





Introduction

A primary goal of science is to provide explanations of phenomena, be they natural, social, or otherwise (Many other goals: forecasting, diagnosis, technology development). Ideally, the choice of scientific actions should rely on some road map for how to make progress towards at least one scientific goal. When the goal is to explain, this road map requires some idea of what type of explanation we're looking for.

There is little agreement or formalization of what would be acceptable explanations of the phenomena of intelligence. Panelists at CCN2017 gave varied answers to the question "what is your definition of success?" The quest for interpretable machine learning is ultimately a similar question, what explanations of AI systems will we accept?

Claim: The integration of deep learning and neuroscience will require a common theory of explanation that applies to both artificial and biological intelligence

Goal: Equip scientists of intelligence to interrogate and justify the theories of explanation that underlie their definitions of scientific progress

The role of a theory of scientific explanation

Characterize the structure of explanations in science

• Distinguish between explanations that are scientific and those that are

• Distinguish between explanations and non-explanations. Sometimes presented as the difference between explanation and description. For example, a set of claims about the appearance of a particular species may be true, accurate and supported by evidence without being explanatory in any way. They are "merely" descriptive.

Reveal criteria for assessing explanations

Phenomenological models (aka "descriptive models"): describe or "save" a phenomena, summarize data compactly

Descriptive and normative goals of this field of philosophy:

- Descriptive: characterize explanations in contemporary science
- Normative: clarify the distinction between good and bad explanations

Particular challenges to explaining intelligence

The form of a "good" explanation of a given phenomenon depends on the nature of the phenomenon itself. Therefore, we cannot separate the ontological question (what is intelligence? what is cognition? what is computation?) from the epistemological question (how to explain intelligence? how to explain neural function?)

Philosophical theories of mind:

 Classical computational theory of mind

- Behaviourism
- Type-identity theory
- Functionalism
- Representational theory of mind
- Connectionism
- Embodied Dynamicism

Philosophical theories of computation:

- Formal
- Mechanistic
- Information processing
- Modeling

The Deductive-Nomole

(aka Covering Law Model) A scientific explanation is "law of nature" that shows explained) is to be expect thing that does the expla "law" is used to different generalizations that are Predictivism: "any pred force" (Kaplan 2011)

The Statistical Releva Explanations must cite c captured by statistical re

P(Pregnancy|Male.Tak P(Pregnancy|Female.7

Functional Explanation Phenomena are explained

organization of their sub complex capacity into si reductionist manner until stages: analysis and inst

The Unification Model

explanation is "a matter argument patterns" "Science advances our ι derive descriptions of m derivation again and aga to reduce the number of ultimate." (Kitcher 1989

The Causal Mechanica

Constitutive (or compone mechanisms underlying a phenomenon---the physical entities and activities organized such that they exhibit the explanadum---and thereby shows how a phenomenon is produced by its causes. Organization is important: not just the sum of parts but their causal interaction.

- Aspects of mechanistic explanation:
- 1. the nature of the phenomenon to be explained
- 2. the constitutive relationship between a phenomenon and its components 3. the difference between real components and useful fictions
- 4. the nature of capacities or activities
- 5. the nature of mechanistic organization
- 6. the nature of constitutive explanatory relevance

Simulation as explanation phenomenon.

Towards a common philosophy of explanation for artificial and biological intelligence

Jessica A.F. Thompson

International Laboratory for Brain, Music and Sound Research Montréal, Canada

Montreal Institute for Learning Algorithms (Mila) Montreal, Canada

Theories of Scientific Explanation and Their Weakne	
Theories	• Criticisms and
bgical (DN) Model (Hempel 1959, 1965) I) s a deductive argument based on at least one vs that the explanadum (the phenomenon to be cted given the premises of the explanans (the aning.) ntiate deterministic laws from other true only "accidentally true" dictively adequate model possesses explanatory	 <i>Explanatory asymmetries</i>: some explanations are directional. (e.g. length shadow cast by a flagpole) <i>Explanatory Irrelevancies</i>: A derivation nor relying on a true generalization that is in (L) All males who take birth control pills (K) John Jones is a male who has been tagen (E) John Jones fails to get pregnant
nce (SR) model (Salmon 1971) ausal relationships and causal relationships are elevance relationships, i.e. conditional dependence es birth control pills) = P(Pregnancy Male) = 0 Takes birth control pills) \neq P(Pregnancy Female)	 Objective homogeneity condition: "the variables that would affect the probabilit Casual relationships are greatly underconstruction (Cartwight, 1979 and Spirte 2000.)
n (Cummins 1975, 1983) ed by reference to the functional role and components. This process of decomposing a mpler subcapacities can be repeated in a I the subcomponents are well-understood. Two antiation.	 Most easily applied when assuming fur theory of mind and computation. Under obvious that subcapacities should explai Does not establish causal relevance of In practice, scientists often present fur Computational chauvinism: in line with explained independently of the physical
(Kitcher 1989) of unifying diverse beliefs under a few simple nderstanding of nature by showing us how to any phenomena, using the same patterns of ain, and in demonstrating this, it teaches us how facts that we have to accept as 423)	 The most unifying statements are not a taxonomies) No notion of causality
I (CM) Model (Salmon, 1984, 1998) ntial) mechanistic explanation reveals the	Not clear how the CM model would apply phenomena, e.g. cognition.

Scientists who develop simulations often claim that simulations represent or demonstrate progress towards the goal of explaining the simulated

Philosophers seems to agree that simulations are, in general, not explanatory. From a causal perspective, identifying one "how possibly" mechanism out of potentially infinitely many, does not in itself constitute progress towards finding the true mechanism. However, Grune-Yanoff (2009) suggests that simulations may be seen as candidate functional explanations.

