

**UCLA**

**Computer  
Science**

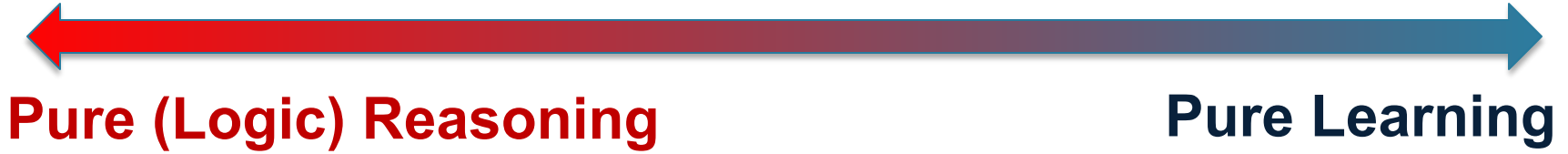


# Reasoning about Learned Models' Behavior

Guy Van den Broeck

Human-Centered AI Conference, Pepperdine - Oct 23, 2021

# The AI Dilemma



# The AI Dilemma



**Pure (Logic) Reasoning**

**Pure Learning**

- Slow thinking: deliberative, cognitive, model-based, extrapolation
- Human-centered
- Amazing achievements until this day
- “*Pure logic is brittle*”  
noise, uncertainty, incomplete knowledge, ...



# The AI Dilemma



**Pure (Logic) Reasoning**

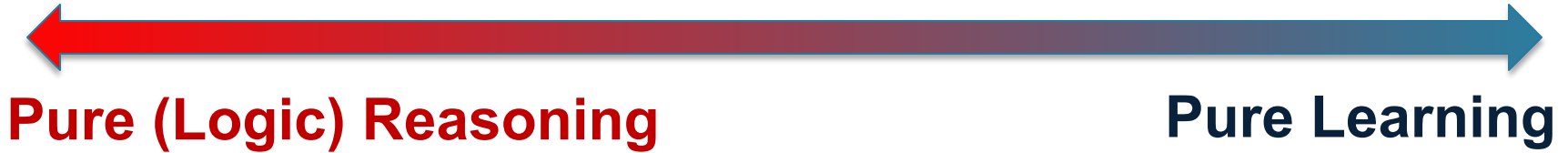
**Pure Learning**

- Fast thinking: instinctive, perceptive, model-free, interpolation
- Data-centered
- Amazing achievements recently
- “*Pure learning is brittle*”

bias, algorithmic fairness, interpretability, explainability, adversarial attacks, unknown unknowns, calibration, verification, missing features, missing labels, data efficiency, shift in distribution, general robustness and safety fails to incorporate a sensible model of the world



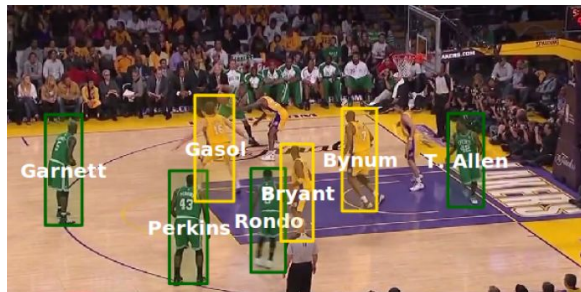
# The AI Dilemma



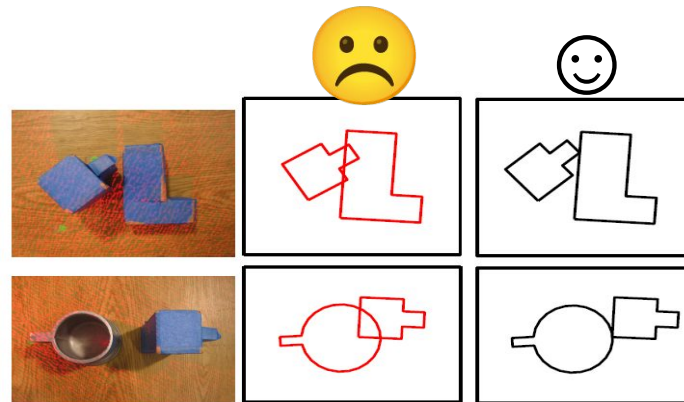
- Learn statistical models subject to logical knowledge
- Integrate reasoning into modern learning algorithms
- Reason about learned models' behavior
  - Algorithmic Fairness - Explainability

