

1 Submitted to *Journal of Vision*.

2 EMERGING TREND IN VISION SCIENCE

3 **Deep learning: Using machine learning** 4 **to study biological vision**

5 Najib J. Majaj¹ and Denis G. Pelli^{1,2}

6 ¹Center for Neural Science and ²Department of Psychology, New York University

7

8

9 **ABSTRACT**

10 Today many vision-science presentations employ machine learning and especially “deep
11 learning”, one of the more recent and successful variants. Many neuroscientists use machine
12 learning to decode neural responses. Many perception scientists try to understand how living
13 organisms recognize objects. To them, deep neural networks offer several benchmark
14 accuracies for recognition of learned stimuli. Originally machine learning was inspired by the
15 brain. Today, machine learning is used as a statistical tool to decode brain activity. Tomorrow,
16 deep neural networks might become our best model of brain function. This brief overview of the
17 use of machine learning in biological vision touches on its strengths, weaknesses, milestones,
18 controversies, and current directions. Here, we hope to help vision scientists assess what role
19 machine learning should play in their research.

20

21 INTRODUCTION

22 What does machine learning offer to biological-vision
23 scientists? Machine-learning was developed as a tool for
24 automated classification, optimized for accuracy.
25 Physiologists use it to identify stimuli based on neural
26 activity. Physiologists and psychophysicists are starting to
27 consider deep learning as a model for object recognition
28 by human and nonhuman primates (Cadieu et al., 2014;
29 Ziskind et al., 2014; Yamins et al., 2014; Khaligh-Razavi &
30 Kriegeskorte, 2014; Testolin, Stoianov, & Zorzi, 2017). We
31 suppose that most of our readers have heard of machine
32 learning but are wondering whether it would be useful in
33 their own research. We begin by describing some of its
34 pluses and minuses.

35 PLUSES: WHAT IT'S GOOD FOR

36 Deep learning is very popular (Fig. 1). Is it just a fad? At
37 the very least, machine learning is a powerful tool for
38 interpreting biological data. For computer vision, the old
39 paradigm was: feature detection, followed by
40 segmentation, and then grouping (Marr, 1982). With
41 machine learning tools, the new paradigm is to just define
42 the task and provide a set of labeled examples, and the
43 algorithm builds the classifier. Unlike the handcrafted
44 pattern recognition (including segmentation and

GLOSSARY

Machine learning is a computer algorithm that learns how to perform a task directly from examples, without a human providing explicit instructions or rules for how to do so. Correctly labeled examples are provided to the learning algorithm, which is then “trained” (i.e. its parameters are gradually adjusted) to be able to perform the task correctly on its own and generalize to unseen examples.

Deep learning is a newly successful and popular version of machine learning that uses backprop neural networks with multiple hidden layers. The 2012 success of AlexNet, then the best machine learning network for object recognition, was the tipping point. It is now ubiquitous in the internet. The idea is to have each layer of processing perform successively more complex computations on the data to give the full “multi-layer” network more expressive power. The drawback is that it is much harder to train multi-layer networks (Goodfellow et al. 2016). Deep learning ranges from discovering the weights of a multilayer network to parameter learning in hierarchical belief networks. Note that the complexity of deep learning may be unwarranted for simple problems that are well handled by, e.g. SVM. Try shallow networks first, when they fail, go deep.

45 grouping) popular in the 70's and 80's, machine-learning
46 algorithms are generic, with little domain-specificity. They
47 replace hand-engineered feature detectors with filters that
48 can be learned from the data. Advances in the mid 90's in
49 machine learning made statistical learning theory useful
50 for practical classification, e.g. handwriting recognition
51 (LeCun et al., 1989; Vapnik, 1999).

52 Machine learning allows a neurophysiologist to decode
53 neural activity without knowing the receptive fields (Seung
54 & Sompolinsky, 1993; Hung et al., 2005). Machine learning
55 shifts the emphasis from how the cells encode to what
56 they encode, i.e. what that code tells us about the
57 stimulus. Mapping a receptive field is the foundation of
58 neuroscience (beginning with Weber's 1834/1996
59 mapping of tactile "sensory circles"), but many young
60 scientists are impatient with the limitations of single-cell
61 recording: looking for minutes or hours at how one cell
62 responds to each of perhaps a hundred different stimuli.
63 New neuroscientists are the first generation for whom it is
64 patently clear that characterization of a single neuron's
65 receptive field, which was invaluable in the retina and V1,
66 fails to characterize how higher visual areas encode the
67 stimulus. Statistical learning techniques reveal "how
68 neuronal responses can best be used (combined) to

Neural nets are computing systems inspired by biological neural networks that learn tasks by considering examples.

