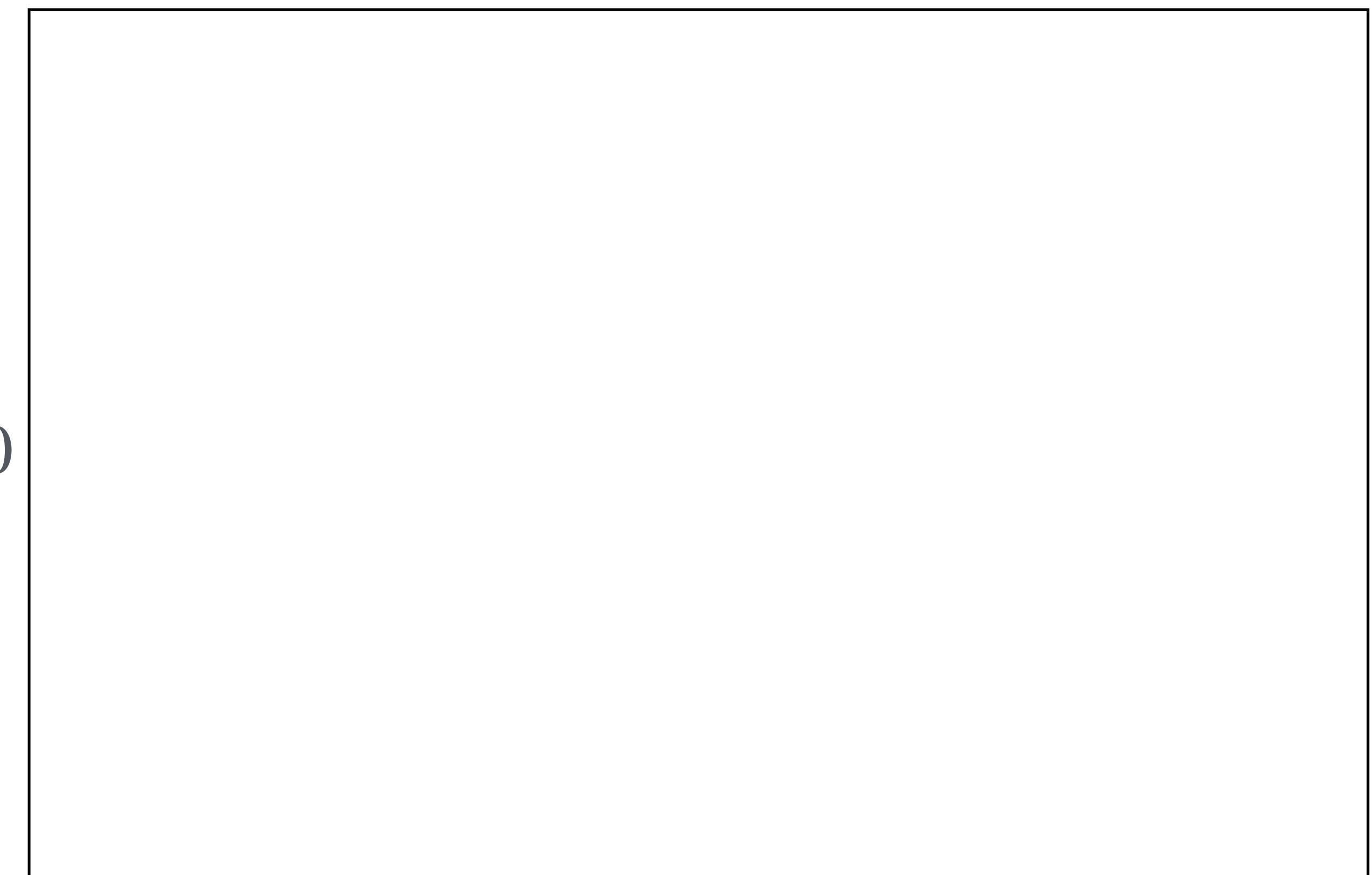
Output
Constraint

Constraint enforced on all feasible points the input can take

Constraint is enforced whenever a relevant property is satisfied by the input

The output is constrained in a certain way

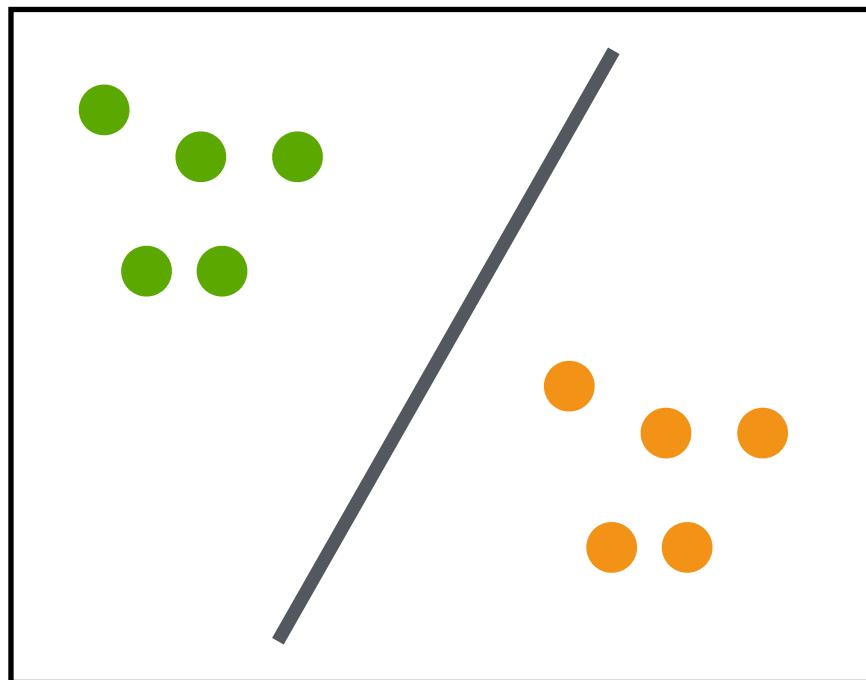
What are we doing here?



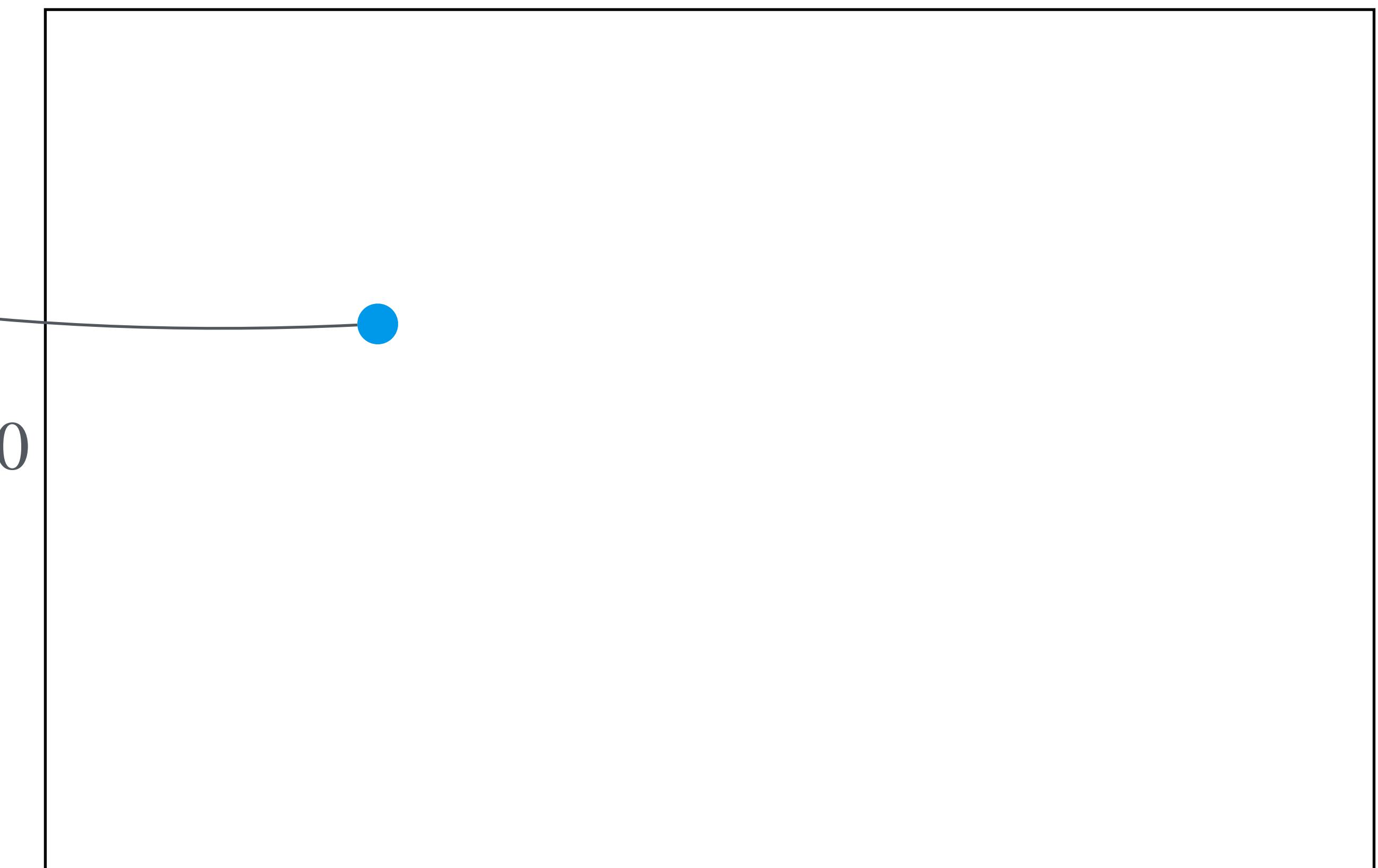
$$w_2 = 0$$

$$w_1 = 0$$

What are we doing here?