# Deep Learning with Output Constraints

# Knowledge in Vision, Robotics, NLP



People appear at most once in a frame

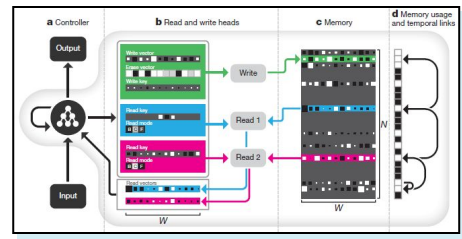
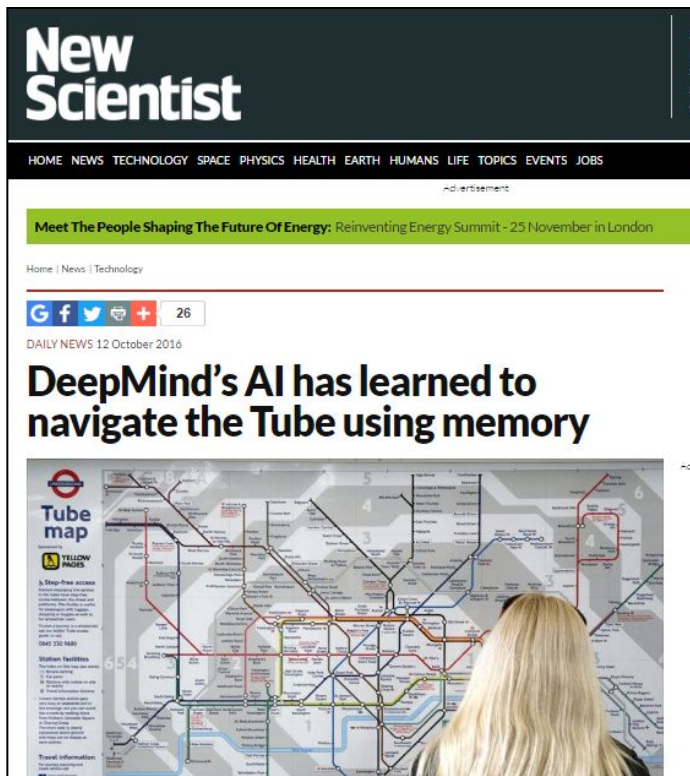


Rigid objects don't overlap

At least one verb in each sentence.  
If X and Y are married, then they are people.

[Lu, W. L., Ting, J. A., Little, J. J., & Murphy, K. P. (2013). Learning to track and identify players from broadcast sports videos.], [Wong, L. L., Kaelbling, L. P., & Lozano-Perez, T., Collision-free state estimation. ICRA 2012], [Chang, M., Ratinov, L., & Roth, D. (2008). Constraints as prior knowledge], [Ganchev, K., Gillenwater, J., & Taskar, B. (2010). Posterior regularization for structured latent variable models]... and many many more!

# Motivation: Deep Learning



[Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., et al.. (2016). Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626), 471-476.]



# Motivation: Deep Learning

DeepMind's latest technique uses external memory to solve tasks that require **logic** and reasoning — a step toward more human-like AI.

... but ...



optimal planner recalculating a shortest path to the end node. To ensure that the network always moved to a valid node, the output distribution was renormalized over the set of possible triples outgoing from the current node. The performance

it also received input triples during the answer phase, indicating the actions chosen on the previous time-step. This makes the problem a 'structured prediction'