Supervised learning refers to any algorithm that accepts a set of labeled stimuli — a training set — and returns a classifier that can label stimuli similar to those in the training set.

Unsupervised learning works without labels. It is less popular, but of great interest because labeled data are scarce while unlabeled data are plentiful. Without labels, the algorithm discovers structure and redundancy in the data.

Cost function. A function that assigns a real number representing cost to a candidate solution, i.e. a set of weights. Solving by optimization means minimizing cost.

Gradient descent: An algorithm that minimizes cost by incrementally changing the parameters in the direction of steepest descent of the cost function.

69 inform perceptual decision-making” (Graf, Kohn, Jazayeri,
70 & Movshon, 2010). The simplicity of the machine
71 decoding can be a virtue as it allows us to discover what
72 can be easily read-out (e.g. by a single downstream
73 neuron) (Hung et al. 2005). Achieving psychophysical
74 levels of performance in decoding a stimulus object’s
75 identity and location from the neural response shows that
76 the measured neural performance has all the information
77 needed for the subject to do the task (Majaj et al. 2015;
78 Hong et al. 2016).

79 For psychophysics, Signal Detection Theory (SDT) proved
80 that the optimal classifier for a known signal in white noise
81 is a template matcher (Peterson, Birdsall, & Fox, 1954;
82 Tanner & Birdsall, 1958). Of course, SDT solves only a
83 simple version of the general problem of object
84 recognition, which includes variation in viewing conditions
85 and diverse objects within a category (e.g. a chair can be
86 any object that affords sitting). SDT introduces the very
87 useful idea of a mathematically defined ideal observer,
88 providing a reference for human performance (e.g.
89 Geisler, 1989; Pelli et al., 2006). However, one drawback
90 is that it doesn’t incorporate learning. Deep learning, on
91 the other hand, provides a pretty good observer that
92 learns, which may inform studies of human learning.

Convexity: A problem is convex if there are no local minima other than the global minimum (or minima if there are several equally good solutions). This guarantees that gradient-descent will converge to a global minimum. There might be more than one global minimum, with equal cost, e.g. in problems with symmetric solutions.

Generalization is how well a classifier performs on new, unseen examples that it did not see during training.

Cross validation assesses the ability of the network to generalize, from the data that it trained on, to new data.

Backprop, short for "backward propagation of errors", is widely used to apply gradient-descent learning to multi-layer networks. It uses the chain rule from calculus to iteratively compute the gradient of the cost function for each layer.

Hebbian learning and spike-timing-dependent plasticity (**STDP**). According to Hebb’s rule, the efficiency of a synapse increases after correlated pre- and post-synaptic activity. In other words, neurons that fire together, wire together (Löwel & Singer, 1992).

93 These networks might reveal the constraints imposed by
94 the training set on learning. Further, unlike SDT, deep
95 neural networks cope with the complexity of real tasks. It
96 can be hard to tell whether behavioral performance is
97 limited by the set of stimuli, their neural representation, or
98 the observer's decision process (Majaj et al. 2015).
99 Implications for classification performance are not readily
100 apparent from direct inspection of families of stimuli and
101 their neural responses. SDT specifies optimal
102 performance for classification of known signals but does
103 not tell us how to generalize beyond a training set.
104 Machine learning does.

105 106 **MINUSES: COMMON COMPLAINTS** 107

108 Some biologists point out that neural nets do not match
109 what we know about neurons (e.g., Crick, 1989; Rubinov,
110 2015). Biological brains learn on the job, while neural
111 networks need to converge before they can be used. A
112 state-of-the-art deep neural network needs five thousand
113 labelled object images per category to match human
114 recognition accuracy (Goodfellow et al., 2016). But
115 children and adults need only a hundred labelled letters of
116 an unfamiliar alphabet to reach the same accuracy as
117 fluent native readers (Pelli et al. 2006). In particular, it is
118 not clear, given what we know about neurons and neural
119 plasticity, whether a backpropagation network can be implemented using biologically plausible

Support Vector Machine (SVM)

is a learning machine for classification. SVMs generalize well. An SVM can quickly learn to perform a nonlinear classification using what is called the “kernel trick”, mapping its input into a high-dimensional feature space (Cortes & Vapnik, 1999).