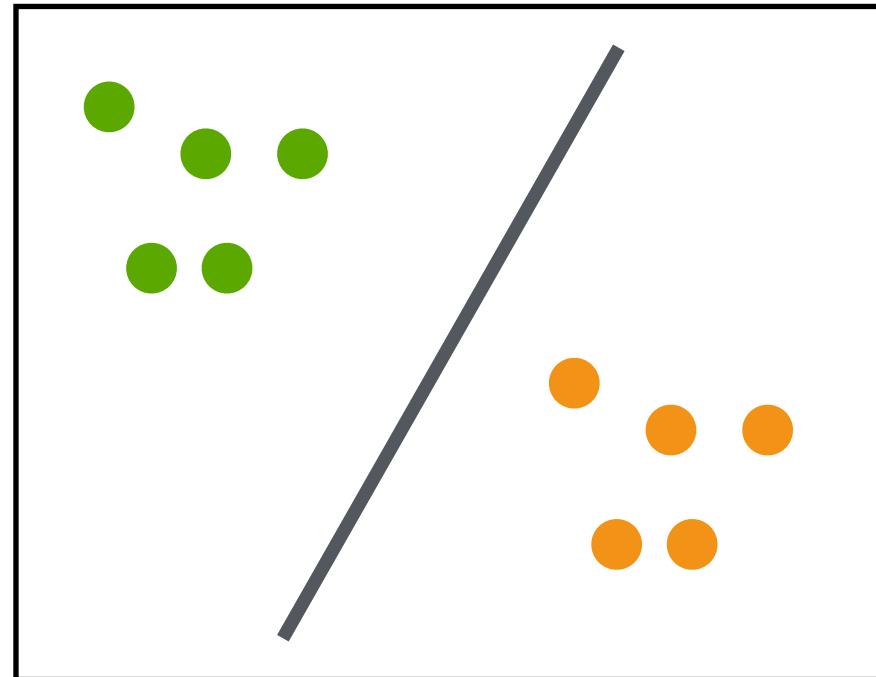


A point in hypothesis
space is a model

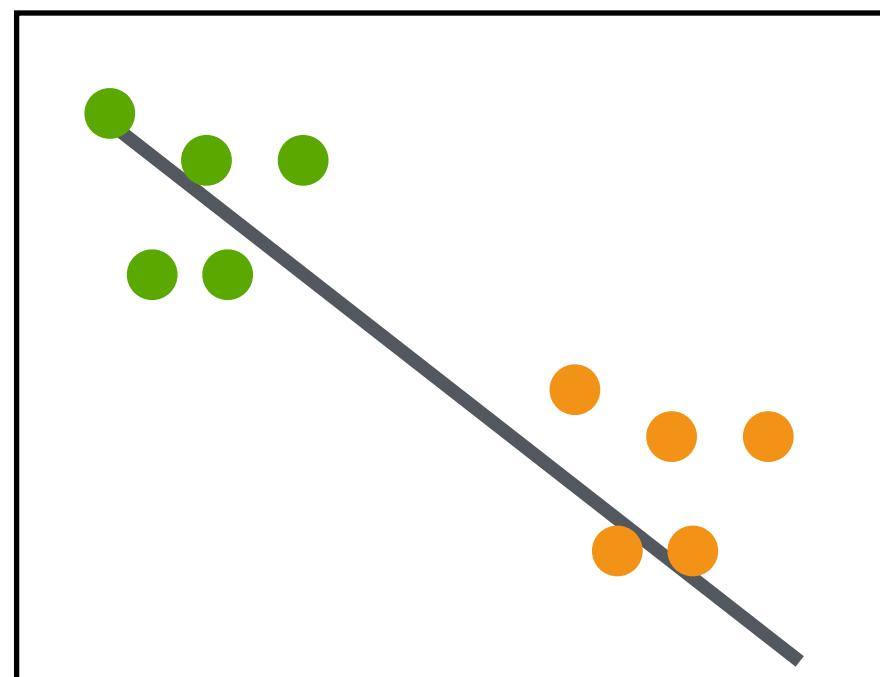


$$w_1 = 0$$

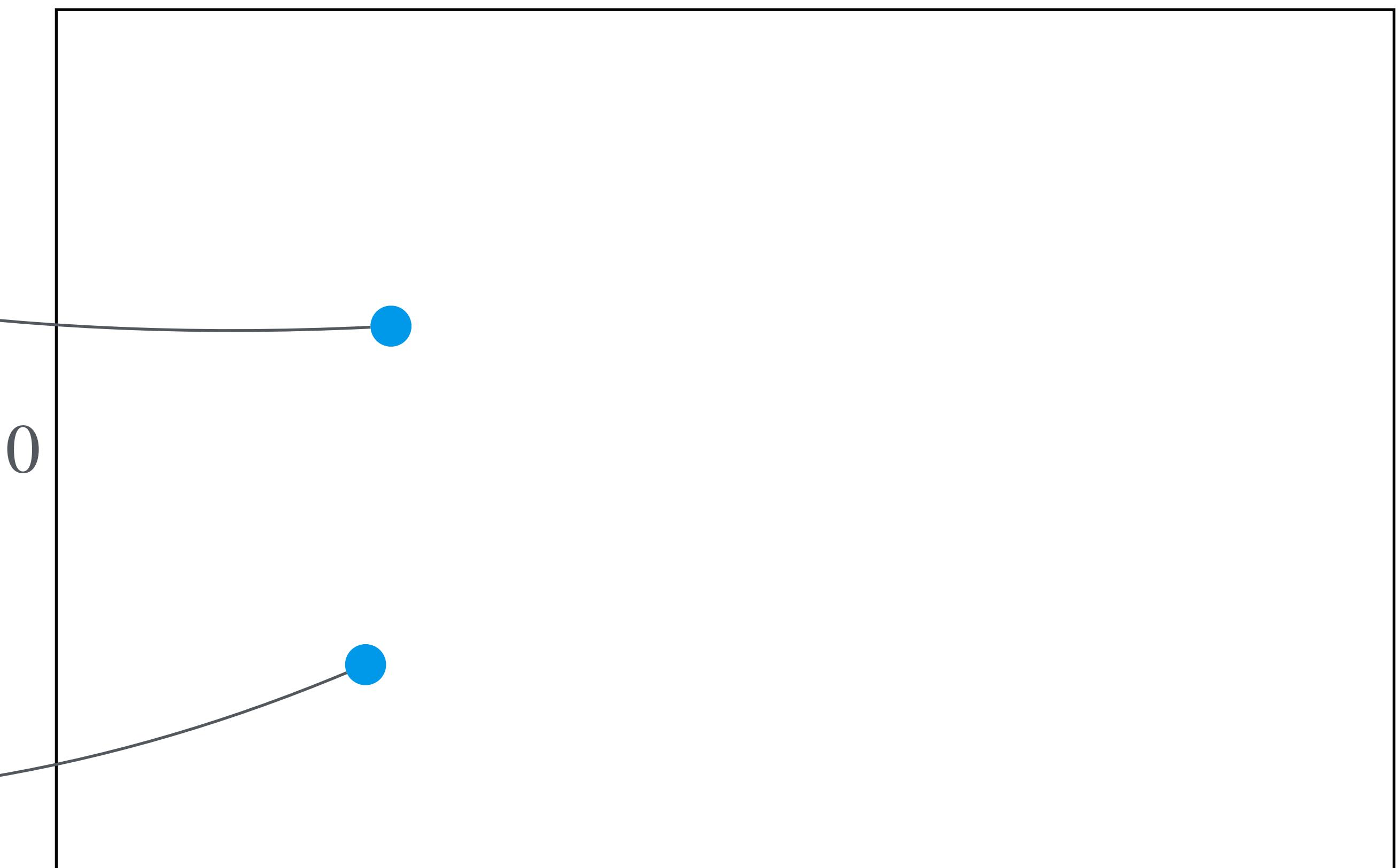
What are we doing here?



A point in hypothesis space is a model



$$w_2 = 0$$



$$w_1 = 0$$

What are we doing here?

Hard constraints partition
your hypothesis space into
feasible and infeasible regions

$$w_2 = 0$$

$$w_1 = 0$$

What are we doing here?

Hard constraints partition
your hypothesis space into
feasible and infeasible regions

$$x + y > 0$$

$$w_2 = 0$$



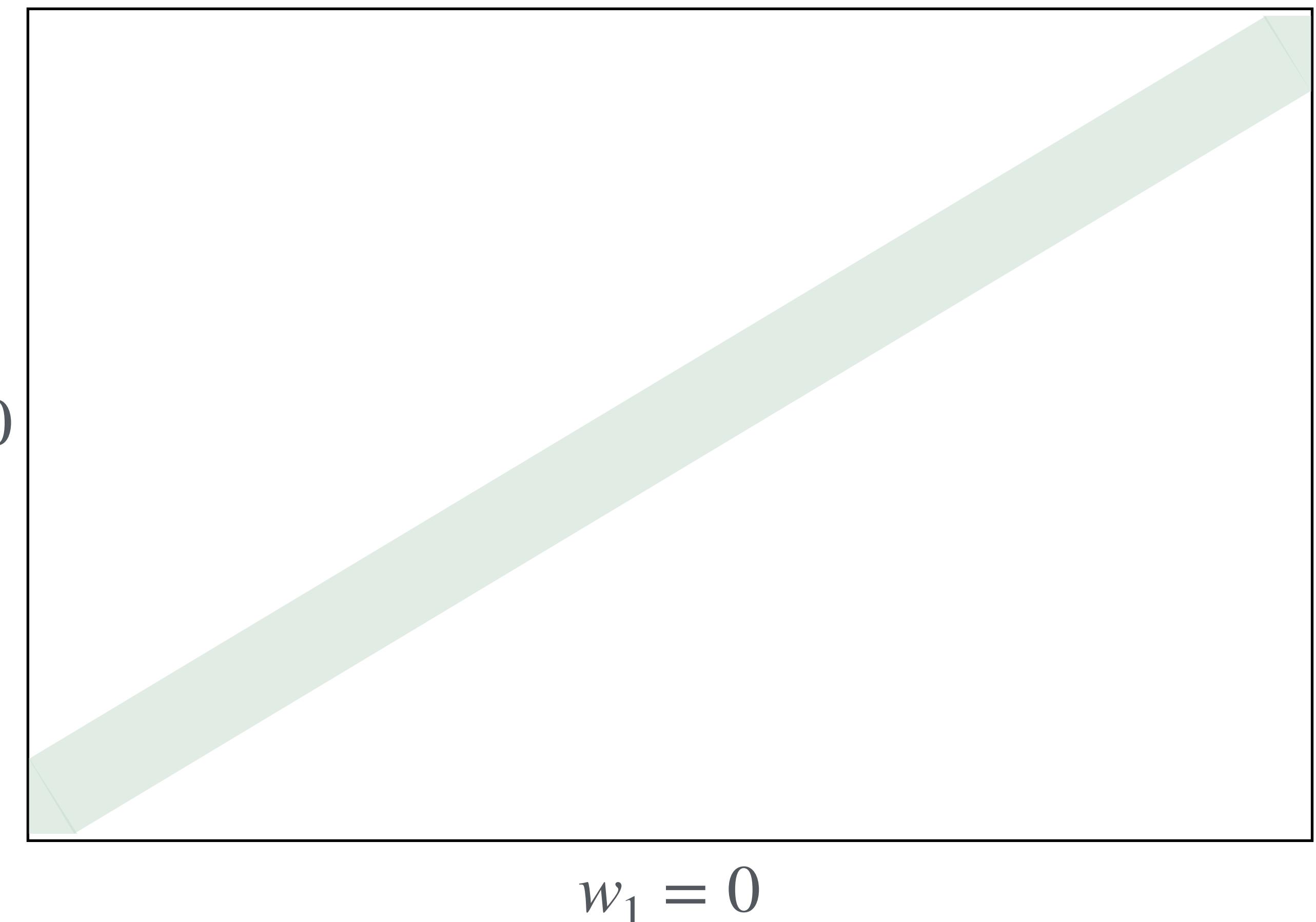
$$w_1 = 0$$

What are we doing here?

Hard constraints partition
your hypothesis space into
feasible and infeasible regions

$$x - y < 3$$

$$w_2 = 0$$



What are we doing here?

Hard constraints partition
your hypothesis space into
feasible and infeasible regions

$$x + y > 0$$

$$x - y < 3$$

$$w_2 = 0$$

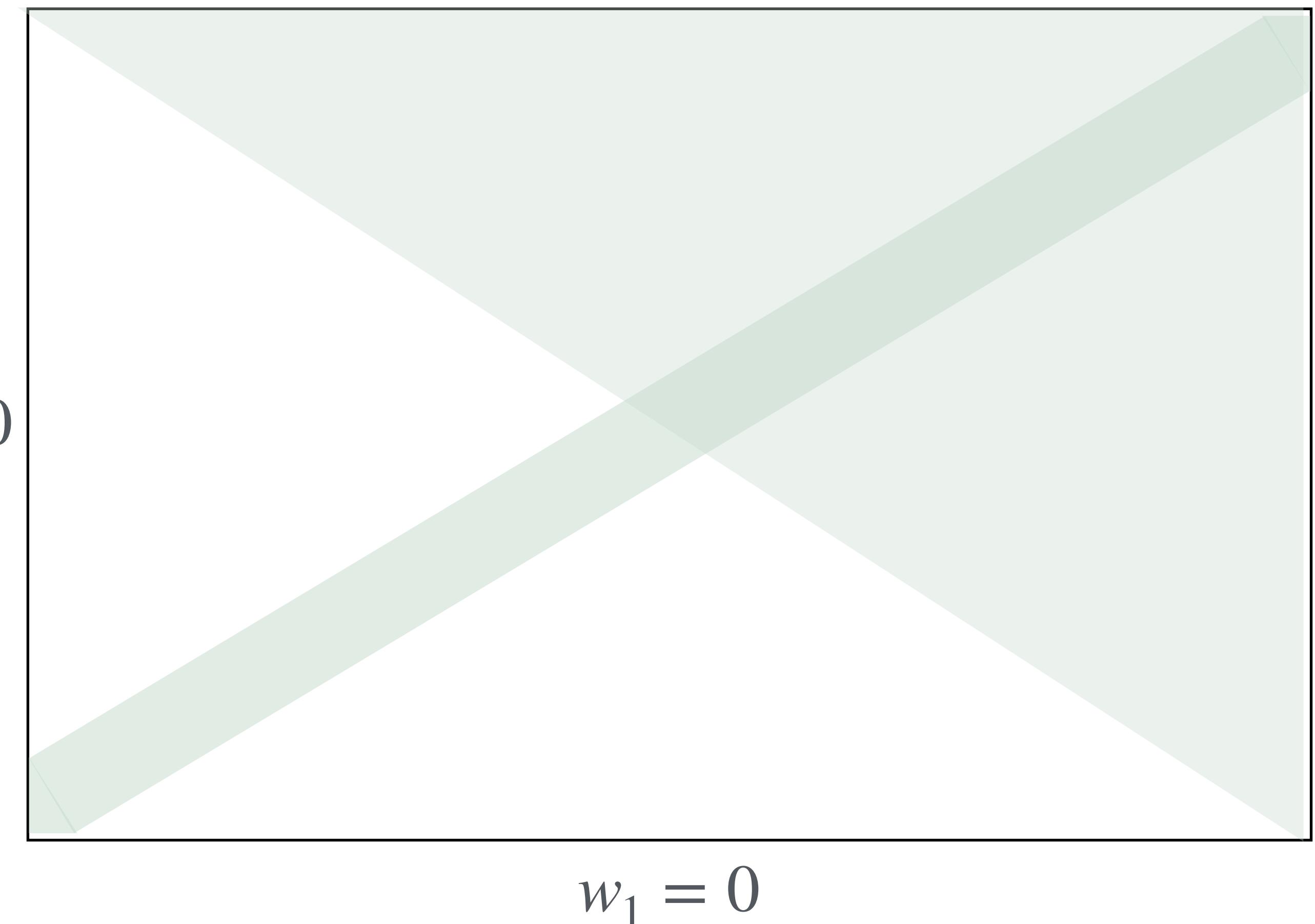
$$w_1 = 0$$



What are we doing here?

Soft constraints introduce
further optional partitioning

$$w_2 = 0$$



What are we doing here?

Soft constraints introduce further optional partitioning

$$-3x + y > 0$$

$$w_2 = 0$$

$$w_1 = 0$$



What are we doing here?

Soft constraints introduce further optional partitioning

$$-3x + y > 0$$

$$4x + y < 0$$

$$w_2 = 0$$

$$w_1 = 0$$



What are we doing here?