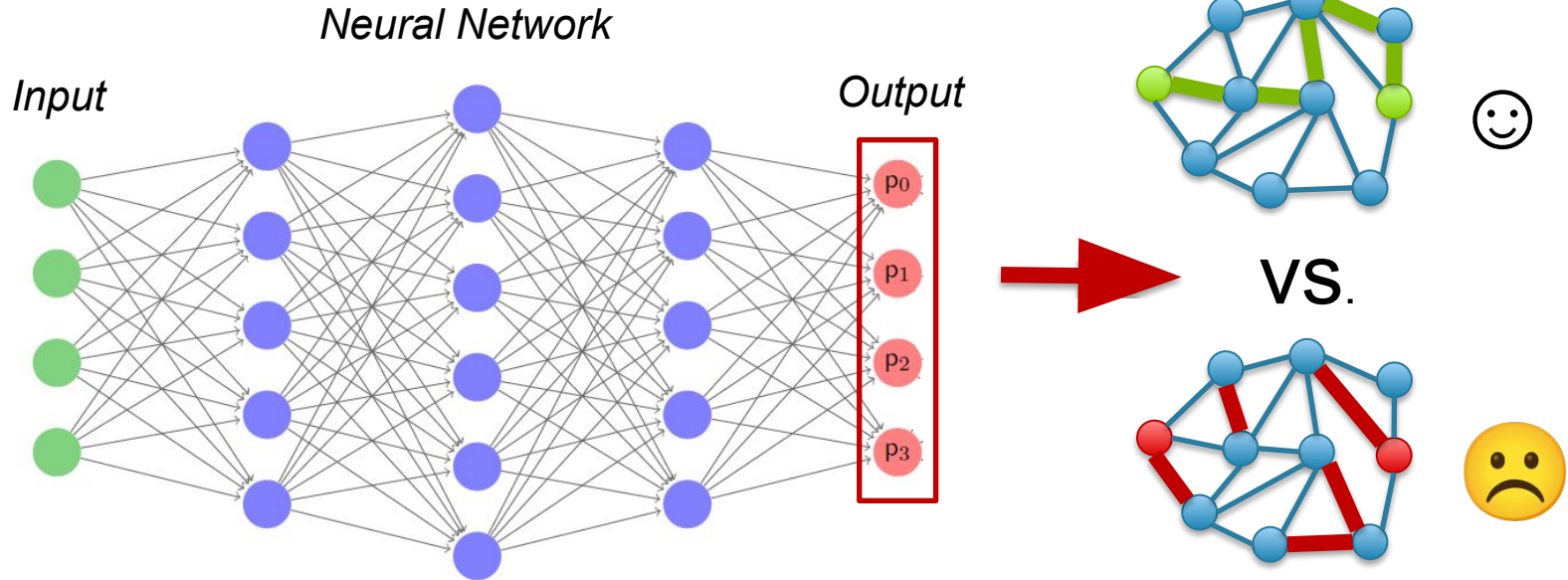
# Knowledge vs. Data

- Where did the world knowledge go?
  - Python scripts
    - Decode/encode cleverly
    - Fix inconsistent beliefs
  - Rule-based decision systems
  - Dataset design
  - “a big hack” (with author’s permission)

# Knowledge vs. Data

- Where did the world knowledge go?
  - Python scripts
    - Decode/encode cleverly
    - Fix inconsistent beliefs
  - Rule-based decision systems
  - Dataset design
  - “a big hack” (with author’s permission)
- In some sense we went backwards
  - Less principled, scientific, and intellectually satisfying ways of incorporating knowledge

# Deep Learning with Symbolic Knowledge



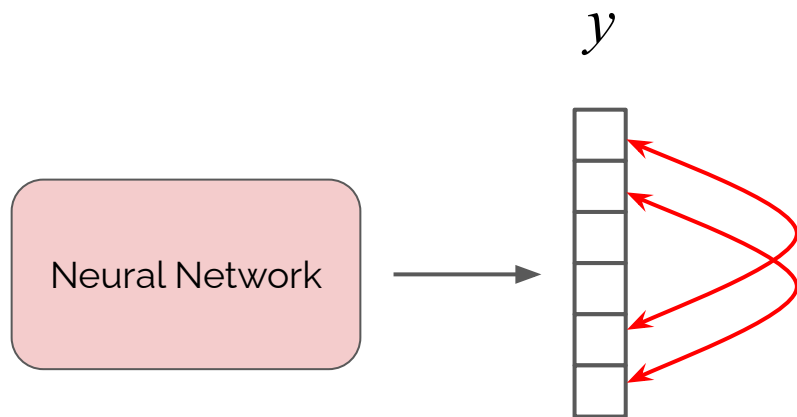
# pylon

A PyTorch Framework for Learning with Constraints

Kareem Ahmed   Tao Li   Thy Ton   Quan Guo,  
Kai-Wei Chang   Parisa Kordjamshidi   Vivek Srikumar  
Guy Van den Broeck   Sameer Singh

<http://pylon-lib.github.io>

# Declarative Knowledge of the Output



How is the output structured?  
Are all possible outputs valid?  
How are the outputs related to each other?

Learning this from data is inefficient  
Much easier to express this declaratively

How can do we inject declarative knowledge into PyTorch training code?

# pylon

Library that extends PyTorch to allow injection of declarative knowledge

- **Easy to Express Knowledge:** users write **arbitrary constraints** on the output
- **Integrates with PyTorch:** **minimal change** to existing code
- **Efficient Training:** compiles into loss that can be **efficiently optimized**

# pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)
```

1

Specify knowledge as a predicate

```
def check(y):  
    ...  
    return isValid
```



# pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)  
    loss += constraint_loss(check)(py)
```

1

Specify knowledge as a predicate

```
def check(y):  
    ...  
    return isValid
```

2

Add as loss to training

```
loss += constraint_loss(check)
```

# pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)  
    loss += constraint_loss(check)(py)
```

1 Specify knowledge as a predicate