Convolutional neural networks (ConvNets)

have their roots in the Neocognitron (Fukushima 1980) and are inspired by the simple and complex cells described by Hubel and Wiesel (1962). ConvNets apply backprop learning to multilayer neural networks based on convolution and pooling (LeCun et al., 1989; LeCun et al., 1990; LeCun et al., 1998).

120 circuits (but see Mazzone et al., 1991, and Bengio et al., 2015). However, there are several
121 promising efforts to implement more biological plausible learning rules, e.g. spike-timing-
122 dependent plasticity (Mazzone et al., 1991; Bengio et al., 2015; Sacramento, Costa, Bengio, &
123 Senn 2017).

124 Unlike the biologists' desire to model, engineers and computer scientists, while inspired by
125 biological vision, focus on what works. To this, one might counter that every biological model is
126 an abstraction and can be useful even while failing to capture all the details of the living
127 organism.

128 Some physiologists note that decoding neural activity to recover the stimulus is interesting and
129 useful but falls short of explaining what the neurons do.

130 Some visual psychophysicists note some salient differences between performance of human
131 observers and deep networks on tasks like object recognition and image distortion (Ullman et al.
132 2016; Berardino et al. 2017).

133 Some biological modelers complain that neural nets have alarmingly many parameters. Deep
134 neural networks continue to be opaque. Before neural-network modeling, a model was simpler
135 than the data it explained. Deep neural nets are typically as complex as the data, and the
136 solutions are hard to visualize (but see Zeiler & Fergus, 2013). However, while the training sets
137 and learned weights are long lists, the generative rules for the network (the computer programs)
138 are short. Traditionally, having very many parameters has often led to overfitting, i.e. a failure to
139 generalize beyond the training set, but the breakthrough is that deep-learning networks with a
140 huge number of parameters nevertheless generalize well.

141 Some cognitive psychologists dismiss deep neural networks as unable to “master some of the
142 basic things that children do, like learning the past tense of a regular verb” (Marcus et al., 1992).

143 Some statisticians worry that rigorous statistical tools are being displaced by machine learning,
144 which lacks rigor (Friedman, 1998; Matloff, 2014, but see Breiman, 2001; Efron & Hastie, 2016).
145 Assumptions are rarely stated. There are no confidence intervals on the solution. However,
146 performance is typically cross-validated, showing generalization, and it has been proven that
147 convex networks can compute posterior probability (e.g. Rojas, 1996). Furthermore, machine
148 learning and statistics seem to be converging to provide a more general perspective on
149 probabilistic inference that combines complexity and rigor.

150 These current limitations drive practitioners to enhance the scope and rigor of deep learning.
151 But bear in mind that some of the best classifiers in computer science were inspired by
152 biological principles (Rosenblatt, 1957; 1958; Rumelhart et al., 1986; LeCun, 1985; LeCun et al.
153 1989; LeCun et al. 1990; Riesenhuber & Poggio, 1999; and see LeCun, Bengio, Hinton 2015).
154 Some of those classifiers are now so good that they occasionally exceed human performance
155 and might serve as rough models for how biological systems classify (e.g. Yamins, et al. 2014;
156 Khaligh-Razavi & Kriegeskorte, 2014; Ziskind, Hénaff, LeCun, & Pelli, 2014; Testolin, Stoianov,
157 & Zorzi, 2017).

158 **MATHEMATICS VS. ENGINEERING**

159 The history of machine learning has two threads: mathematics and engineering. In the
160 *mathematical* thread, two statisticians, Fisher and later Vapnik, developed mathematical
161 transformations to uncover categories in data, and proved that they give unique answers. They
162 assumed distributions and proved convergence.

163 In the *engineering* thread, a loose coalition of psychologists, neuroscientists, and computer
164 scientists (e.g. Turing, Rosenblatt, Minsky, Fukushima, Hinton, Sejnowski, LeCun, Poggio,
165 Bengio) sought to reverse-engineer the brain to build a machine that learns. Their algorithms
166 are typically applied to stimuli with unknown distributions and lack proofs of convergence.

