
AI as statistical methods for imperfect theories

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Science has progressed by reasoning on what models could not predict because
2 they were missing important ingredients. And yet without correct models, standard
3 statistical methods for scientific evidence are not sound. Here, I argue that machine-
4 learning methodology provides solutions to ground reasoning about empirically
5 evidence more on models' predictions, and less on their ingredients.

6 Science uses false models as means for truer theory [Wimsatt, 1987]. How can statistical tools
7 ground valid reasoning on empirical evidence without true models? Generalization is the key. Here
8 I develop the argument that, unlike popular belief, reasoning from black-box models is good for
9 science, because it builds the validity of inferences on prediction of observables.

10 **1 Science has progressed by refining relevant constructs from wrong models**

11 **1.1 Observing motions of bodies, working out laws of physics**

12 Early scientists, such as Aristotle, did not conceive mechanics in terms of acceleration and forces.
13 Rather, they thought in terms of natural motion of objects, proportional to their weight. The notion
14 of force made its way, as discussed by Ibn Sīnā, but motion was seen as proportional to external
15 forces. The Copernican revolution motivated the importance of acceleration. Increasingly precise
16 astronomical observations led to formulate planetary motion as elliptical trajectories. Scientists such
17 as Kepler were seeking simple phenomenological rules, “harmonies” in his words, to explain the
18 observations, *eg* that across the different planets the square of the period is proportional to the cube of
19 the major diameter of the orbit. By introducing acceleration via differential calculus, Newton could
20 propose laws of mechanics that explained observations of both celestial and earthly motion.

21 The birth of Newtonian mechanics illustrates how better *observations* and *statistical models* lead
22 to better theories, even when starting *without the right theoretical framework*. It shows how *new*
23 *ingredients* may be needed, such as introducing the construct of acceleration. It shows that progress
24 is driven by seeking theories that *generalize* across many settings. The importance of acceleration
25 was revealed by uniting motion of bodies on Earth and in astronomy. Indeed, as friction is ubiquitous
26 on Earth, applying a force to an object often leads to a velocity roughly proportional to this force.

27 Later, better observations called for new frameworks, quantum or relativistic. Irregularities in the
28 orbit of Mercury were first explained by adding a planet to the solar system, Vulcan. But observations
29 of this planet turned out to be flawed, and the irregularities in Mercury's orbit are now understood
30 as relativistic corrections. The Vulcan hypothesis illustrates how theoretical frameworks shape the
31 interpretations of empirical results: observations are “theory laden” [Boyd and Bogen, 2021].

32 Today, the fundamental laws of physics are incredible precise. Are phenomenological models still
33 important for their empirical validation? From a statistical perspective, the Neyman-Pearson lemma

34 tells us that the optimal way to compare models is to use their likelihood [van Dyk, 2014]. Indeed,
35 particle physics has long polished probabilistic models, minute stochastic description of observations
36 built from first principles [Sjöstrand et al., 2001, Aaltonen et al., 2008]. And yet, recent statistical
37 analysis of Higgs bosons is powered by black-box machine learning models –such as boosted decision
38 trees– as they capture best background sensor noise [Aaltonen et al., 2009, Radovic et al., 2018].

39 1.2 Cognitive neuroscience: uncovering the functional units of human vision

40 Cognitive neuroscience strives to explain cognitive functions from neural activity. Which ingredients
41 to include in such a model is a more open-ended question in than in physics. Breaking down high-
42 level functions into units of investigation is particularly challenging. This endeavor has made much
43 progress for the specific problem of vision. Studying early visual cortex response to specially-crafted
44 stimuli, Hubel and Wiesel [1959] revealed neurons that form localized edge detectors. Slightly more
45 complex shapes isolated other brain units [Logothetis et al., 1995]. These findings are tied to the
46 stimuli presented, themselves motivated by cognitive theories used to decompose mental processes.
47 Theories of visual processing break down it into successive operations tuned to specific aspects of
48 the stimuli [Marr, 1982]. As any cognitive theory, their empirical neuroscience validation is then
49 bound to this decomposition. Even with modern neural measurements, a decomposition into invalid
50 ingredients, such as “alimentiveness” or “philoprogenitiveness” of 19th century phrenology, would
51 lead to a brain mapping valid from the statistical standpoint [Poldrack, 2010].

