*Explainable Artificial Intelligence*
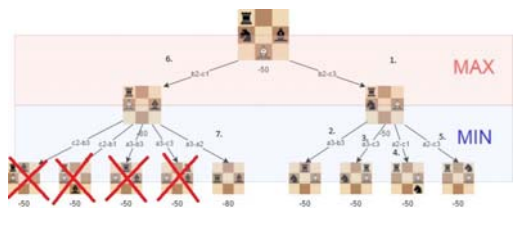
Johan Kwisthout, DCC / AI

---

Thank you for an enjoyable game!



---

Minimax with alpha-beta pruning



---

I'm sorry, Frank, I think you missed it

- Frank missed, that the minimax algorithm returns -1, indicating a winning path for black

- Frank missed, that $\exists m_B \forall m_W \exists m_B \forall m_W$ ... black wins

- In the (science fiction) movie, HAL **explained** to Frank why he played a sure loss by **exemplifying** a 'natural' sequence of moves that made it **obvious** to **him** that Frank's position was hopeless

- (There's a whole lot more going on in this scene that is not relevant for this talk…)

---

Explanation in MYCIN expert system



- Diagnosis of infections and suggested antibiotics

- Developed in 1972 (!)

- Set of 600 rules, reasoning with uncertainty

- Sometimes outperforming human diagnoses

- Can explain and answer 'why' and 'how' questions

---

Explanation in MYCIN expert system

## Explanation in MYCIN expert system

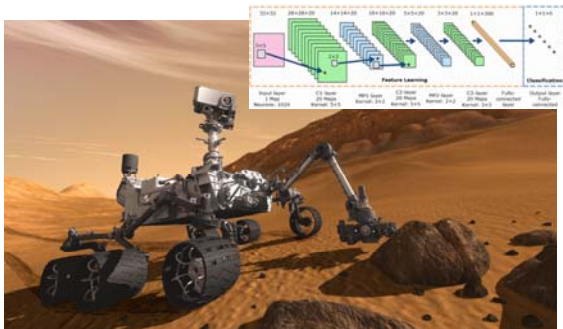- Maybe not so sophisticated as it looks…

```
(defun print-why (rule parm)
  "Tell why this rule is being used.  Print what is known,
  what we are trying to find out, and what we can conclude."
  (format t "~&[Why is the value of ~a being asked for?]" parm)
  (if (member rule '(initial goal))
      (format t "~&~a is one of the ~a parameters."
              parm rule)
      (multiple-value-bind (knowns unknowns)
          (partition-if #'(lambda (premise)
                            (true-p (eval-condition premise nil)))
                        (rule-premises rule))
        (when knowns
          (format t "~&It is known that:")
          (print-conditions knowns)
          (format t "~&Therefore,"))
        (let ((new-rule (copy-rule rule)))
          (setf (rule-premises new-rule) unknowns)
          (print new-rule)))))
```

https://norvig.com/paip/mycin.lisp

---

## MYCIN revisited

- MYCIN was a **rule-based** expert system within a very **specific domain**, all rules **hand-coded** based on expert knowledge elicitation

- Explanation in this system useful, but mechanistic

- In modern AI, information is mostly *machine learned* by discovering statistical patterns in data

- How can we let such systems explain their decisions to us?

---

## Curiosity Mars Rover



---

## Explanations from Curiosity

- Under the hood, autonomous vehicles might base their decisions on (e.g.) deep neural networks

- Sub-symbolic AI: information decoded in weights between artificial neurons

- When asked 'why did you make that decision?' we don't want an answer like:

"because $f\left(\sum_{k=1}^{n} i_k \cdot W_k\right) > 0$"!

---

## Gunning (DARPA): XAI project



- Train the net to associate semantic attributes with hidden layer nodes and known ontologies

---

## Generating examples

**Learning to Generate Chairs with Convolutional Neural Networks**

Alexey Dosovitskiy    Jost Tobias Springenberg    Thomas Brox
Department of Computer Science, University of Freiburg
{dosovits, springj, brox}@cs.uni-freiburg.de

## Why is this a cat?

- Specific nodes in a deep layer of the neural network capture statistical abstractions – that may or may not correspond to salient features we recognize…



## Pattern recognition ≠ explanation!

- Specific nodes in a deep layer of the neural network capture statistical abstractions – that may or may not correspond to salient features we recognize…