Soft constraints introduce further optional partitioning

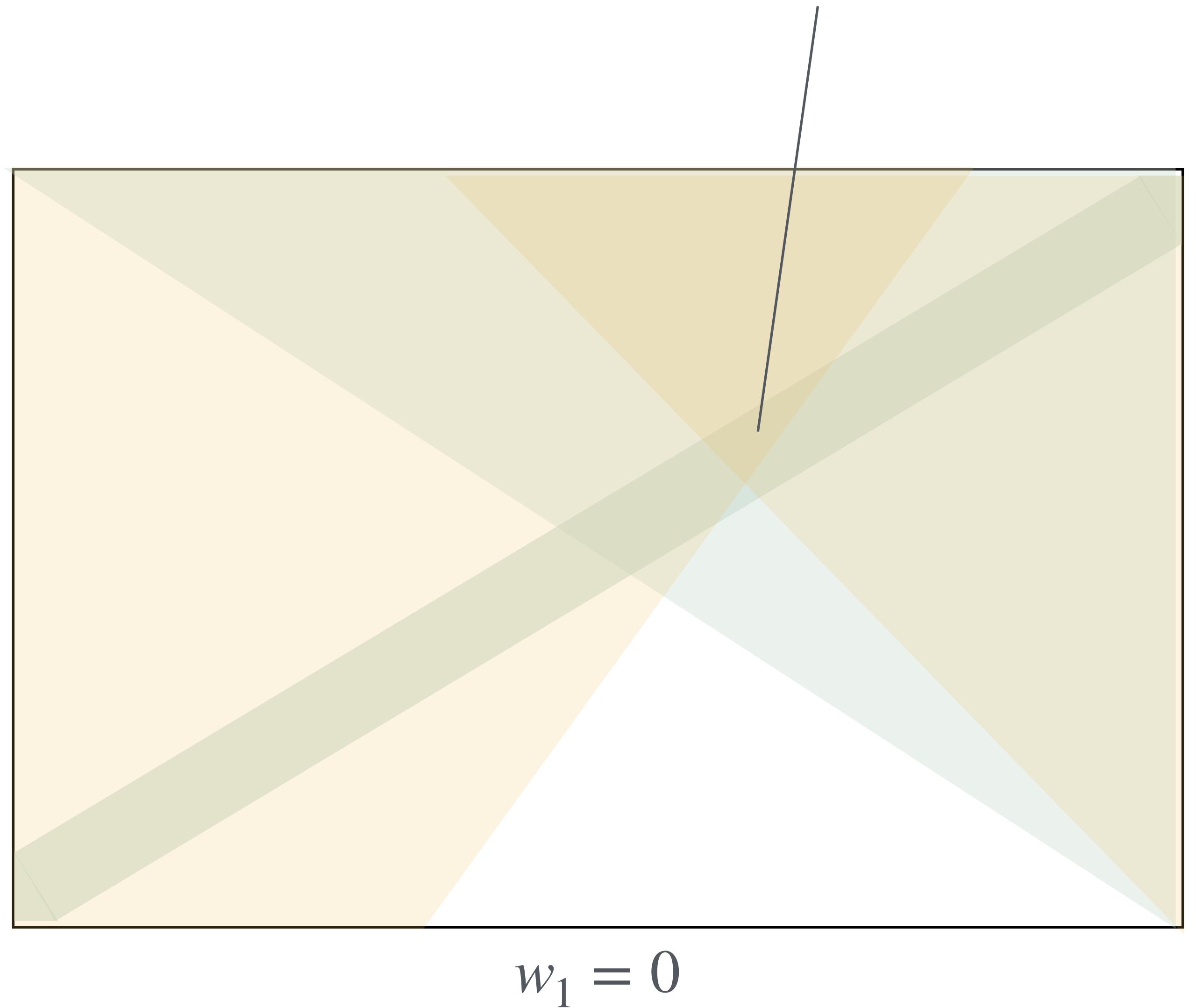
$$-3x + y > 0$$

$$4x + y < 0$$

$$w_2 = 0$$

$$w_1 = 0$$

This is the place where your model lies







This works only on small models!

Why do neural networks scale?

Why do neural networks scale?