52 Complete models of cortical visual processing assemble brain functional units, each implementing
53 specific operations [Riesenhuber and Poggio, 1999]. They derive from many studies of neural
54 responses to elementary manipulations of visual stimuli. But their neuroscience validity faced a
55 chicken-and-egg problem as long as each functional unit had been studied in isolation: each study
56 had investigated only one aspect of otherwise very complex stimuli, natural images. Models of vision
57 can be derived without invoking neuroscience arguments, as in computer vision where computational
58 models are optimized directly on natural images, *eg* for object recognition [Pinto et al., 2009, Sermanet
59 et al., 2014]. In fact, *encoding* studies showed that pure computational models explain better neural
60 activity than models based on hand-crafted reductions of natural images [Yamins et al., 2014]. These
61 computer-vision models, based on artificial neural networks, extract intermediate representations of
62 natural images, which can be mapped to brain responses, confirming functional units obtained in
63 more hypothesis-laden neuroscience experiments [Eickenberg et al., 2017].

64 The large computational models do not answer some cognitive-neuroscience debates, such as the
65 specific semantic tuning of functional areas. For instance, a brain area crucial to recognizing human
66 faces is known as the *fusiform face area* [Kanwisher et al., 1997]. Yet, some researchers claim that
67 its role is best explained by implementing visual expertise, rather than face recognition [Tarr and
68 Gauthier, 2000]. As the corresponding brain area responds to both types of stimuli, the debate became
69 trapped in a ontological disagreement: which of visual expertise or face recognition is a more central
70 mental function? One side argues visual expertise leads to face recognition, and the other that face
71 recognition is innate to the social human.

72 Encoding studies use as ingredients to map brain responses the internals of large computational
73 models of vision. As such, they circumvent questions related to finding valid ontologies of cognitive
74 processes: on the one hand, they cannot bring evidence in favor of ontological choices, but on
75 the other hand they enable empirical evidence without buying into one framework. There are two
76 ingredients to this robustness. First, encoding studies can work on more ecological and richer stimuli.
77 Hence they capture all facets of cognition, but must rely on computational models of the stimuli,
78 typically borrowing from artificial intelligence [Varoquaux and Poldrack, 2019]. Second, they model
79 brain responses using high-dimensional statistical models focused on prediction. These can fit more
80 ingredients jointly, avoiding difficult modeling choices. As a result, they can generalize findings
81 across stimuli probing different parts of a cognitive ontology: natural images, simplified faces, or
82 wedges traditionally used for retinotopic mappings [Eickenberg et al., 2017]. This is in sharp contrast
83 with conventional brain mapping methodology: based on oppositions between stimuli, it does not
84 lead to formal models bridging results from different experimental paradigms.

85 2 How do statistical tools fit in scientific progress

86 2.1 From scientific evidence to scientific knowledge: more than data

87 **Internal versus external validity** The validity of a study’s findings is more than a statistical
88 question. Internal validity controls inferences about the relations across the quantities in the study,
89 for instance that measurements have no unmodeled errors. External validity, more important but less
90 discussed, asserts that those relations are maintained beyond the study’s settings [Cook and Campbell,
91 1979]. It may for instance fail when running a study on a sample non representative of the population.

92 **Validity of constructs** Scientific theories and models are constructed from abstract ingredients
93 such as “intelligence” or phrenology’s “alimentiveness” in psychology. These *constructs* are central
94 to reasoning about empirical evidence, to position it in a broader context. A good construct is one
95 that is useful to explain many different observations, beyond a single study [Cronbach and Meehl,
96 1955]. Interpreting an empirical study in a theoretical framework requires *construct validity* of its
97 measures and manipulations: that these indeed to relate well to the construct of interest. For instance,
98 to be interpretable as intelligence, IQ tests should not be confounded by cultural knowledge.

99 **Stances on theories** Models, and thus theory, are needed to interpret empirical finding. The
100 acceptance of these theories often builds upon implicit stances on their ingredients. In psychology,
101 Fried [2020] argues that statistical models should build on “strong theories” and provide “explanation
102 of a phenomenon” relating valid psychological constructs, beyond mere data fit. Yarkoni [2020]
103 points out that such a view carries implicit preferences on choices of construct that may be difficult to
104 defend. In particular, such model esthetic assumes realism about psychological constructs: that these
105 have an absolute existence beyond the minds of the scientists. A scientific discourse must position its
106 claims on unobservable constructs, for instance centers of gravity in mechanics. *Realism* accepts to
107 build scientific knowledge on unobservable entities only if they are objective and mind-independent.
108 *Instrumentalism*, rather, accepts that some ingredients of theories are mere instruments needed to tie
109 together observable outcomes, and that the success of a theory is asserted solely on these observables.