- Even granted that, and granted that we would provide the learning algorithm with a full ontology, just listing (selected) hidden nodes that are activated on an input is not so useful…

- **Inference to the best explanation**: what is the **cause** (or reason, principle, etc.) of the phenomena I can **observe**? Also called **abduction**

## Deduction, induction, and abduction

- **Deduction**: from premises and a rule deduce a logically valid conclusion
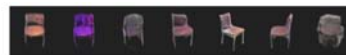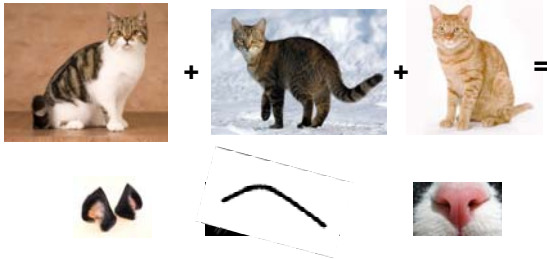
$$(\forall x H(x) \rightarrow M(x), H(s) \vDash M(s))$$

- **Induction**: from specific observations derive general principles

  (all biological life forms we know depend on liquid water to exist, so it is likely that there cannot be life without water)

- **Abduction**: from an observation derive a cause that best explains the observation

  (When I leave my house in the morning the grass is wet; given weather forecasts the best explanation is that it rained tonight)

## Statistical association ≠ explanation!

- From data alone we cannot learn causal relations ("correlation is not causation")

- Judea Pearl (2017): vital for AI to be able to ask questions "what if I do X" (intervention) and "what I had done Y" (counterfactuals)



*Pearl at NIPS 2017*

- These questions are also important for Explainable AI!

## What is a good explanation?

- Explanation not just answers "**why this**", but "why this, **rather than** that" (parsimoniously)

- Q "Why did Alice got tenure (while Bob didn't)?"

- A1 "Alice had a good publication record"
  - But Bob had a good publication record as well! That doesn't explain why she got tenure!

- A2: "Alice had a good publication record and did quality teaching"
  - Bob was a poor teacher, so this explains why Alice got tenure and Bob was denied tenure!

## What is a good explanation?

- Explanation must be based on *relevant* information

- Q "Why did Alice got tenure (while Bob didn't)?"

- A3 "Alice had a good publication record and wore glasses"
  - But Bob had a good publication record as well! And wearing glasses ought to be irrelevant for getting tenure! That doesn't explain why she got tenure!

- But how do we decide that wearing glasses is not relevant *even if this might be statistically significant?*

## Abduction ingredients

- Abduction is **making sense** of your observations in order to act accordingly and **motivate** your actions
- **Generate** candidate hypotheses that might explain the phenomena that have been observed
- Decide what is **relevant** in its sensory input and what is not
- Decide when to gather **new evidence** (e.g., re-orient your sensors, do additional tests) to reduce uncertainty
- From a set of candidate hypotheses, **select the best** one
- In a given context, **determine** what constitutes 'best'
- Try to infer **causal relationships** and test hypotheses by **interventions**, i.e., acting in the world
- Generate and reason through '**what if**' scenarios

## Important



- Some or all of these 'ingredients' might be 'implemented' in modern sub-symbolic models (e.g. convolutional deep neural networks)

- Yet, for Explainable AI the challenges are the same as for (symbolic) GOFAI and philosophy of mind (e.g., symbol grounding, frame problem)

- It is one thing for AlphaGO to beat the world champion; explaining the rationale behind its moves in a way we can understand is a different ball game

## Sherlock Holmes without Watson is 'useless'



## Explainable AI as research method



- "What I cannot create, I do not understand"

- Try to implement a theory in order to identify its ambiguities, test consistency and completeness, and identify gaps

- Learn about human cognition by trying to implement your favorite theory of human sense-making in AI

- See also: Otworowska et al. (BNAIC 2015). "The Robo-havioral Methodology"

## Conclusion and summary

- Explainable AI becomes more and more important particularly when the AI becomes a 'black box' as in deep neural networks and machine learning

- The challenges in creating Explainable AI are similar to the challenges in understanding human sense-making (e.g. frame problem)

- Apart from a research goal on its own, Explainable AI can be a research method to put theories of human sense-making to the test