```
def check(y):  
    ...  
    return isValid
```

2 Add as loss to training

```
loss += constraint_loss(check)
```

3 pylon derives the gradients  
(solves a combinatorial problem)

# Warcraft Shortest Path

Predicting the min-cost simple-path in a grid



# Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	<b>97.7</b>	56.9

*Is prediction  
the shortest path?*  
**This is the real task!**

*Are individual  
edge predictions  
correct?*

*Is output  
a path?*

## Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	<b>97.7</b>	56.9
+ Semantic loss	<b>50.9</b>	<b>97.7</b>	<b>67.4</b>

## Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	<b>97.7</b>	56.9
Semantic loss	<b>50.9</b>	<b>97.7</b>	<b>67.4</b>
+ Entropy All	51.5	97.6	67.7
+ Entropy Circuit	<b>55.0</b>	<b>97.9</b>	<b>69.8</b>

# pylon

- Joint entity-relation extraction in natural language processing
- Semantic role labeling in natural language processing
- Training MNIST recognition network from arithmetic supervision
- Training neural net to solve Sudoku
- Learning to rank
- etc.

# How do you get the loss function? What magic is this?

$$L^s(\alpha, \mathbf{p}) \propto -\log \underbrace{\sum_{\mathbf{x} \models \alpha} \prod_{i: \mathbf{x} \models X_i} p_i \prod_{i: \mathbf{x} \models \neg X_i} (1 - p_i)}_{\text{Probability of satisfying constraint } \alpha \text{ after sampling from neural net output layer } \mathbf{p}}$$

Probability of satisfying constraint  $\alpha$  after sampling from neural net output layer  $\mathbf{p}$