110 Questions on the validity of basic modeling ingredients are less discussed in a well-established
111 science such as physics, as there is a consensus on the ingredients: forces, acceleration, temperature
112 –which has a non-trivial definition–... And yet, this consensus was achieved through iterations.
113 Planetary observations in the times of Kepler were analyzed with phenomenological models lacking
114 the ingredients of dynamics, but were fundamental to nourishing Newtonian mechanics.

115 2.2 Reasoning with statistical tools

116 Statistics gives the scientist tools to reason from noisy observations. The prevailing approach is
117 **model reasoning**: a probabilistic model describing data generation is built, encompassing ingredients
118 of the application domain. Parameters estimated using this model are interpreted within its logic [Cox,
119 2006, chap 9]. Cox [2001] goes as far as saying that statistical models are “efforts to establish data
120 descriptions that are potentially causal”. Another form of reasoning –design-based [Cox, 2006, chap
121 9] or **warranted reasoning** [Cook, 1991, Baiocchi and Rodu, 2021]– relies on specific experimental
122 design, as randomization, for causal inferences without a model of the data-generating mechanism.
123 Finally, Breiman [2001] famously noted that increasingly many statistical tools forgo data modeling,
124 to focus on algorithmic capacity to approximate relations. Their success is established by **outcome**
125 **reasoning**: gauging predictions on observables [Baiocchi and Rodu, 2021], key to machine learning.

126 3 Grounding more statistical reasoning on output rather than models

127 With a historical emphasis on data modeling, statistics has an implicit realism stance. Yet, as we have
128 seen in physics or vision neuroscience, scientific progress is achieved despite analyzing observations
129 without the right conceptual framework. Outcome reasoning, with tools of machine learning, gives a
130 robust statistical framework for science: given imperfect premises, it fails less.

131 **3.1 Robustness to model mis-specification**

132 With model reasoning, parameters can be interpreted only conditional to the choice of model, which
133 is outside of statistical control. Statisticians often assume that many hard modeling questions can be
134 resolved by domain experts. Yet science is performed by limited beings [Wimsatt, 2007] and even
135 experts have finite resources to dedicate to a given problem [Simon, 1955]. Model imperfections can
136 have vast consequences on statistical conclusions. Botvinik-Nezer et al. [2020] asked 70 different
137 teams of experts to analyze the same brain imaging data. Variations in modeling choices –all based on
138 linear models– lead to vastly different parameters, and qualitatively different neuroscience findings.

139 Controlling predictions instead of model parameters leads to a different statistical regime. Even the
140 simple case of the linear model changes drastically: with learning theory, analysis is possible even
141 in the miss-specified setting, showing that multi-collinearity in the design is not an issue [Hsu et al.,
142 2014], unlike when performing inference on model parameters. Higher-dimensional settings are
143 possible, which means that the analyst no longer has to cherry-pick a small number of descriptors.
144 In neuroscience, it has enabled studying richer descriptions of the stimuli, generated by artificial
145 intelligence techniques rather than set in a specific reductionist theoretical framework. Switching to
146 output reasoning requires reinventing analytical paradigms: in brain imaging switching to *decoding*
147 models that gauge the ability to *predict* neural responses.

148 **3.2 Putting explicit generalization at the center of the inference**

149 Judging a model by its predictions is good science. It shifts the burden on validity on observables.
150 These may suffer biases, such as censoring, which must be accounted for even in machine-learning
151 settings [Ishwaran et al., 2008]. But in the long run, the validity of scientific theories is established by
152 their ability to generalize across many settings.

153 Cross-validation on a study sample is however not a test of a strong ability to generalize; it gives
154 no evidence of external validity. Machine-learning models may easily create local approximations
155 which do not generalize to new settings, bad scientific models. Yet, their ability to generalize can
156 be explicitly tested. This is unlike model-based tests of qualitative theories, as in psychology or
157 sociology. Indeed, a methodology based on machine learning can be applied to rich descriptions
158 of the objects under study –the raw images presented–, while model-based reasoning is applied to
159 a small number of features, specially crafted to represent the constructs of interest –a face-place
160 opposition. In the former, the generalization is readily tested on data from different settings via a
161 quantitative prediction error. In the latter, the finding is more conceptual and given a new setting it
162 must be instantiated with a new modeling effort.