(mini) batching

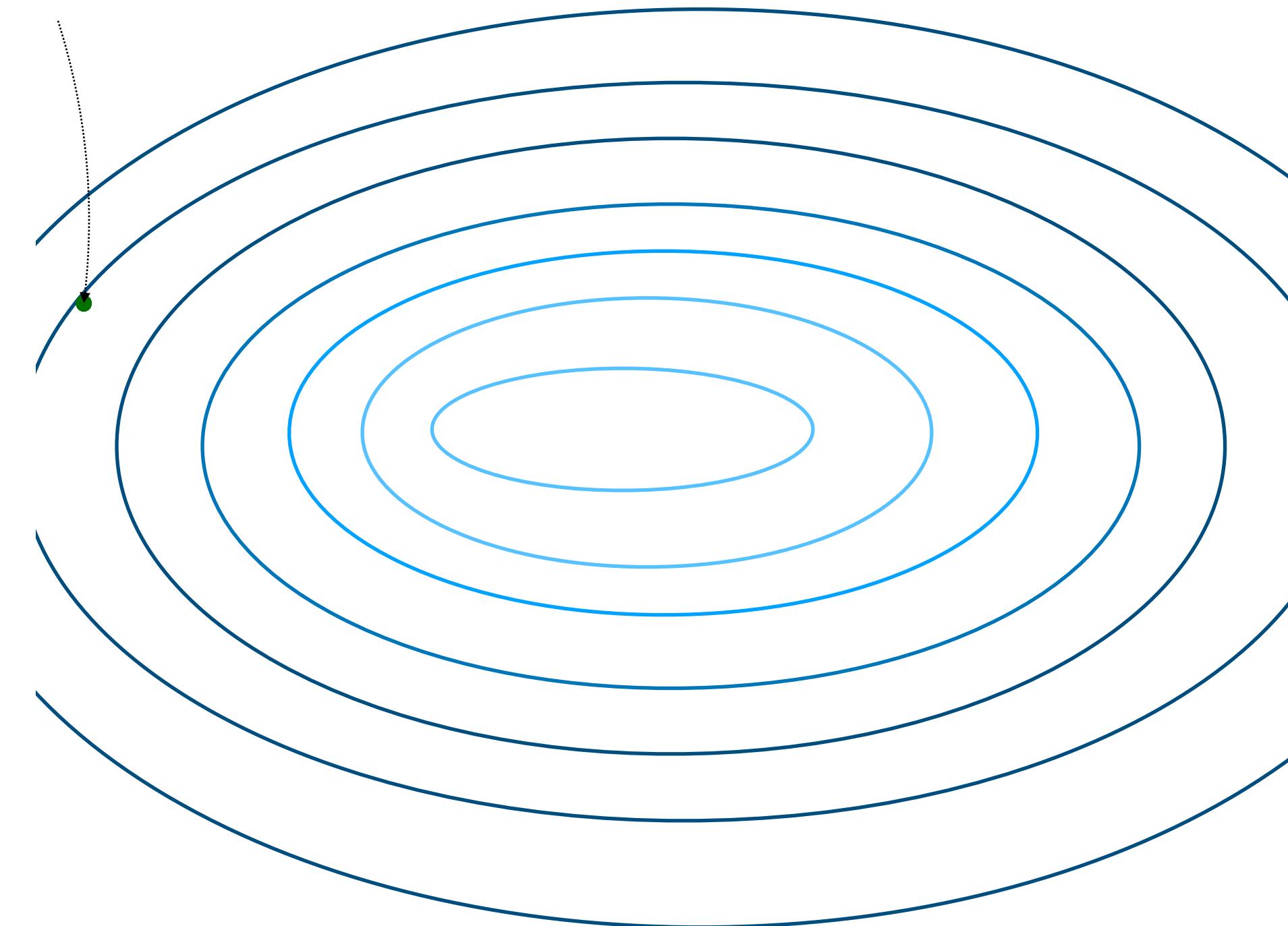
stochastic gradient descent

small local changes

Satisfiability descent

Gradient Descent
+
Maximum Satisfiability

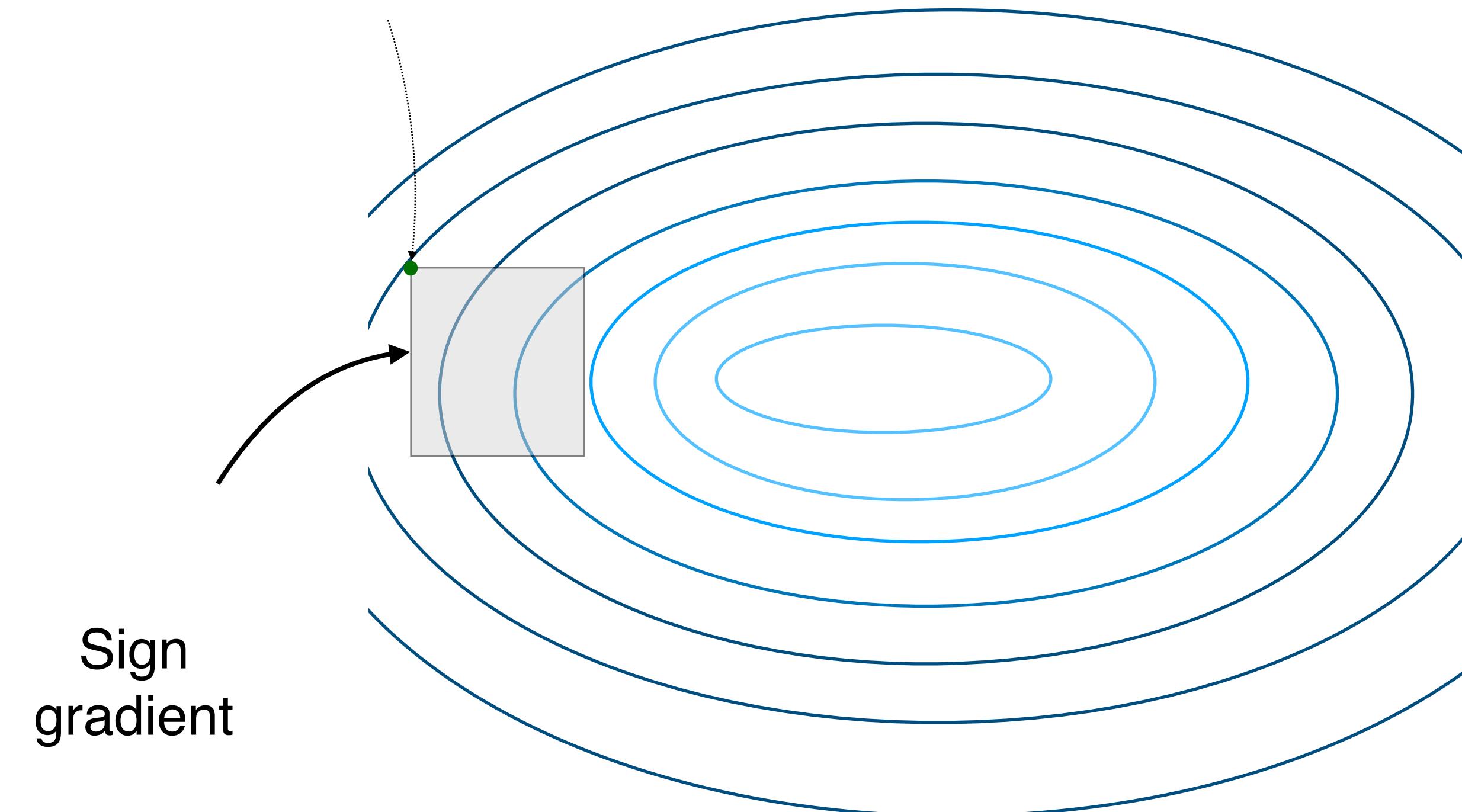
Current best
solution



Satisfiability descent

Gradient Descent
+
Maximum Satisfiability

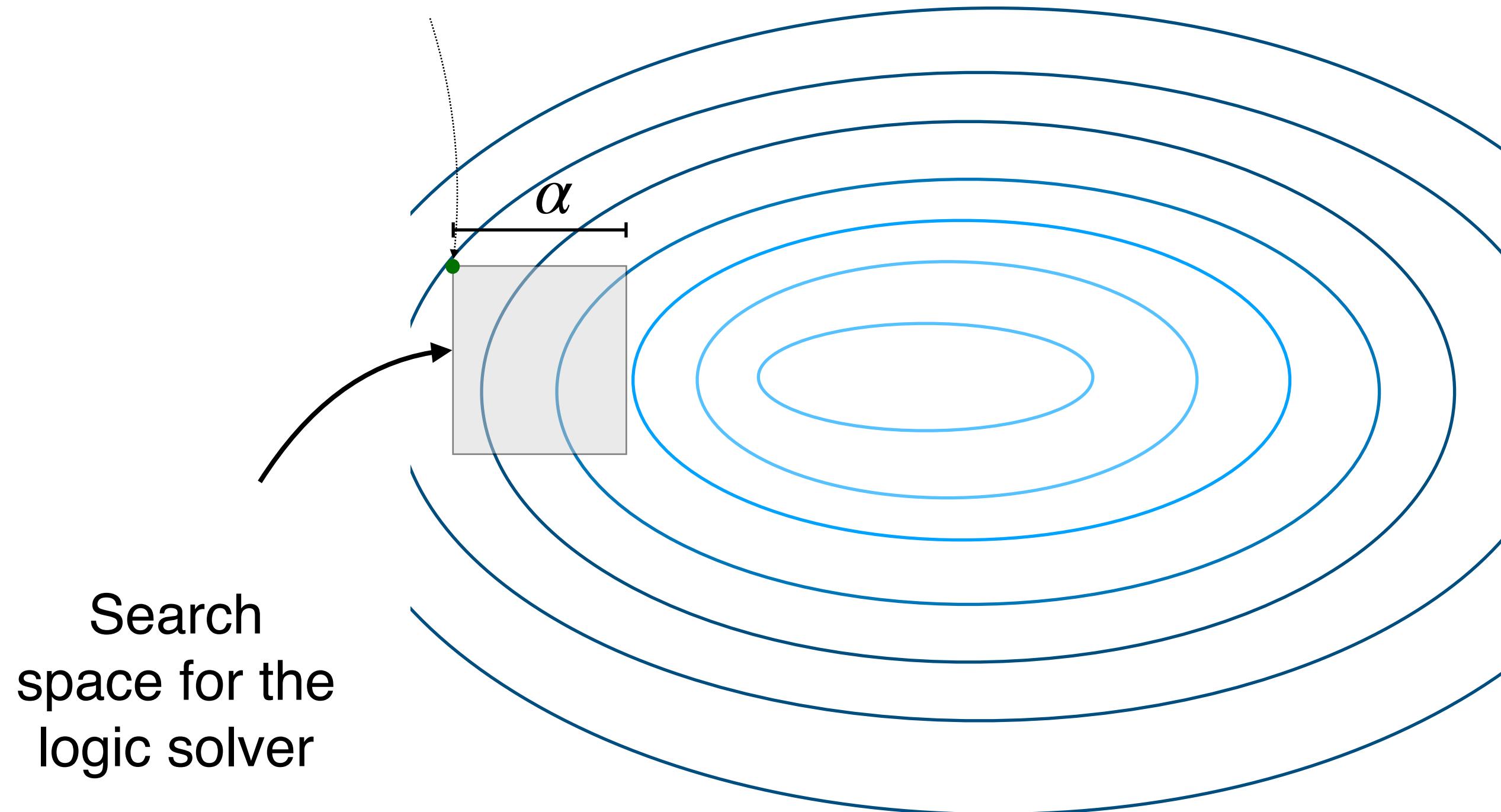
Solution
learned by
MaxSMT



Satisfiability descent

Gradient Descent
+
Maximum Satisfiability

Solution
learned by
MaxSMT



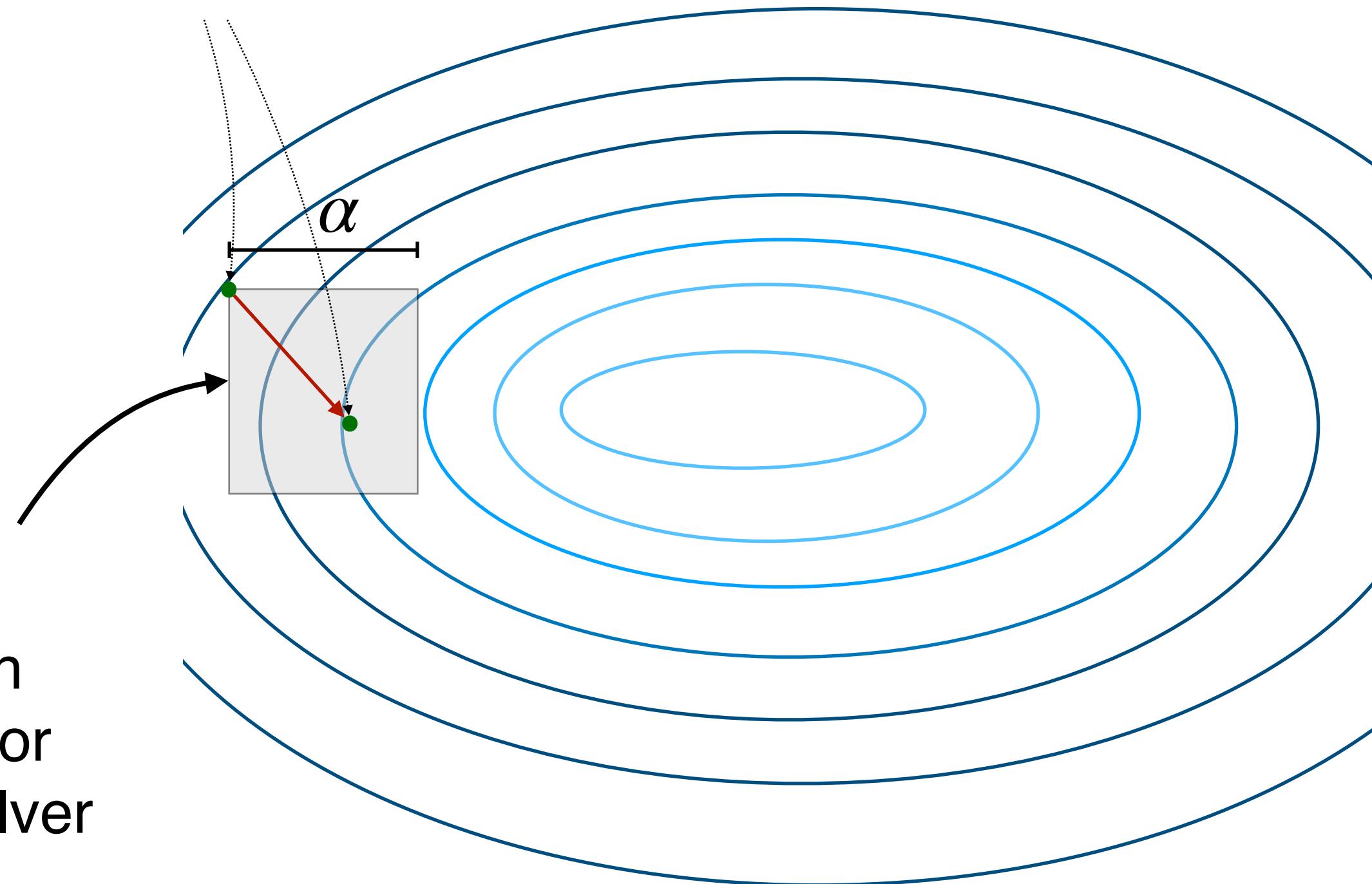
Satisfiability descent

Gradient Descent
+
Maximum Satisfiability

- ▶ Keeps the MaxSMT problem small
- ▶ Constraint is always satisfied by the learned solution.

Solution
learned by
MaxSMT

Search
space for
Logic Solver



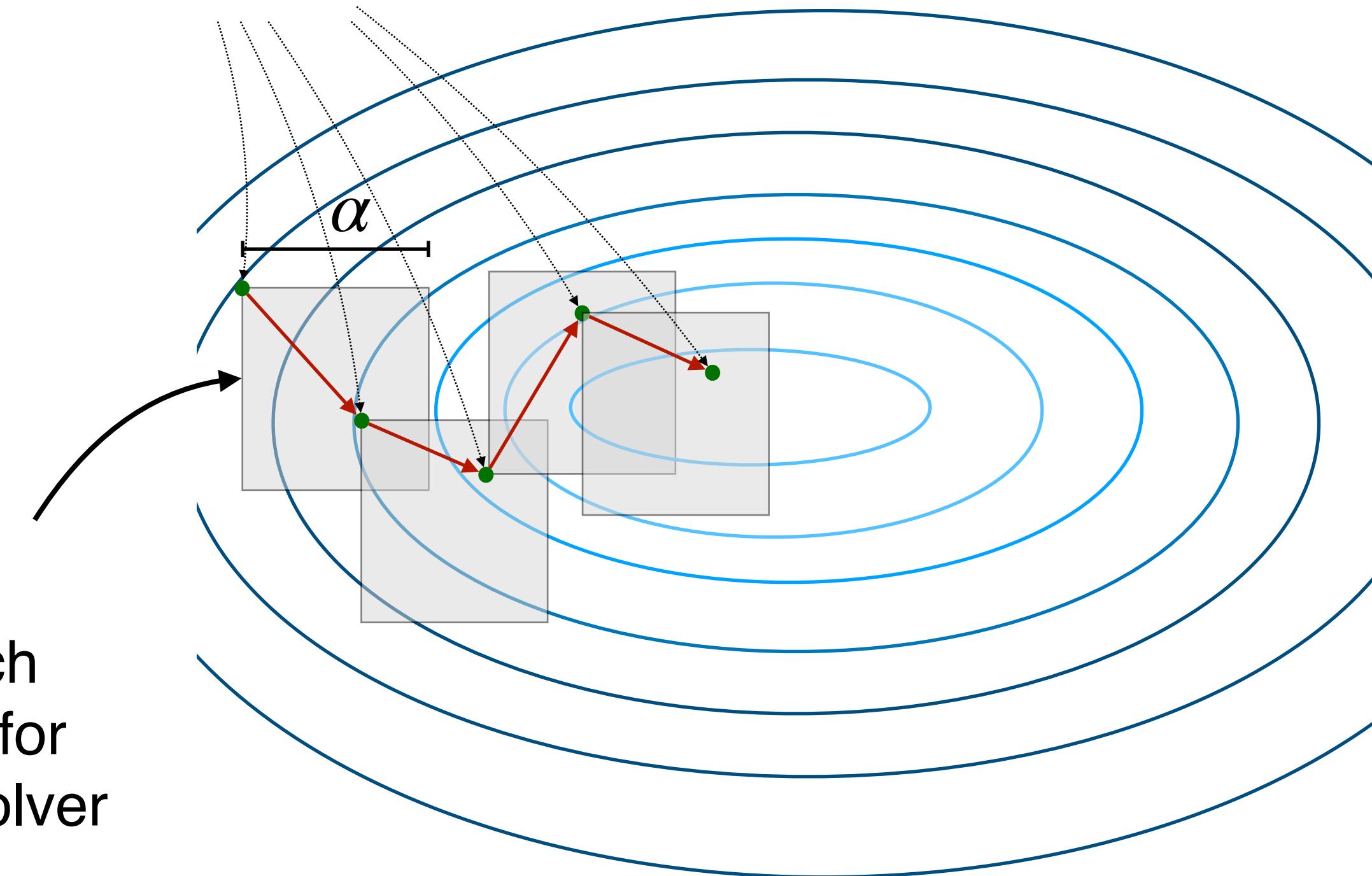
Satisfiability descent

Gradient Descent
+
Maximum Satisfiability

- ▶ Keeps the MaxSMT problem small
- ▶ Constraint is always satisfied by the learned solution.

Solution
learned by
MaxSMT

Search
space for
Logic Solver







SaDe works as long as we don't have too many features



SaDe works as long as we don't have too many features

(Perhaps surprisingly, works well for $k \times 100$ without universal constraints)



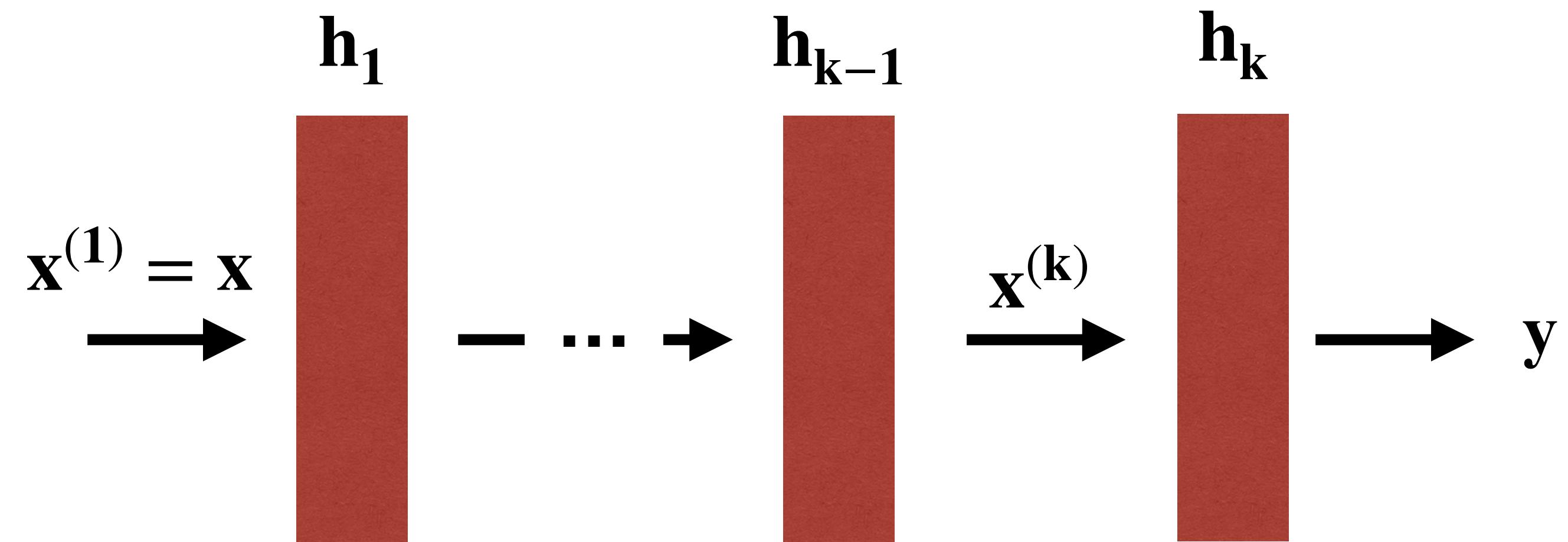
SaDe works as long as we don't have too many features