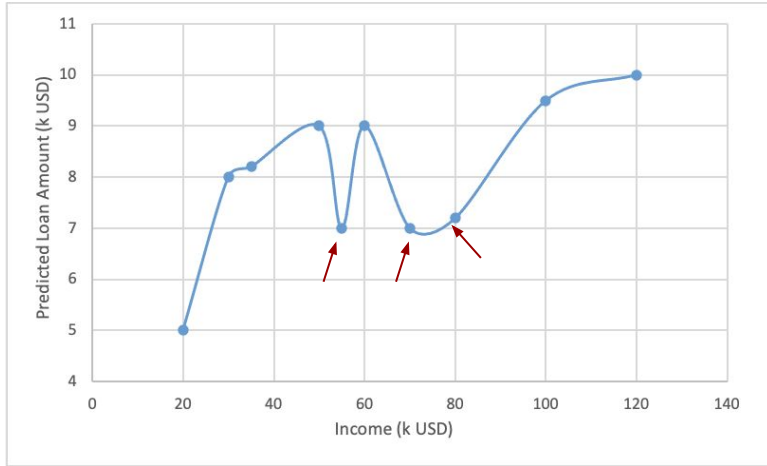
In general: #P-hard 😞

We do this probabilistic-logical reasoning during learning in a computation graph



# Monotonicity Invariants for Neural Networks

# Predict Loan Amount

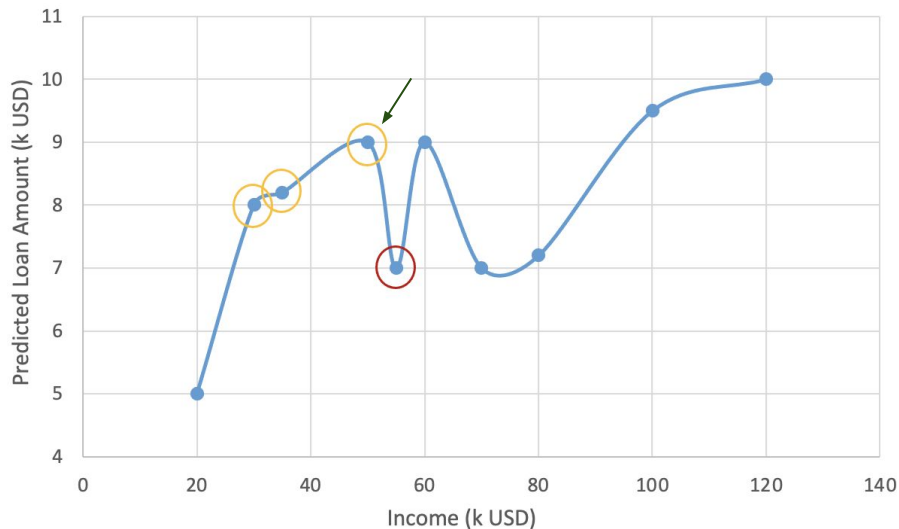


Neural Network Model: **Increasing income can decrease the approved loan amount**

Monotonicity (Prior Knowledge):

Increasing income should increase the approved loan amount

# Counterexamples

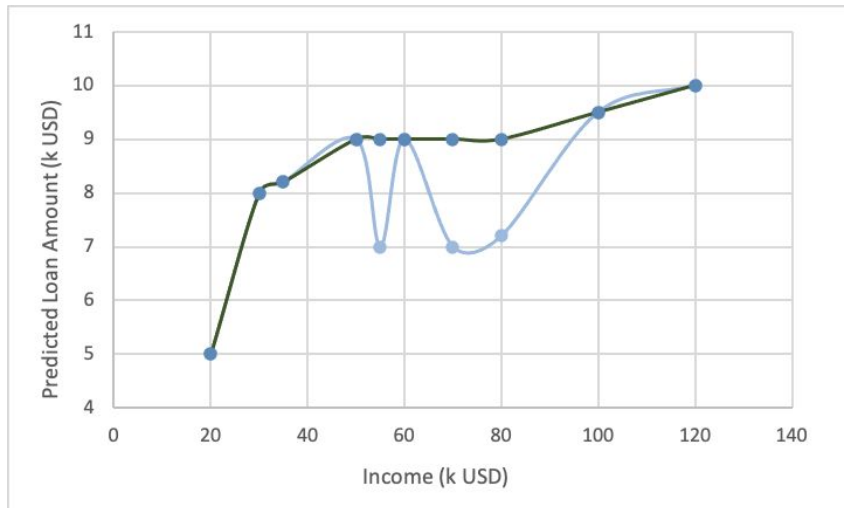


$$\exists x, y \ x \leq y \implies f(x) > f(y)$$

Computed using SMT(LRA)  
logical reasoning solver

Maximal counterexamples  
(largest violation) using OMT

# Counterexample-Guided Predictions



## Monotonic Envelope:

- Replace each prediction by its maximal counterexample
- Envelope construction is online (during prediction)
- Guarantees monotonic predictions for any ReLU neural net
- Works for high-dimensional input
- Works for multiple monotonic features

# Monotonic Envelope: Performance

Dataset	Feature	NN <sub>b</sub>	Envelope
Auto-MPG	Weight	9.33±3.22	<b>9.19±3.41</b>
	Displ.	9.33±3.22	9.63±2.61
	W,D	9.33±3.22	9.63±2.61
	W,D,HP	9.33±3.22	9.63±2.61
Boston	Rooms	14.37±2.4	<b>14.19±2.28</b>
	Crime	14.37±2.4	<b>14.02±2.17</b>

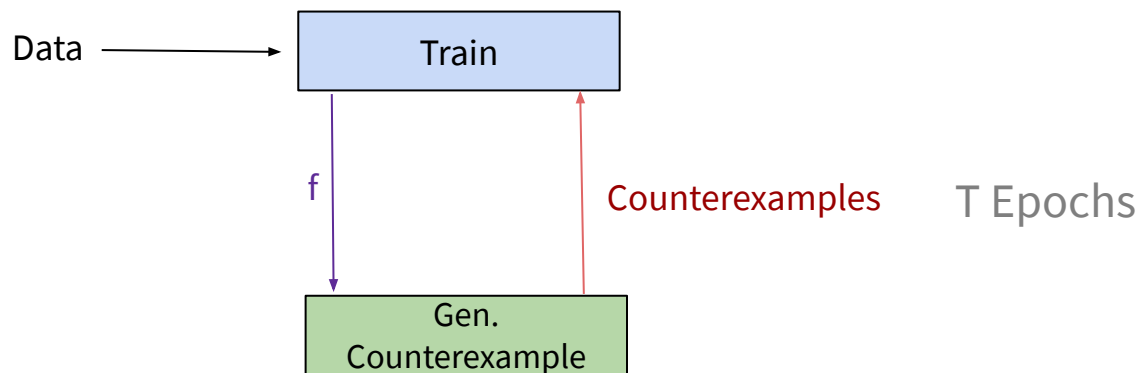
Dataset	Feature	NN <sub>b</sub>	Envelope
Heart	Trestbps	0.85±0.04	0.85±0.04
	Chol.	0.85±0.04	0.85±0.05
	T,C	0.85±0.04	0.85±0.05
Adult	Cap. Gain	0.84	0.84
	Hours	0.84	0.84

Guaranteed monotonicity at little to no cost