167 **MILESTONES IN CLASSIFICATION**

168 1936: Linear discriminant analysis

169 1953: Machine learning

170 1958: Perceptron

171 1969: Death of the perceptron

172 1974: Backprop

173 1980: Neocognitron

174 1987: NETtalk

175 1989: ConvNets

176 1995: Support Vector Machine (SVM)

177 2006: Backprop, revived

178 2012: Deep learning

179

180 **1936: Linear discriminant analysis.** Fisher (1936) introduced linear discriminant analysis to

181 classify two species of iris flower based on four measurements per flower. When the distribution

182 of the measurements is normal and the covariance matrix between the measurements is known,

183 linear discriminant analysis answers the question: Supposing we use a single-valued function to

184 classify, what linear function $y = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$, of four measurements x_1, x_2, x_3, x_4

185 made on flowers, with free weights w_1 , w_2 , w_3 , w_4 , will maximize discrimination of species?¹

186 Linear classifiers are great for simple problems for which the category boundary is a hyperplane
187 in a small number of dimensions. However, complex problems like object recognition typically
188 require more complex category boundaries in a large number of dimensions. Furthermore, the
189 distributions of the features are typically unknown and not necessarily normal.

190 Cortes & Vapnik (1995) note that the first algorithm for pattern recognition was Fisher's optimal
191 decision function for classifying vectors from two known distributions. Fisher solved for the
192 optimal classifier in the presence of gaussian noise and known covariance between elements of
193 the vector. When the covariances are equal, this reduces to a linear classifier. The ideal
194 template matcher of signal detection theory is an example of such a linear classifier (Peterson et
195 al., 1954). This fully specified simple problem can be solved analytically. Of course, many
196 important problems are not fully specified. In everyday perceptual tasks, we typically know only
197 a "training" set of samples and labels.

198 **1953: Machine learning.** The first developments in machine learning were to play chess and
199 checkers. "Could one make a machine to play chess, and to improve its play, game by game,
200 profiting from its experience?" (Turing, 1953). Arthur Samuel (1959) defined *machine learning*
201 as the "Field of study that gives computers the ability to learn without being explicitly
202 programmed."

203 **1958: Perceptron.** Inspired by physiologically measured receptive fields, Rosenblatt (1958)
204 showed that a very simple neural network, the perceptron, could learn to classify from training
205 samples. Perceptrons combined several linear classifiers to implement piecewise-linear

¹ Linear discriminant analysis is an outgrowth of regression which has a much longer history. Regression is the optimal least-squares linear combination of given functions to fit given data and was applied by Legendre (1805) and Gauss (1809) to astronomical data to determine the orbits of the comets and planets around the sun. The estimates come with confidence intervals and the fraction of variance accounted for, which rates the goodness of the explanation.

206 separating surfaces. The perceptron learns the weights to use in a linear combination of feature-
207 detector outputs. The perceptron transforms the stimulus into a binary feature vector and then
208 applies a linear classifier to the feature vector. The perceptron is piecewise linear and has the
209 ability to learn from training examples without knowing the full distribution of the stimuli. Only the
210 final layer in the perceptron learns.

211 **1969: Death of the perceptron.** However, it quickly became apparent that the perceptron and
212 other single-layer neural networks cannot learn tasks that are not linearly separable, i.e. cannot
213 solve problems like connectivity (Are all elements connected?) and parity (Is the number of
214 elements odd or even?); people solve these readily (Minsky & Papert, 1969). On this basis
215 Minsky and Papert announced the death of artificial neural networks.

216 **1974: Backprop.** The death of the perceptron showed that learning in a one-layer network was
217 too limited. This impasse was broken by the introduction of the backprop algorithm, which
218 allowed learning to propagate through multiple-layer neural networks. The history of backprop is
219 complicated (see Schmidhuber, 2015). The idea of minimization of error through a differentiable
220 multi-stage network was discussed as early as the 1960s (e.g. Bryson, Denham, & Dreyfus,
221 1963). It was applied to artificial neural networks in the 1970s (e.g. Werbos, 1974). In the 1980s,
222 efficient backprop first gained recognition, and led to a renaissance in the field of artificial neural
223 network research (LeCun, 1985; Rumelhart, Hinton, & Williams, 1986). During the 2000s
224 backprop neural networks fell out of favor, due to four limitations (Vapnik, 1999): **1. No proof of**
225 **convergence.** Backprop uses gradient descent. Gradient descent with a nonconvex error
226 function with multiple minima is only guaranteed to find a local, not the global minimum of the
227 error function. This has long been considered a major limitation, but Yann LeCun et al. (2015)
228 claim that it hardly matters in practice in current implementations of deep learning. **2. Slow.**
229 Convergence to a local minimum can be slow due to the high dimensionality of the weight
230 space. **3. Poorly specified.** Backprop neural networks had a reputation for being ill-specified,