(Perhaps surprisingly, works well for $k \times 100$ without universal constraints)

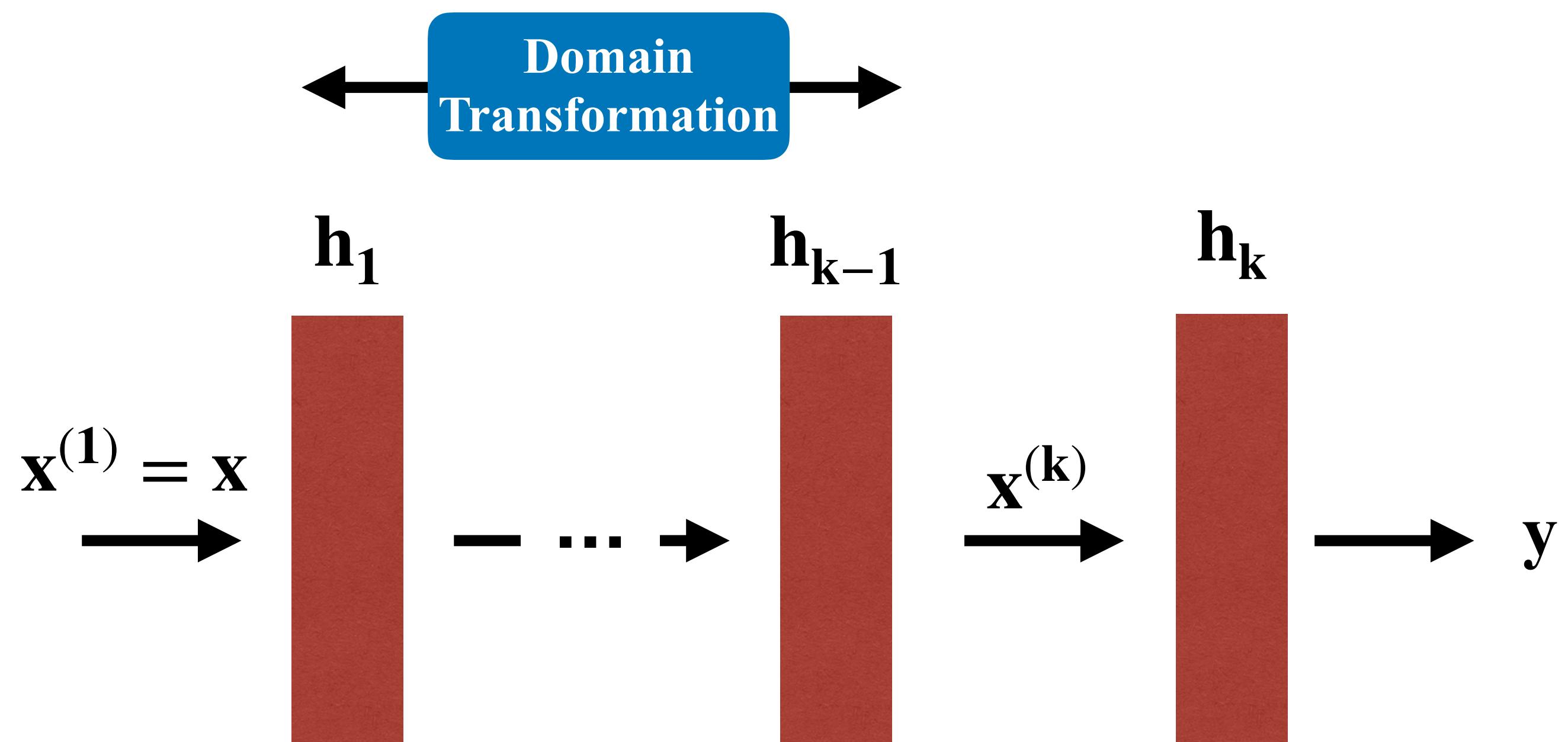
How do we get it to work on much bigger models?