# Counterexample-Guided Learning

How to use monotonicity to improve model quality?

“Monotonicity as inductive bias”



# Counterexample-Guided Learning: Performance

Dataset	Feature	NN <sub>b</sub>	CGL
Auto-MPG	Weight	9.33±3.22	<b>9.04±2.76</b>
	Displ.	9.33±3.22	<b>9.08±2.87</b>
	W,D	9.33±3.22	<b>8.86±2.67</b>
	W,D,HP	9.33±3.22	<b>8.63±2.21</b>
Boston	Rooms	14.37±2.4	<b>12.24±2.87</b>
	Crime	14.37±2.4	<b>11.66±2.89</b>

Dataset	Feature	NN <sub>b</sub>	CGL
Heart	Trestbps	0.85±0.04	<b>0.86±0.02</b>
	Chol.	<b>0.85±0.04</b>	<b>0.85±0.05</b>
	T,C	0.85±0.04	<b>0.86±0.06</b>
Adult	Cap. Gain	<b>0.84</b>	<b>0.84</b>
	Hours	<b>0.84</b>	<b>0.84</b>

Monotonicity is a *great* inductive bias for learning

# Counterexample-Guided Monotonicity Enforced Training (COMET)

Table 4: Monotonicity is an effective inductive bias. COMET outperforms Min-Max networks on all datasets. COMET outperforms DLN in regression datasets and achieves similar results in classification datasets.

Dataset	Features	Min-Max	DLN	COMET
Auto-MPG	Weight	9.91±1.20	16.77±2.57	<b>8.92±2.93</b>
	Displ.	11.78±2.20	16.67±2.25	<b>9.11±2.25</b>
	W,D	11.60±0.54	16.56±2.27	<b>8.89±2.29</b>
	W,D,HP	10.14±1.54	13.34±2.42	<b>8.81±1.81</b>
Boston	Rooms	30.88±13.78	15.93±1.40	<b>11.54±2.55</b>
	Crime	25.89±2.47	12.06±1.44	<b>11.07±2.99</b>

Dataset	Features	Min-Max	DLN	COMET
Heart	Trestbps	0.75±0.04	0.85±0.02	<b>0.86±0.03</b>
	Chol.	0.75±0.04	0.85±0.04	<b>0.87±0.03</b>
	T,C	0.75±0.04	<b>0.86±0.02</b>	<b>0.86±0.03</b>
Adult	Cap. Gain	0.77	<b>0.84</b>	<b>0.84</b>
	Hours	0.73	<b>0.85</b>	0.84