163 Beyond broad generalization, an oft-requested feature of an analytical model is to provide “under-
164 standing”. For domain reasoning, it is helpful to try to tease out the contribution of various ingredients.
165 An emerging non-parametric statistical toolbox is catering to this purpose: black-box explanation
166 techniques [Molnar, 2020], such as partial dependency plot [Friedman, 2001] or the knock-off [Barber
167 and Candès, 2019]. These tools ground their inferences on model outputs, the quantities amenable
168 to strong empirical validation. Demanding more from an analytical model, for instance opposing
169 phenomenological data explanations with valid theoretical understanding, forces buying into a given
170 theoretical framework Yarkoni [2020], with the risk of circular reasoning on the evidence.

171 Parametric models are appealing for intuitive counterfactual reasoning [Angrist and Pischke, 2008]:
172 they appear as “data descriptions that are potentially causal” [Cox, 2001]. Yet, more than a para-
173 metric model, valid causal inference needs a structural characterization of variables, distinguishing
174 confounders, colliders, mediators... [Greenland et al., 1999]. In such settings, machine learning
175 models shine by their potential robustness to mismodeling [Rose and Rizopoulos, 2020].

176 **Black-box models for thinking outside the box** Empirical validation of a theory tied to its
177 ingredients smells of self-fulfilling prophecies. This is the risk of model-based statistical reasoning.
178 Science needs statistical reasoning based more on model predictions. Machine learning will provide
179 the building blocks, for broad generalization and counterfactual reasoning.

180 References

- 181 T Aaltonen, J Adelman, T Akimoto, Michael G Albrow, B Álvarez González, S Amerio, D Amidei, A Anastassov,
182 Alberto Annovi, J Antos, et al. Measurement of the single-top-quark production cross section at cdf. *Physical*
183 *review letters*, 101(25):252001, 2008.
- 184 Terhi Aaltonen, J Adelman, Tb Akimoto, B Álvarez González, Sara Amerio, Da Amidei, Aa Anastassov,
185 A Annovi, J Antos, G Apollinari, et al. Observation of electroweak single top-quark production. *Physical*
186 *review letters*, 103(9):092002, 2009.
- 187 Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics*. Princeton university press, 2008.
- 188 Michael Baiocchi and Jordan Rodu. Reasoning using data: Two old ways and one new. *Observational Studies*, 7
189 (1):3–12, 2021.
- 190 Rina Foygel Barber and Emmanuel J Candès. A knockoff filter for high-dimensional selective inference. *The*
191 *Annals of Statistics*, 47(5):2504–2537, 2019.
- 192 Rotem Botvinik-Nezer, Felix Holzmeister, Colin F Camerer, Anna Dreber, Juergen Huber, Magnus Johannesson,
193 Michael Kirchler, Roni Iwanir, Jeanette A Mumford, R Alison Adcock, et al. Variability in the analysis of a
194 single neuroimaging dataset by many teams. *Nature*, 582(7810):84–88, 2020.
- 195 Nora Mills Boyd and James Bogen. Theory and Observation in Science. In Edward N. Zalta, editor, *The Stanford*
196 *Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2021 edition, 2021.
- 197 Leo Breiman. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical*
198 *science*, 16(3):199–231, 2001.
- 199 TD Cook and DT Campbell. *Quasi-experimentation: Design and analysis issues for field settings 1979 Boston*.
200 MA Houghton Mifflin, 1979.
- 201 Thomas D Cook. Clarifying the warrant for generalized causal inferences in quasi-experimentation. In *Evaluation*
202 *and education: At quarter century*. 1991.
- 203 David R Cox. [statistical modeling: The two cultures]: Comment. *Statistical science*, 16(3):216–218, 2001.
- 204 David Roxbee Cox. *Principles of statistical inference*. Cambridge university press, 2006.
- 205 Lee J. Cronbach and Paul E. Meehl. Construct validity in psychological tests. *Psychological Bulletin*, 52:281,
206 1955.
- 207 Michael Eickenberg, Alexandre Gramfort, Gaël Varoquaux, and Bertrand Thirion. Seeing it all: Convolutional
208 network layers map the function of the human visual system. *NeuroImage*, 152:184–194, 2017.
- 209 Eiko I Fried. Lack of theory building and testing impedes progress in the factor and network literature.
210 *Psychological Inquiry*, 31(4):271–288, 2020.
- 211 Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages
212 1189–1232, 2001.
- 213 Sander Greenland, Judea Pearl, and James M Robins. Causal diagrams for epidemiologic research. *Epidemiology*,
214 pages 37–48, 1999.
- 215 Daniel Hsu, Sham Kakade, and Tong Zhang. Random design analysis of ridge regression. *Foundations of*
216 *Computational Mathematics*, 14, 2014.
- 217 D H Hubel and T N Wiesel. Receptive fields of single neurones in the cat’s striate cortex. *J. Physiol.*, 148:
218 574–591, 1959.
- 219 Hemant Ishwaran, Udaya B Kogalur, Eugene H Blackstone, and Michael S Lauer. Random survival forests. *The*
220 *annals of applied statistics*, 2(3):841–860, 2008.
- 221 N Kanwisher, J McDermott, and M M Chun. The fusiform face area: a module in human extrastriate cortex
222 specialized for face perception. *J. Neurosci.*, 17(11):4302–4311, 1997.
- 223 Nikos K Logothetis, Jon Pauls, and Tomaso Poggio. Shape representation in the inferior temporal cortex of
224 monkeys. *Current biology*, 5(5):552–563, 1995.
- 225 David Marr. *Vision: A computational investigation into the human representation and processing of visual*
226 *information*. The MIT press, Cambridge, 1982.

