

Beneficial and harmful explanatory machine learning

Lun Ai, Imperial College London

Challenges and issues

- Under-specified and ambiguous definitions
 - A lack of empirical data
 - Explanation effects are “subjective”
 - Limited references to social science literature
 - Little accounting for humans' perspective
 - Focus on exclusively on beneficiality
-

- 1) Operational measure of machine learning explanatory effect
 - 2) Cognitive window framework
 - 3) Demonstration of beneficial/harmful explanatory effect on human comprehension
-

Human out-of-sample predictive accuracy after studying training materials

Effect = machine-aided comprehension - self-learning comprehension

Beneficial = positive effect

Harmful = negative effect

Otherwise, no observable effect

Bound on hypothesis space size (**learning**):

Effect is negative when $|S| > B(M(E), H)$

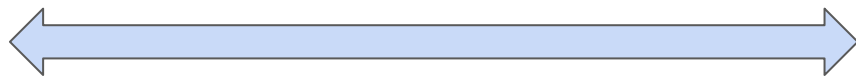
Kolmogorov complexity bound on execution cost (**runtime**):

Effect is non-positive when $Cog(M(E), x) \geq CogP(E, \bar{M}, \phi, x)$

Noughts and Crosses isomorphism:

- Avoid ceiling effect
 - Same complexity
 - 2-ply initialisation
-

Easiest



Hardest

Last move to win as Player X

(win_1)

X		X
O	X	
O		O

Two moves to win as Player X

(win_2)

X		
O	X	
		O

Three moves to win as Player X

(win_3)

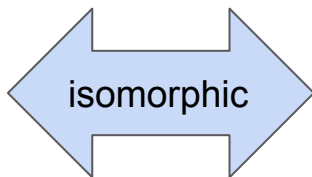
O	X	

Noughts and Crosses

X to move

Win - player has three pieces in a line

X		X
O	X	
O		O

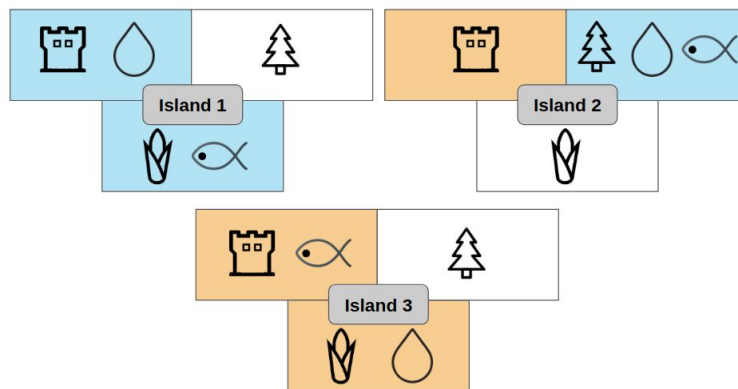


The Island Game

Blue to move

Win - 1) player has three cells on the same island

2) player has three of the same resource



MIGO:

Sufficient and necessary BK

Positive examples only

Learns minimax algorithm

MIPlain (adapted MIGO with Metaopt):

Additional BK

Positive and negative examples

Learns smaller program with less inferential cost

Depth	Rules
1	<code>win_1(A,B):-move(A,B),won(B)</code>
2	<code>win_2(A,B):-move(A,B),win_2_1(B)</code> <code>win_2_1(A):-number_of_pairs(A,x,2), number_of_pairs(A,o,0)</code>
3	<code>win_3(A,B):-move(A,B),win_3_1(B)</code> <code>win_3_1(A):-number_of_pairs(A,x,1),win_3_2(A)</code> <code>win_3_2(A):-move(A,B),win_3_3(B)</code> <code>win_3_3(A):-number_of_pairs(A,x,0),win_3_4(A)</code> <code>win_3_4(A):-win_2(A,B),win_2_1(B)</code>

Control group :

pre test => training with examples => post test

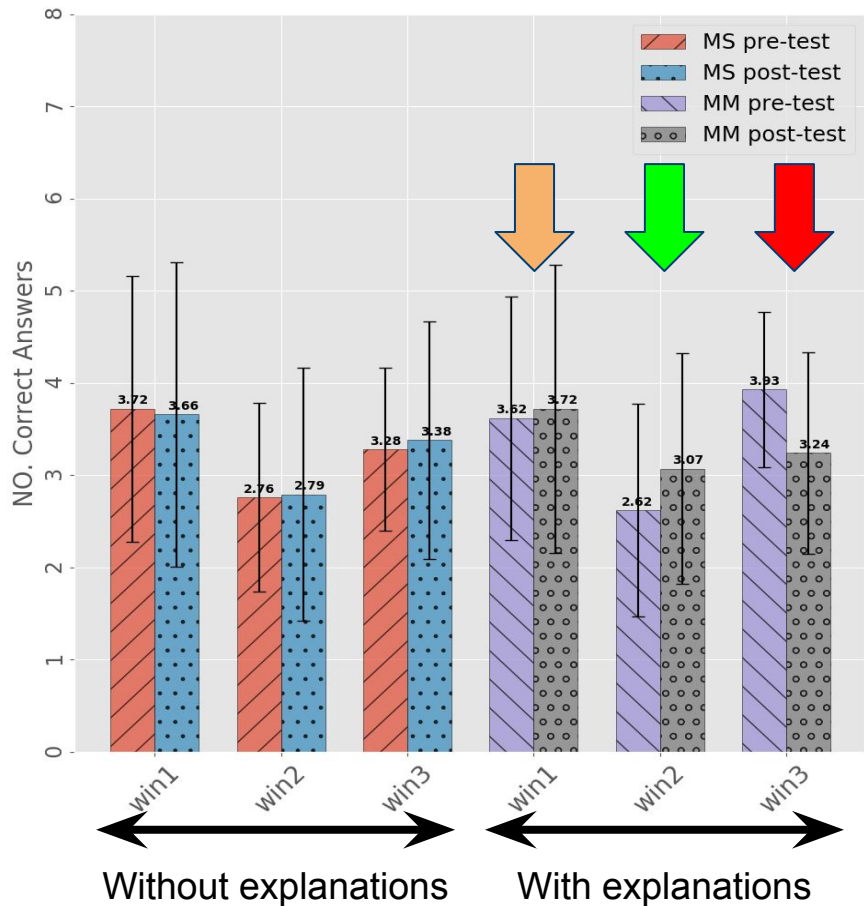
Treatment group:

pre test => training with examples & explanations => post test

English translation of MIPlain's learned theory

Visualisation of game states

Contrastive consequences of good and bad moves



	Orange	Green	Red
Effect	No observable	Beneficial	Harmful

- 1) Bound on the program size learnable
 - 2) Short-term comprehension improved by heuristics
 - 3) Evidence for a *cognitive window*
-

How does the ordering of concepts affect comprehension?



How does machine-aided knowledge corrections affect comprehension?

Thank you

Bound on hypothesis space size (**learning**):

Effect is negative when $|S| > B(M(E), H)$

Kolmogorov complexity bound on execution cost (**runtime**):

Effect is non-positive when $Cog(M(E), x) \geq CogP(E, \bar{M}, \phi, x)$

A balance between memory and computational complexity

More training examples improves predictive accuracy

Symbolic output

Humans achieve better performance via teaching