COMET = Provable Guarantees + SotA Results



# Reasoning about the Feature Distribution

# Reasoning about World Model + Classifier

- “*Pure learning is brittle*”

bias, **algorithmic fairness**, interpretability, **explainability**, adversarial attacks, unknown unknowns, calibration, verification, **missing features**, missing labels, data efficiency, shift in distribution, general robustness and safety  
fails to incorporate a sensible model of the world



- Given a learned predictor  $F(x)$  over features  $x$
- Given a probabilistic world model  $P(x)$  - a feature distribution
- How does the world act on learned predictors?

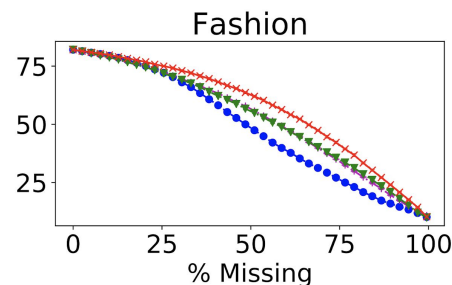
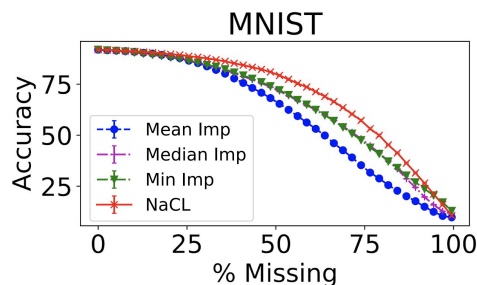
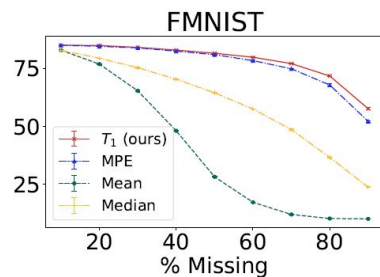
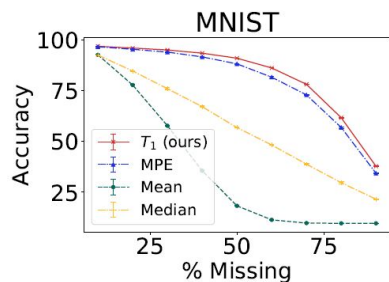
*Can we solve these hard problems?*

# What to expect of classifiers?

- Missing features at prediction time
- What is expected prediction of  $F(x)$  in  $P(x)$ ?

$$E_{\mathcal{F}, P}(\mathbf{y}) = \mathbb{E}_{\mathbf{m} \sim P(\mathbf{M}|\mathbf{y})} [\mathcal{F}(\mathbf{y}\mathbf{m})]$$

**M**: Missing features  
**y**: Observed Features



# Explaining classifiers on the world

*If the world looks like  $P(x)$ ,  
then what part of the data is **sufficient** for  
 $F(x)$  to make the prediction it makes?*

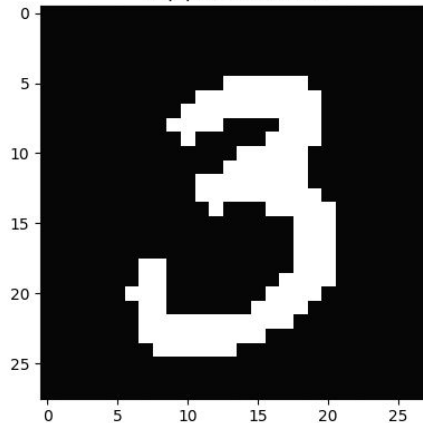
Probabilistic Sufficient  
Explanations

# Correctly Classified Examples

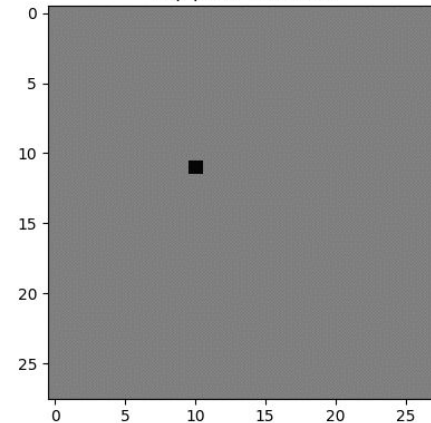
Binary classification: 3 vs 5