- 227 Christoph Molnar. *Interpretable machine learning*. Lulu.com, 2020.
- 228 Nicolas Pinto, David Doukhan, James J DiCarlo, and David D Cox. A high-throughput screening approach to
229 discovering good forms of biologically inspired visual representation. *PLoS computational biology*, 5(11):
230 e1000579, 2009.
- 231 Russell A Poldrack. Mapping mental function to brain structure: how can cognitive neuroimaging succeed?
232 *Perspectives on psychological science*, 5:753, 2010.
- 233 Alexander Radovic, Mike Williams, David Rousseau, Michael Kagan, Daniele Bonacorsi, Alexander Himmel,
234 Adam Aurisano, Kazuhiro Terao, and Taritree Wongjirad. Machine learning at the energy and intensity
235 frontiers of particle physics. *Nature*, 560(7716):41–48, 2018.
- 236 Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex. *Nature*
237 *neuroscience*, 2(11):1019–1025, 1999.
- 238 Sherri Rose and Dimitris Rizopoulos. Machine learning for causal inference in biostatistics. *Biostatistics*, 21(2):
239 336–338, 2020.
- 240 Pierre Sermanet, David Eigen, Xiang Zhang, Michael Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Inte-
241 grated recognition, localization and detection using convolutional networks. In *2nd International Conference*
242 *on Learning Representations, ICLR 2014*, 2014.
- 243 Herbert A Simon. A behavioral model of rational choice. *The quarterly journal of economics*, 69(1):99–118,
244 1955.
- 245 Torbjörn Sjöstrand, Patrik Eden, Christer Friberg, Leif Lönnblad, Gabriela Miu, Stephen Mrenna, and Emanuel
246 Norrbin. High-energy-physics event generation with pythia 6.1. *Computer physics communications*, 135(2):
247 238–259, 2001.
- 248 Michael J Tarr and Isabel Gauthier. Ffa: a flexible fusiform area for subordinate-level visual processing
249 automatized by expertise. *Nature neuroscience*, 3:764, 2000.
- 250 David A van Dyk. The role of statistics in the discovery of a higgs boson. *Annual Review of Statistics and Its*
251 *Application*, 1:41–59, 2014.
- 252 Gaël Varoquaux and Russell A Poldrack. Predictive models avoid excessive reductionism in cognitive neu-
253 roimaging. *Current opinion in neurobiology*, 55:1–6, 2019.
- 254 William C Wimsatt. False models as means to truer theories. *Neutral models in biology*, pages 23–55, 1987.
- 255 William C Wimsatt. *Re-engineering philosophy for limited beings: Piecewise approximations to reality*. Harvard
256 University Press, 2007.
- 257 Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo.
258 Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of*
259 *the national academy of sciences*, 111(23):8619–8624, 2014.
- 260 Tal Yarkoni. Implicit realism impedes progress in psychology: Comment on fried (2020). *Psychological Inquiry*,
261 31(4):326–333, 2020.