Resting state

period /

1 2 3 4 5

5 Threshold

Centre for Research on Brain, Music and Language (CRBLM) Montréal, Canada

sses

d Examples



nay satisfy the DN model, while relevant to the explanadum. E.g. s regularly fail to get pregnant aking birth control pills regularly

re are no additional omitted

determined by statistical relevance es, Glymour and Scheines, 1993,

nctionalist/information processing other conceptions of mind, it is not in supercapacities.

subcapacities.

nctional analysis as explanatory. functionalism, cognition can be implementation in the brain

necessarily explanatory (e.g.

y to higher-level, more abstract

Physical components may be irrelevant to explain some phenomena



Three Perspectives in Cog Comp Neuro

What neuroscientists say

What philosophers say

Ilya Nemenman: "It doesn't matter if it's true" theories/models of complex Such models are certainly

Good biological systems: explanatory?

- are phenomenological
- make accurate predictions
- throw away unnecessary details
- are as simple as possible (Occam's Razor)
- explain a limited set oh phenomena
- are falsifiable but not falsified, according
- to Bayesian statistics

The Balmer Formula

 $\frac{1}{\lambda} = R_H \left(\frac{1}{n^2} - \frac{1}{n^2} \right)$

Predictivism



Kendrick Kay: Explain via functional analysis

"Mental operations [...] can be viewed as • information processing operations. The explains cognitive phenomena cognitive neuroscientist asks: for a given via functional analysis of the brain region, what stimulus, cognitive, or brain motor operations are performed by neurons in that region? [...] Models posit that theory of mind specific variables relate to neural activity. As such, models provide explanations of subcomponent is indicative of measurements of the brain [...] With its role in the capacity to be appropriate experimental measurements, explained. we can adjudicate different models and decide which model is most accurate"

To explain, compare interpretable models with relatively small number of parameters. Models are falsified to the extent that they fail to make accurate predictions.

- Cognitive neuroscience
- Information processing Selective response

Jonas Kubilis: "Predict then simplify"

The absurdity of Occam's Razor:

"how could a fixed bias toward simplicity indicate the possibly complex truth any better than a broken thermometer that always reads zero can indicate temperature? " (Kelly, 2007)

 Build deep network-based models that predict neural activity as well as possible for the broadest set of experiments/stimuli Then narrow to small number principles rather than parameters that are integral to the system

• As our models demonstrate increasingly realistic behaviours and mimic neural representations with increasing fidelity, we'll understand the system better

- Predictivism
- Unification model
- Simulation

Abandons commitment to information processing theory of mind (or any theory of References mind?)

 Imprecise about how turning the crank between better models and more naturalistic experiments will ultimately lead to unifying explanation





j.thompson@umontreal.ca

Constraints on explanations (according to Craver 2007)

- Mere temporal sequences are not explanations
- Causes explain effects and not vice versa
- Causally independent effects of common causes to not explain one another
- Causally irrelevant phenomena are not explanatory
- Causes need not make effects probably to explain them

Artificial neuroscience and biological machine learning

Swapping methodologies, approaches and philosophies has the potential to demonstrate the strengths and limitations of our scientific activities and perspectives, helping us to select those that will be most useful towards our common goal of understanding intelligence.

• Empirical analyses of deep learning systems that seem almost neuroscientific (e.g. ablation analyses, receptive field analysis, psychophysics) but with reference to concepts from deep learning theory (e.g. generalization, expressivity).

Machine learning inspired neuroscience

Desiderata for a new theory of explanation for both artificial and biological intelligence:

- Reflect that learning is central to intelligence
- Multiple realizability without computational chauvinism
- Abandon focus on physical computation

• Not concerned with characterizing the specific function that is computed by a network

- 1. because we probably can't glean anything meaningful from that anyway (see "Can neural computation be compressed enough for us to understand it?" Lilicrap & Kording, this meeting)
- 2. because we know that repeated optimizations of the same network lead to solutions that occupy distinct regions in function space, many local minima (Erhan 2010)

• Take the good parts from existing theories: causality (CM model), unifying principles (unification model), abstraction (functional)

Conclusion

• If the nature of a good scientific explanation is dependent on the phenomenon to be explained, then, to the extent that an artificial system and a biological system demonstrate the same phenomenon, their explanations should share the same form.

• Let's be precise with our use of the words "explanation", "mechanism", "description". Not a value judgment—description is important too!

• Room for increased clarity about what phenomena we are studying, what our short term and long term goals are, and how our short term goals will serve the ultimate goal of explaining the phenomenon in question. A we studying visual object recognition or the neural activity elicited during visual object recognition?

[1] Nemenman, I. (2018). Playing Newton: Automatic Construction of Phenomenological, Data-Driven Theories and Models. https://simons.berkeley.edu/talks/ilya-nemenman-4-17-18. [2] Kaplan, D. M. (2011). Explanation and description in computational neuroscience. Synthese, 183

[3] Kay, K. N. (2017). Principles for models of neural information processing. NeuroImage. [4] Kubilius, J. (2017). Predict, then simplify. NeuroImage.

[5] Woodward, J. (2017). Scientific Explanation. In E. N. Zalta (Ed.), The stanford encyclopedia of philosophy (https://pled.). Metaphysics Research Lab, Stanford University.

[6] Kelly, K. T. (2007). Simplicity, Truth, and the Unending Game of Science. In S. Bold, Benedikt Lowe, T. Rasch, & J. van Benthem (Eds.), Foundations of the formal sciences v: Infinite games [7] Grüne-Yanoff, T., & Weirich, P. (2010). The philosophy and epistemology of simulation: A review. Simulation and Gaming, 41(1), 20–50.

useful, but are

Statistical relevance mode

Looking Forward