Used decision forest classifier and probabilistic circuit feature distribution

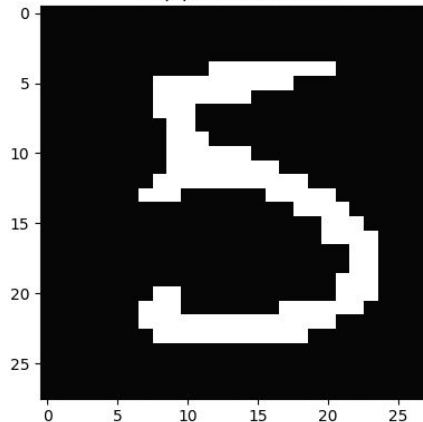
Actual label: 3(+)  
Exp pred: 3.758819



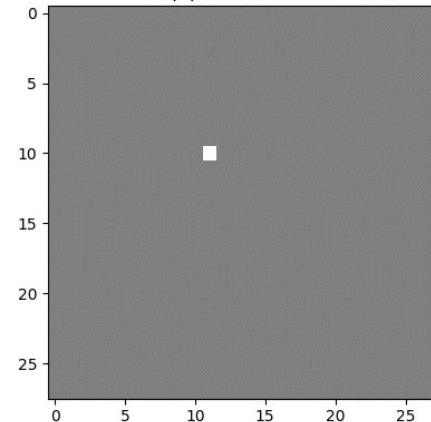
Max Features: 1  
Exp pred: 0.575232



Actual label: 5(-)  
Exp pred: -3.192585



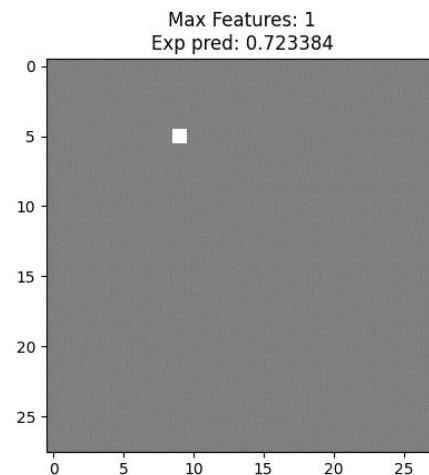
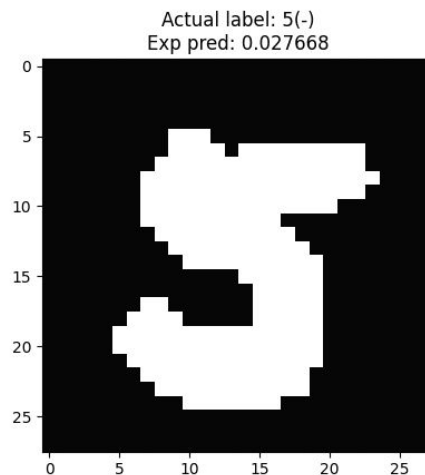
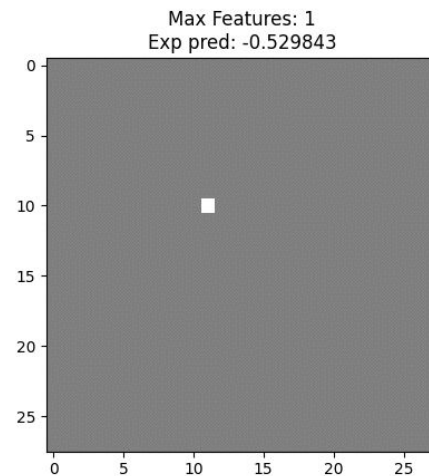
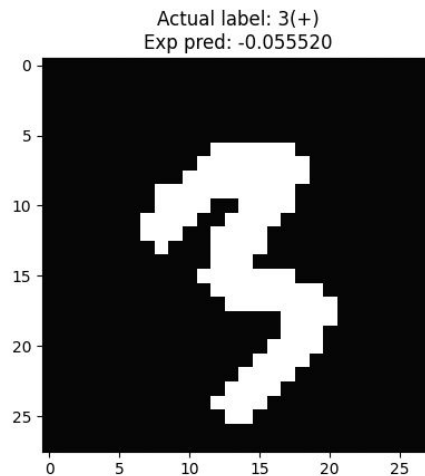
Max Features: 1  
Exp pred: -0.529843



# Misclassified Examples

Binary classification: 3 vs 5

Used decision forest classifier  
and probabilistic circuit feature  
distribution



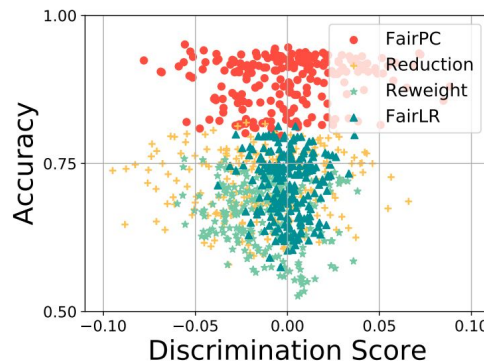
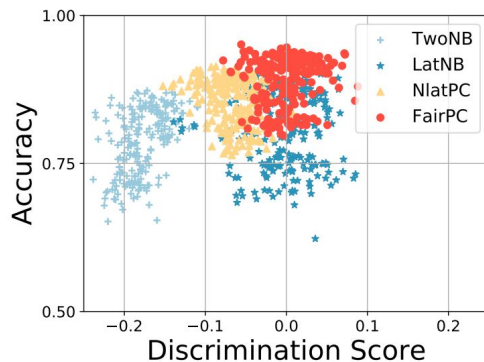
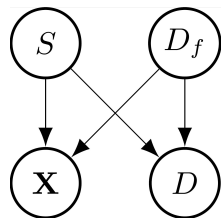
# Algorithmic Fairness: Latent Fair Decisions

Learn classifier given

- features  $S$  and  $X$
- training labels/decisions  $D$

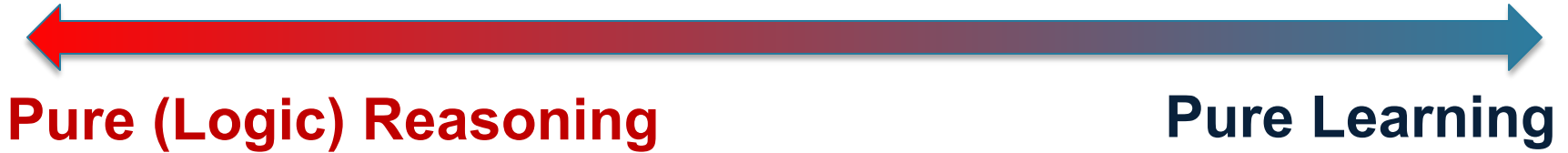
*Unknown fair decision  $D_f$  should be independent of the sensitive attribute  $S$*

Discover the **latent fair decision**  $D_f$  by learning distribution  $P(S, X, D_f, D)$  where  $D_f$  is 'fair'.



competitive  
*classification accuracy*  
and *better fairness guarantee*

# The AI Dilemma



- Learn statistical models subject to logical knowledge
- Integrate reasoning into modern learning algorithms
- Reason about learned models' behavior
  - Algorithmic Fairness - Explainability



# Thanks

*This was the work of many wonderful students/postdoc/collaborators!*

References: <http://starai.cs.ucla.edu/publications/>