DeepSaDe:
SaDe for deep neural networks

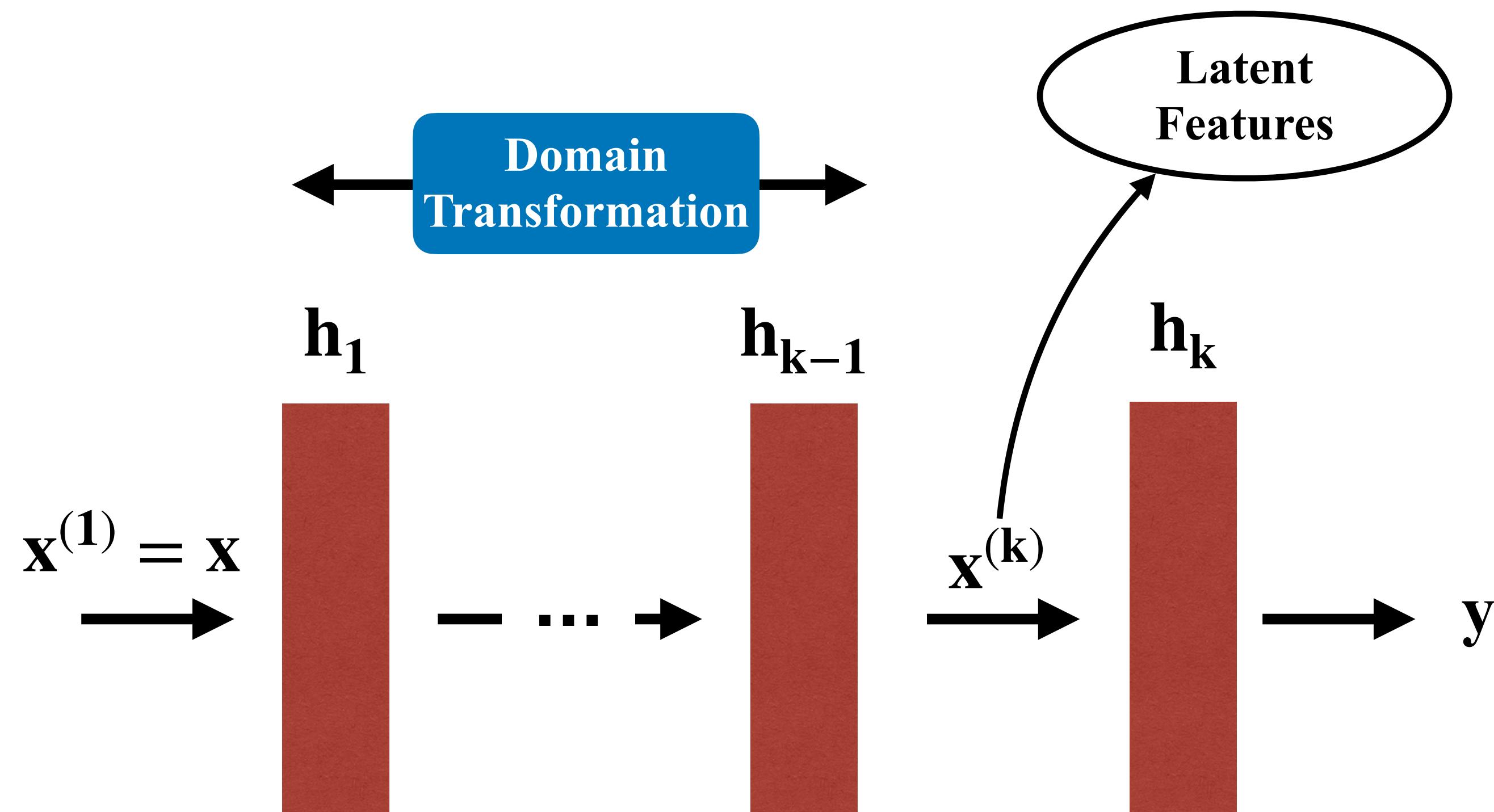
Deep satisfiability descent



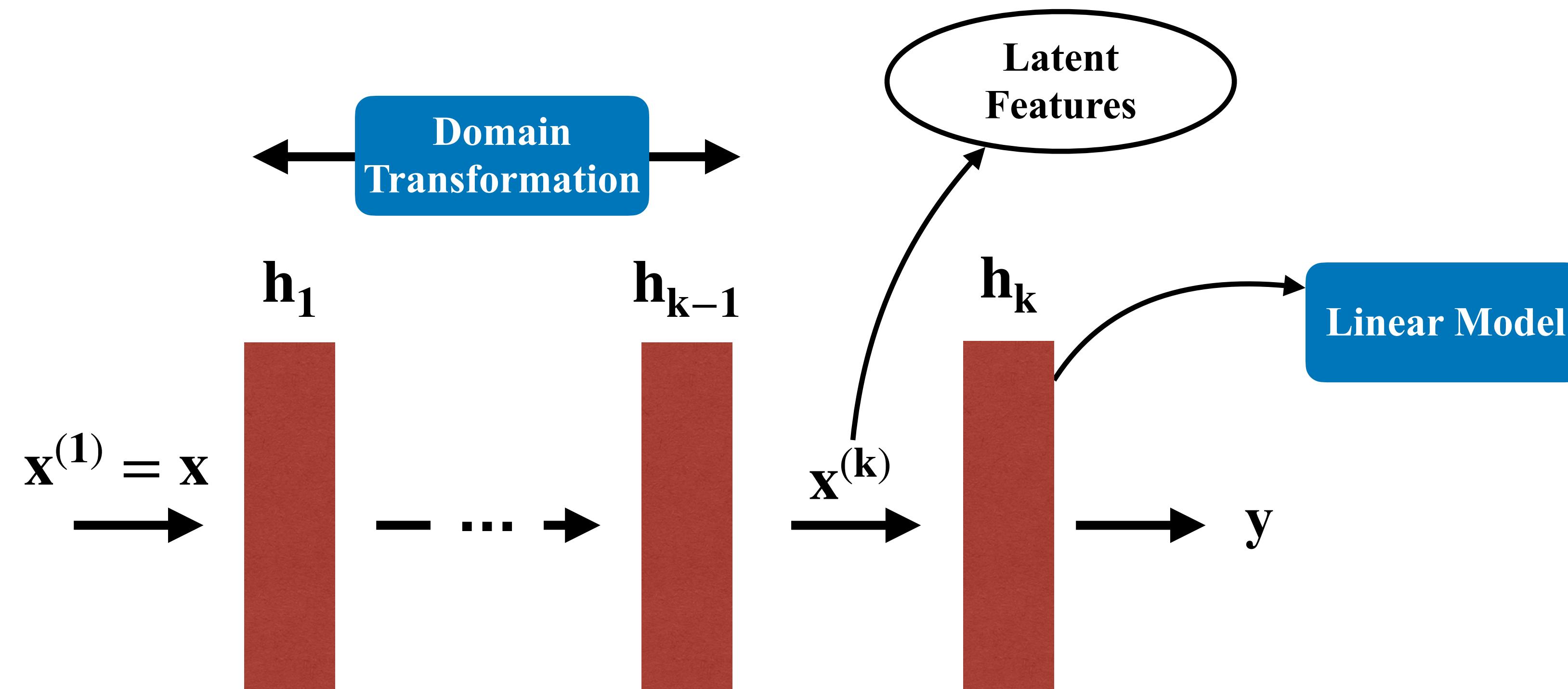
Deep satisfiability descent



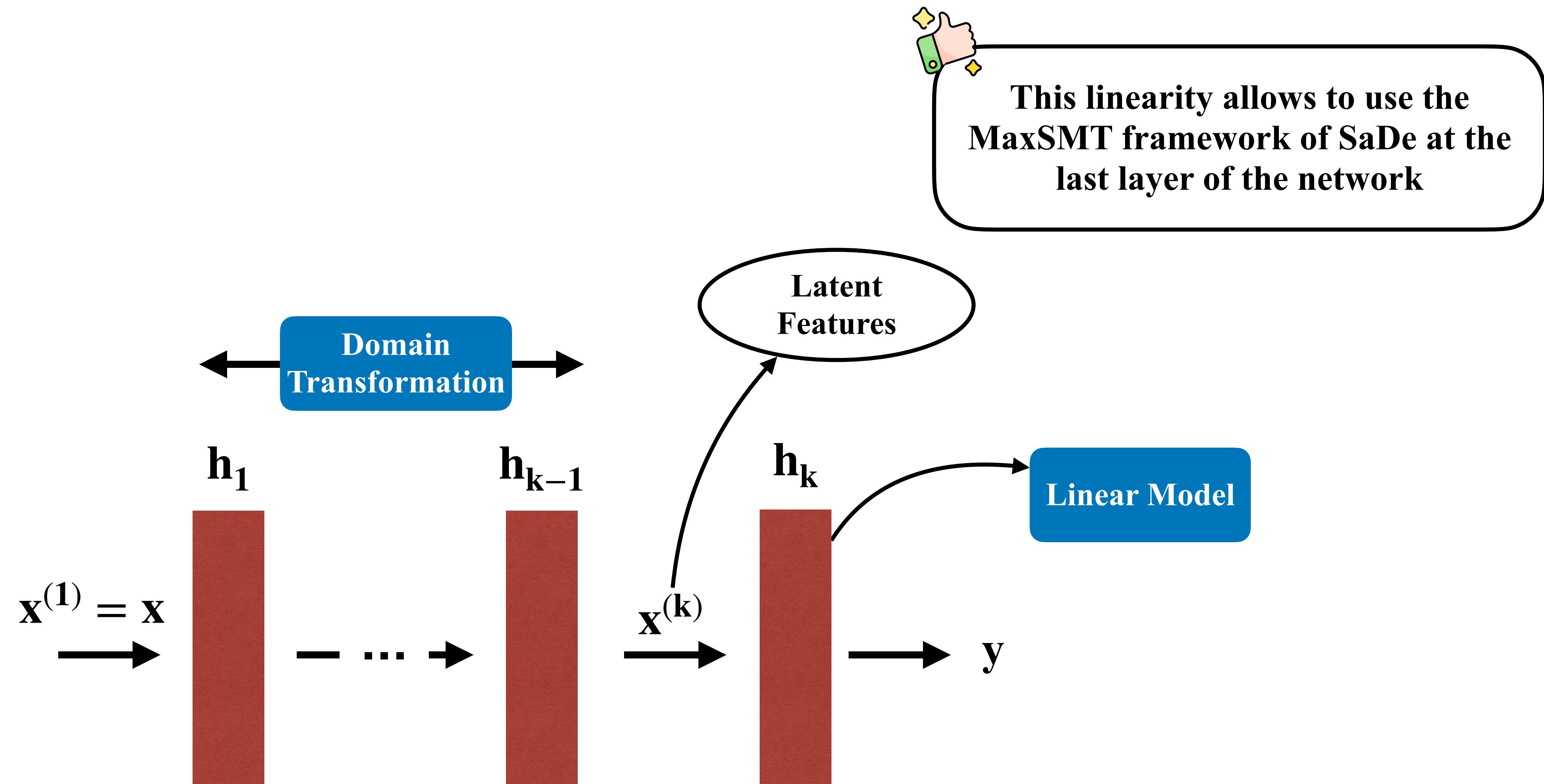
Deep satisfiability descent



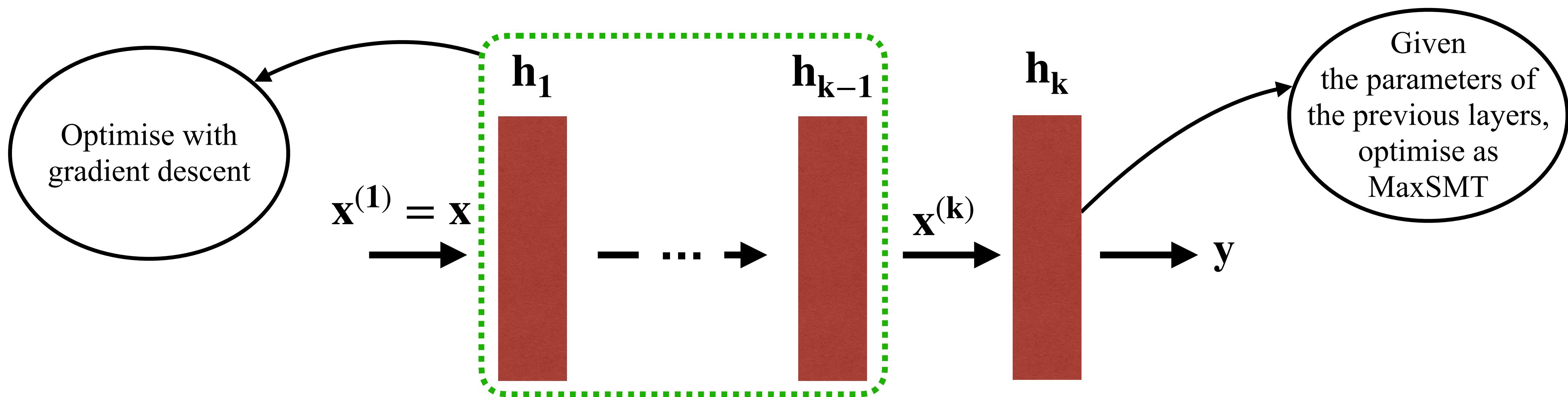
Deep satisfiability descent



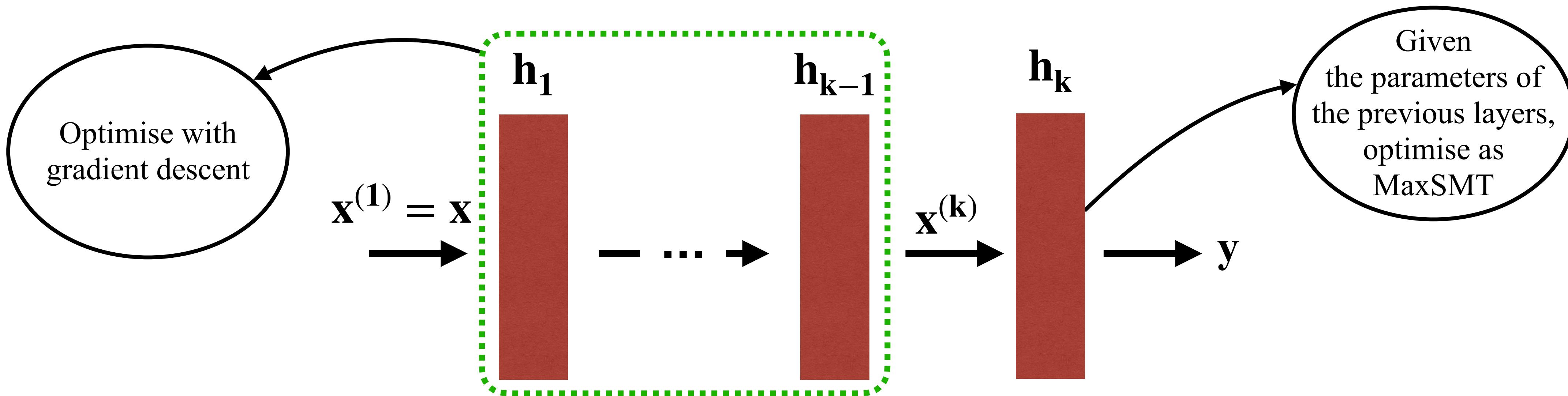
Deep satisfiability descent



Deep satisfiability descent



Deep satisfiability descent



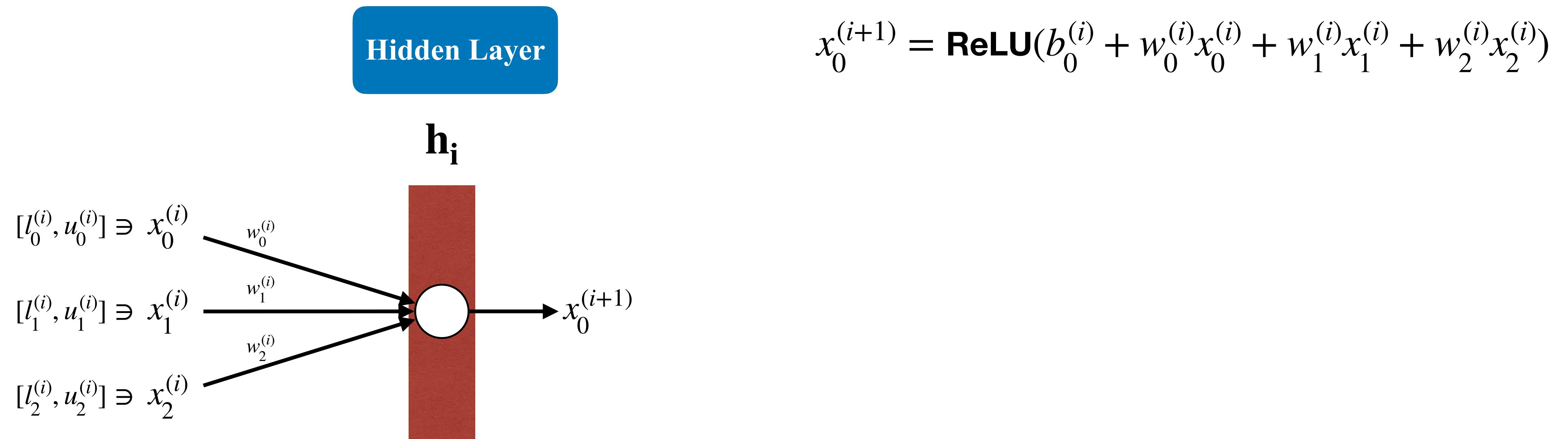
Simply rewrite soft constraints:

$$\begin{aligned} f_w([1,4]) &< 0 \\ f_w([2,3]) &< 0 \\ f_w([-3,2]) &> 0 \end{aligned}$$



$$\begin{aligned} h_k(h_{k-1} \dots ([1,4]) \dots) &< 0 \\ h_k(h_{k-1} \dots ([2,3]) \dots) &< 0 \\ h_k(h_{k-1} \dots ([-3,2]) \dots) &> 0 \end{aligned}$$

Computing domain translation

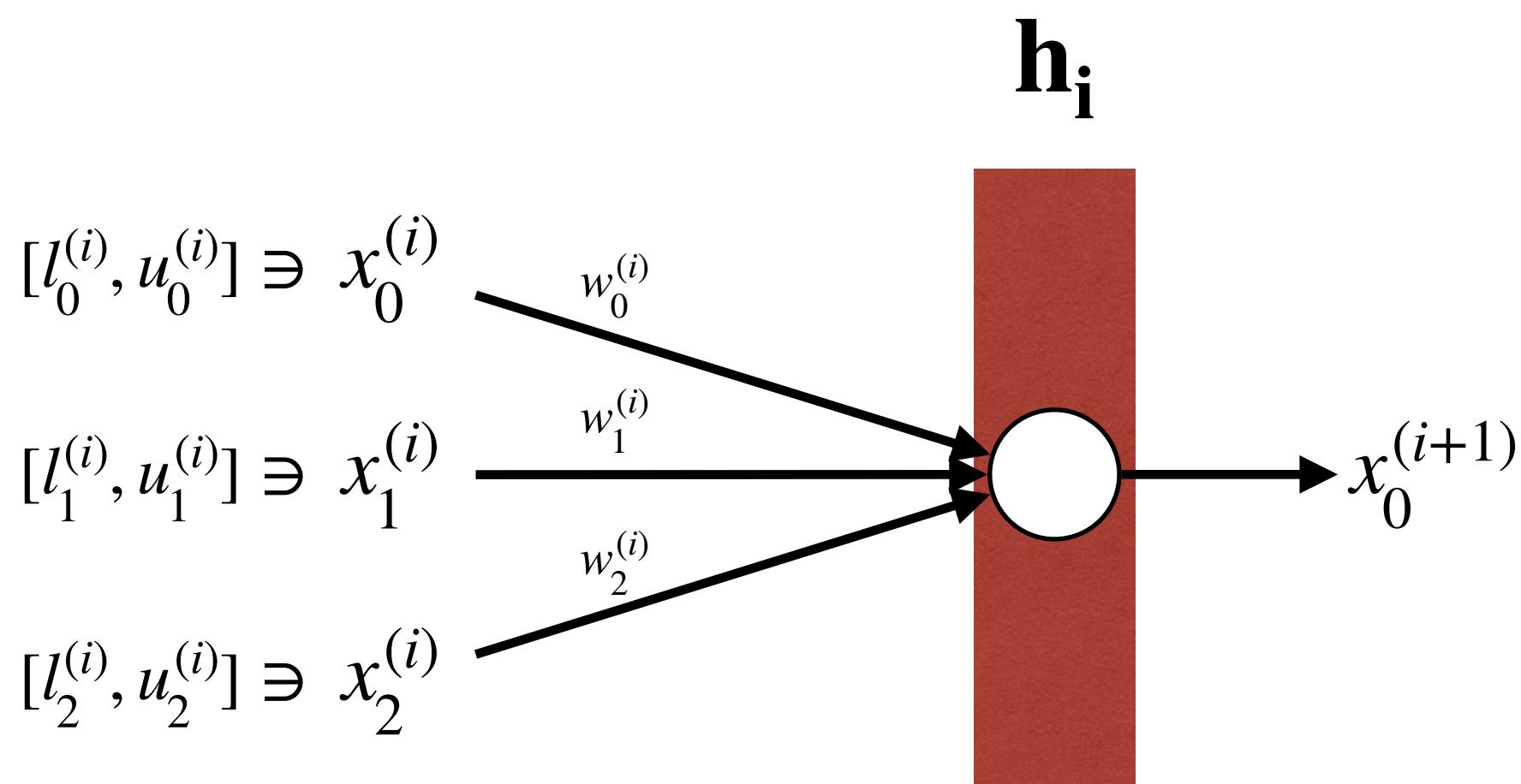


Computing domain translation

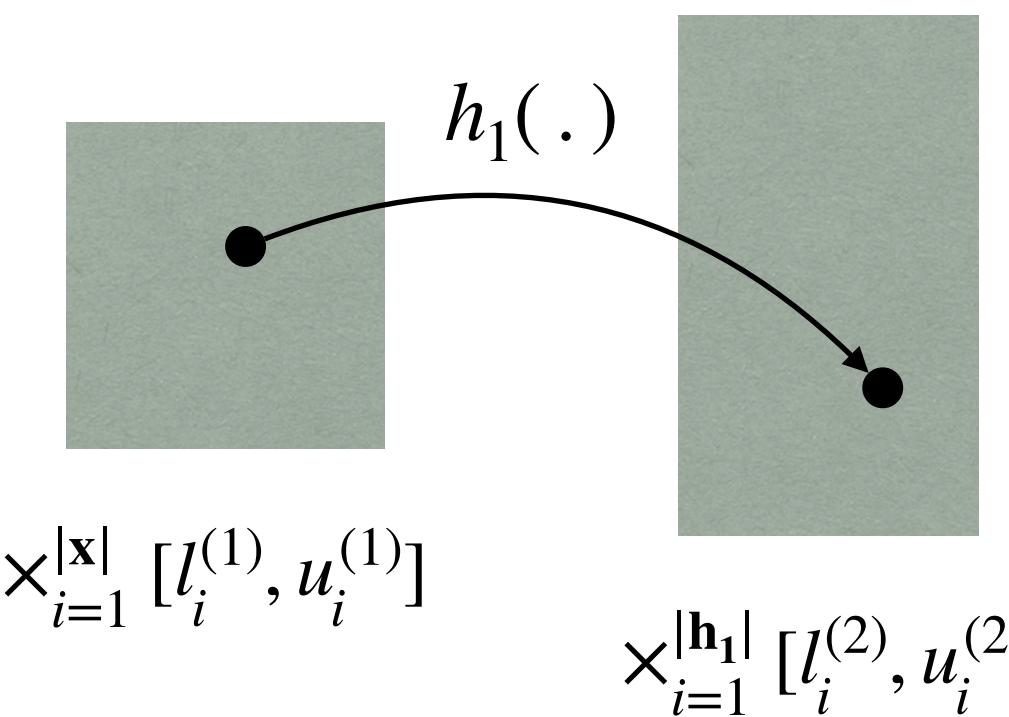
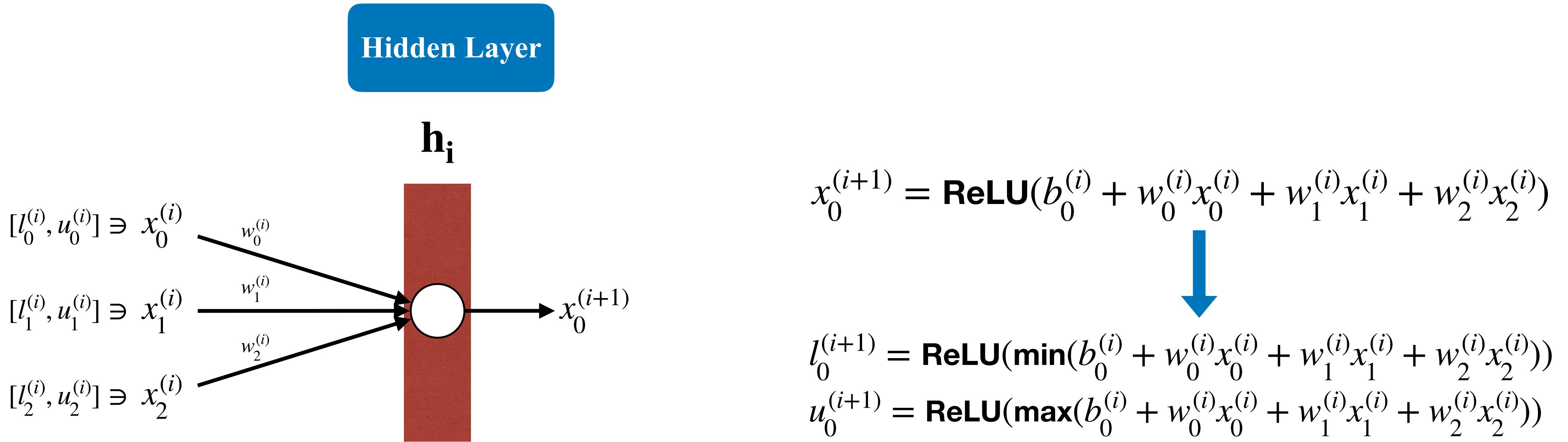
Hidden Layer

$$x_0^{(i+1)} = \mathbf{ReLU}(b_0^{(i)} + w_0^{(i)}x_0^{(i)} + w_1^{(i)}x_1^{(i)} + w_2^{(i)}x_2^{(i)})$$

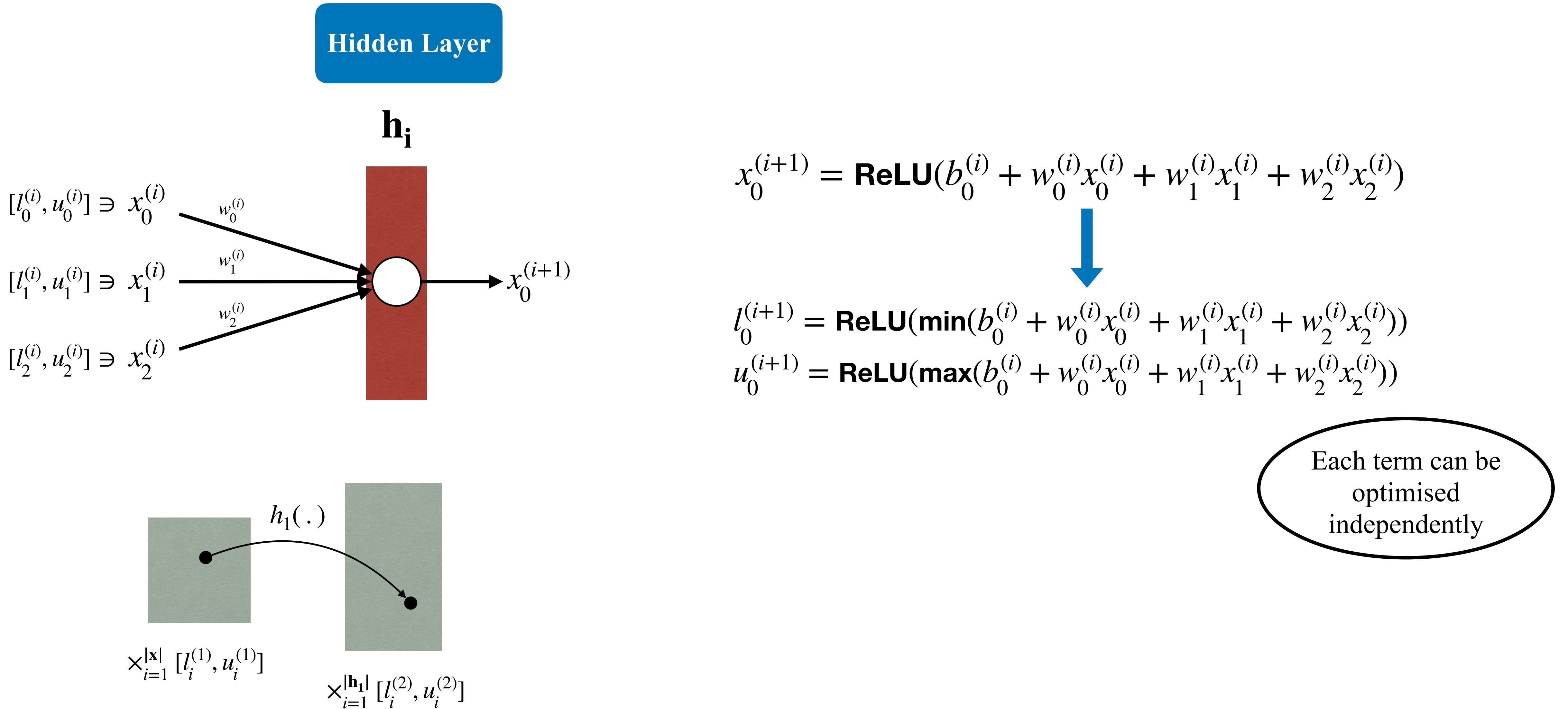
Monotonically
non-decreasing



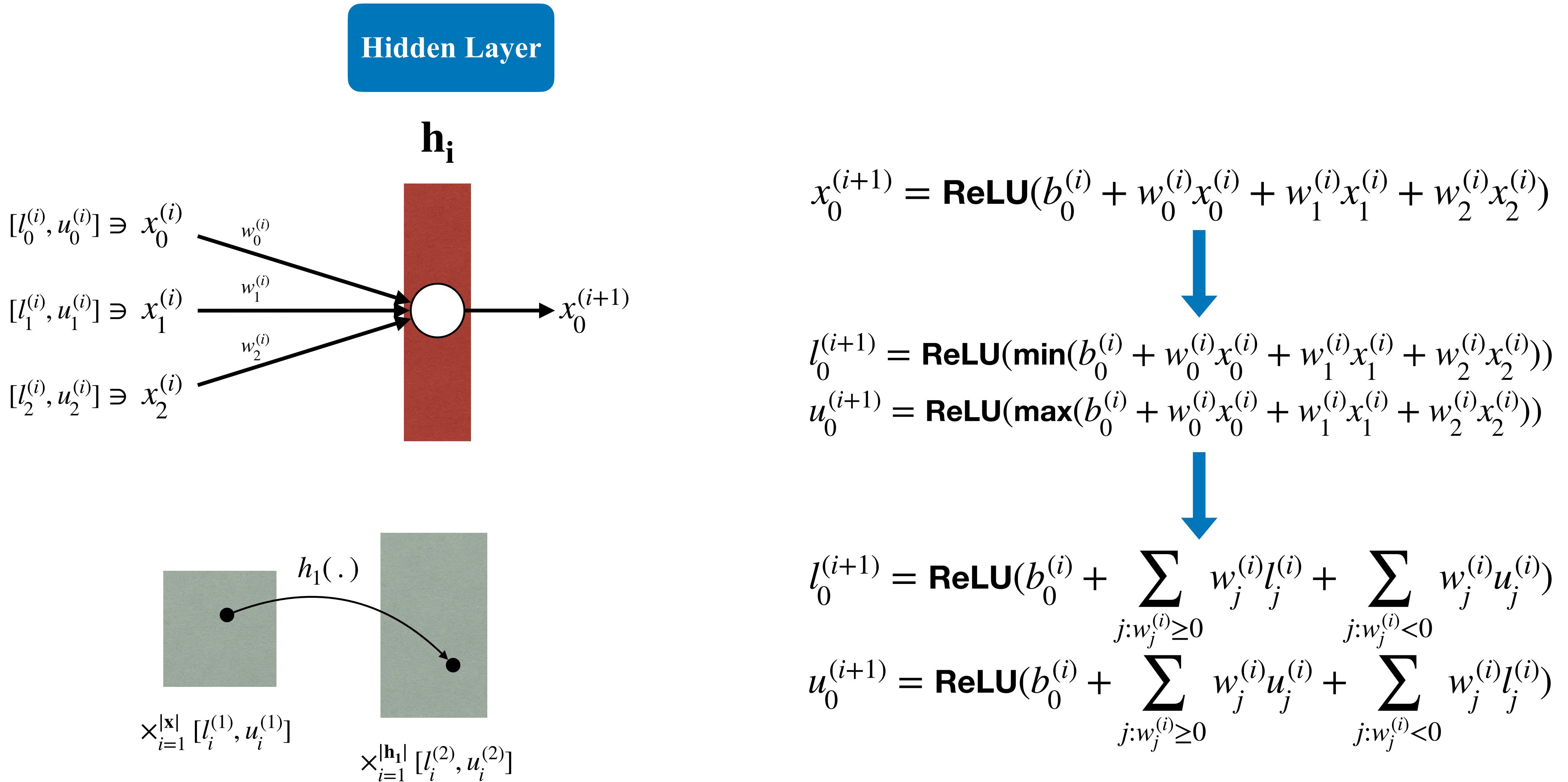
Computing domain translation



Computing domain translation



Computing domain translation



Experiments

1

Task: Multi-target regression to predict 5 different household expenses based on certain features

Constraint:

Sum of the predicted expenses must be smaller than “income” & going out expense must be smaller than 5% of the “income”

2

Task: Binary classification to predict if a person should be given the loan by a bank

Constraint:

People with absent “credit history” and low “income” must be denied the loan

3

Task: Multi-class classification to predict the genre of a song

Constraint:

A song by “the beatles” must only be predicted as either pop or rock

Experiments

4

Task: Multi label classification to identify labels from a sequence of 4 MNIST images

 → [1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0]

Constraint:

Sum of the predicted labels must be greater than 10

5

Task: Preference learning task to predict preference order of sushi.

Sake, Unagi, Ika → [1, 0, 0, 0, 0, 1, 0, 1, 0]

Constraint:

The preference order must be coherent [[Xu et al.](#)]

How do we evaluate these models?

Accuracy?

How do we evaluate these models?

Accuracy?

What if your historical data violates your constraints?

How do we evaluate these models?

Accuracy?

What if your historical data violates your constraints?

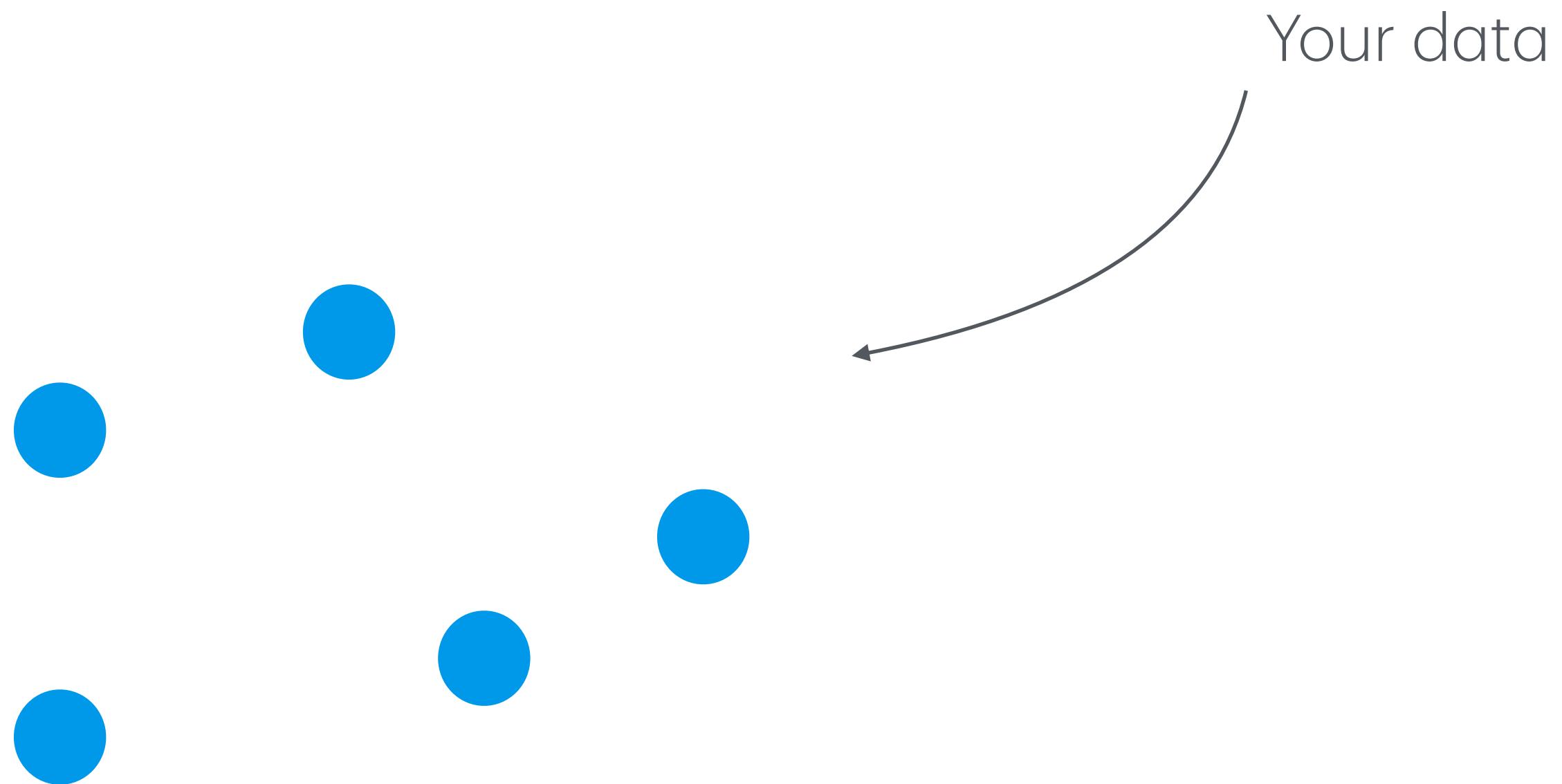
Adversity index

How do we evaluate these models?

Accuracy?

What if your historical data violates your constraints?

Adversity index



How do we evaluate these models?

Accuracy?

What if your historical data violates your constraints?

Adversity index

Look at small
balls around
your data



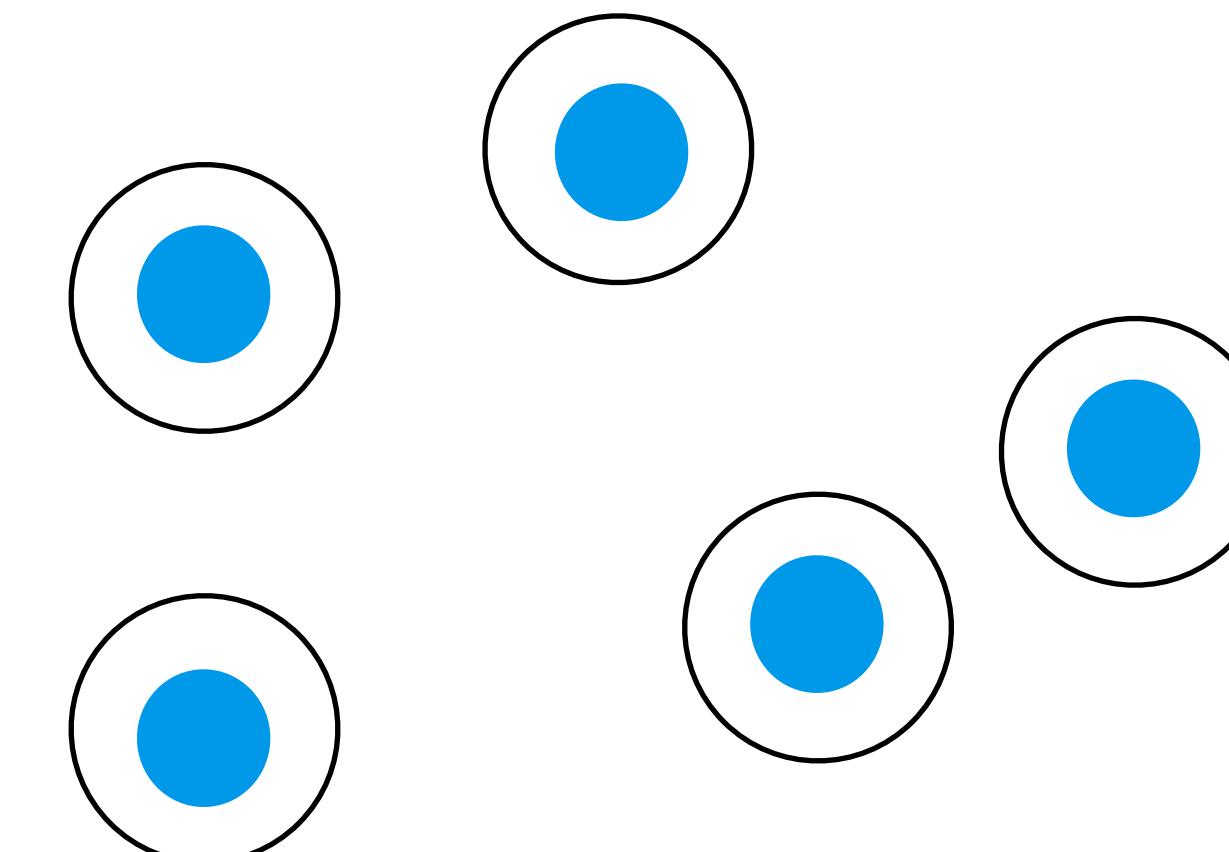
How do we evaluate these models?

Accuracy?

What if your historical data violates your constraints?

Adversity index

Look at small
balls around
your data



In how many of these balls
can I find a counterexample?

Experiments: Results

UC	Approach	Constraint	AdI($\delta = 0.1$)	Accuracy/MSE	Runtime (sec)
UC1	DeepSaDe	100\pm0	0\pm0	*38.36 \pm 4.59	102341 \pm 40198
	REG	93.50 \pm 2.01	0.97 \pm 0.003	*30.50\pm6.49	227 \pm 70
UC2	DeepSaDe	100\pm0	0\pm0	80.04 \pm 4.29	447 \pm 105
	SL	100\pm0	0.002 \pm 0.006	80.17\pm3.88	45 \pm 30
	SBR	100\pm0	0.002 \pm 0.004	80.04 \pm 3.95	45 \pm 27
UC3	DeepSaDe	100\pm0	0\pm0	80.11 \pm 4.99	6580 \pm 1915
	SL	99.97 \pm 0.03	0.42 \pm 0.28	82.53\pm3.58	101 \pm 45
	SBR	99.97 \pm 0.05	0.29 \pm 0.10	76.51 \pm 7.50	196 \pm 68

Experiments: Results

UC	Approach	Constraint	AdI($\delta = 0.1$)	Accuracy/MSE	Runtime (sec)
UC1	DeepSaDe	100 ± 0	0 ± 0	$*38.36 \pm 4.59$	102341 ± 40198
	REG	93.50 ± 2.01	0.97 ± 0.003	$*30.50 \pm 6.49$	227 ± 70
UC2	DeepSaDe	100 ± 0	0 ± 0	80.04 ± 4.29	447 ± 105
	SL	100 ± 0	0.002 ± 0.006	80.17 ± 3.88	45 ± 30
	SBR	100 ± 0	0.002 ± 0.004	80.04 ± 3.95	45 ± 27
UC3	DeepSaDe	100 ± 0	0 ± 0	80.11 ± 4.99	6580 ± 1915
	SL	99.97 ± 0.03	0.42 ± 0.28	82.53 ± 3.58	101 ± 45
	SBR	99.97 ± 0.05	0.29 ± 0.10	76.51 ± 7.50	196 ± 68

Existing methods do not find models
that guarantee domain constraint
satisfaction

Experiments: Results

DeepSaDe's hard guarantees lead to drop in predictive performance in some cases compared to baselines

UC	Approach	Constraint	AdI($\delta = 0.1$)	Accuracy/MSE	Runtime (sec)
UC1	DeepSaDe	100 ± 0	0 ± 0	$*38.36 \pm 4.59$	102341 ± 40198
	REG	93.50 ± 2.01	0.97 ± 0.003	$*30.50 \pm 6.49$	227 ± 70
UC2	DeepSaDe	100 ± 0	0 ± 0	80.04 ± 4.29	447 ± 105
	SL	100 ± 0	0.002 ± 0.006	80.17 ± 3.88	45 ± 30
	SBR	100 ± 0	0.002 ± 0.004	80.04 ± 3.95	45 ± 27
UC3	DeepSaDe	100 ± 0	0 ± 0	80.11 ± 4.99	6580 ± 1915
	SL	99.97 ± 0.03	0.42 ± 0.28	82.53 ± 3.58	101 ± 45
	SBR	99.97 ± 0.05	0.29 ± 0.10	76.51 ± 7.50	196 ± 68

Existing methods do not find models that guarantee domain constraint satisfaction

Experiments: Results

DeepSaDe's hard guarantees lead to drop in predictive performance in some cases compared to baselines

UC	Approach	Constraint	AdI($\delta = 0.1$)	Accuracy/MSE	Runtime (sec)
UC1	DeepSaDe	100 ± 0	0 ± 0	$*38.36 \pm 4.59$	102341 ± 40198
	REG	93.50 ± 2.01	0.97 ± 0.003	$*30.50 \pm 6.49$	227 ± 70
UC2	DeepSaDe	100 ± 0	0 ± 0	80.04 ± 4.29	447 ± 105
	SL	100 ± 0	0.002 ± 0.006	80.17 ± 3.88	45 ± 30
	SBR	100 ± 0	0.002 ± 0.004	80.04 ± 3.95	45 ± 27
UC3	DeepSaDe	100 ± 0	0 ± 0	80.11 ± 4.99	6580 ± 1915
	SL	99.97 ± 0.03	0.42 ± 0.28	82.53 ± 3.58	101 ± 45
	SBR	99.97 ± 0.05	0.29 ± 0.10	76.51 ± 7.50	196 ± 68

Existing methods do not find models that guarantee domain constraint satisfaction

Training time of DeepSaDe is significantly higher than the existing methods

Experiments: Results

UC	Approach	Constraint	Coherent	Flattened	Jaccard	Runtime (sec)
UC4	DeepSaDe	100\pm0	6.62 \pm 1.72	78.52 \pm 1.73	62.97 \pm 2.28	227928 \pm 30559
	FFN	88.00 \pm 3.26	23.94\pm4.25	85.81\pm1.18	71.55\pm2.40	3215 \pm 2663
UC5	DeepSaDe	100\pm0	11.08\pm2.61	67.17 \pm 1.48	25.94\pm2.89	17586 \pm 5074
	FFN	0.04 \pm 0.15	0.01 \pm 0.04	75.69\pm0.15	13.04 \pm 1.06	48 \pm 10
	SL	100\pm0	4.06 \pm 3.33	63.16 \pm 2.62	18.08 \pm 3.33	298 \pm 110

Existing methods have even worse constraint satisfaction for complex constraints

DeepSaDe's hard guarantees lead to drop in predictive performance in some cases compared to baselines

Training time of DeepSaDe is significantly higher than the existing methods

Summary

We still don't know how to learn models that provably obey constraints

Training models with guarantees requires new learning techniques

SaDe combines maximum satisfiability with gradient descent for performance and guarantees

To cover big models, you only need to consider “the end” of the model