231 an unconstrained number of units and training examples, and a step size that varied by
232 problem. “Neural networks came to be painted as slow and fussy to train [,] beset by voodoo
233 parameters and simply inferior to other approaches.” (Cox & Dean, 2014). **4. Not biological.**
234 Lastly, backprop learning may not to be physiological: While there is ample evidence for
235 Hebbian learning (increase of a synapse’s gain in response to correlated activity of the two cells
236 that it connects), such changes are never propagated backwards, beyond the one synapse, to a
237 previous layer. **5. Inadequate resources.** With hindsight it is clear that backprop in the 80’s was
238 crippled by limited computing power and lack of large labeled datasets.

239 **1980: Neocognitron**, the first convolutional neural network. Fukushima (1980) proposed and
240 implemented the Neocognitron, a hierarchical, multilayer artificial neural network. It recognized
241 stimulus patterns (deformed numbers) despite small changes in position and shape.

242 **1987: NETtalk**, the first impressive backprop neural network. Sejnowski et al. (1987) reported
243 the exciting success of NETtalk, a neural network that learned to convert English text to speech:
244 “*The performance of NETtalk has some similarities with observed human performance. (i) The*
245 *learning follows a power law. (ii) The more words the network learns, the better it is at*
246 *generalizing and correctly pronouncing new words. (iii) The performance of the networks*
247 *degrades very slowly as connections in the network are damaged: no single link or processing*
248 *unit is essential. (iv) Relearning after damage is much faster than learning during the original*
249 *training...*”

250 **1989: ConvNets.** Yann LeCun and his colleagues combined convolutional neural networks with
251 backprop to recognize handwritten characters (LeCun et al., 1989; LeCun et al., 1990). This
252 network was commercially deployed by AT&T, and today reads millions of checks a day
253 (LeCun, 1998). Later, adding half-wave rectification and max pooling greatly improved its
254 accuracy in recognizing objects (Jarrett et al., 2009).

255 **1995: Support Vector Machine (SVM).** Cortes & Vapnik (1995) proposed the support vector
256 network, a learning machine for binary classification problems. SVMs generalize well and are
257 free of mysterious training parameters. Many versions of the SVM are convex (e.g. Lin, 2001).

258 **2006: Backprop, revived.** Hinton & Salakhutdinov (2006) sped up backprop learning by
259 unsupervised pre-training. This helped to revive interest in backprop. In the same year, a
260 supervised backprop-trained convolutional neural network set a new record on the famous
261 MNIST handwritten-digit recognition benchmark (Ranzato et al., 2006).

262 **2012: Deep learning.** Geoff Hinton says, “It took 17 years to get deep learning right; one year
263 thinking and 16 years of progress in computing, praise be to Intel.” (Cox & Dean, 2014; LeCun,
264 Bengio, & Hinton, 2015). It is not clear who coined the term “deep learning”.² In their book, *Deep*
265 *Learning Methods and Applications*, Deng & Yu (2014) cite Hinton et al. (2006) and Bengio
266 (2009) as the first to use the term. However, the big debut for deep learning was an influential
267 paper by Krizhevsky et al. (2012) describing AlexNet, a deep convolutional neural network that
268 classified 1.2 million high-resolution images into 1000 different classes, greatly outperforming
269 previous state-of-the-art machine learning and classification algorithms.

270 **CONTROVERSIES**

271 The field is growing quickly, yet certain topics remain hot. For proponents of deep learning, the
272 ideal network is composed of simple elements and learns everything from the training data. On
273 the other extreme, computer vision scientists argue that we know a lot about how the brain
274 recognizes objects that we can engineer into the networks before learning (e.g. gain control and
275 normalization). Some engineers look to the brain only to copy strengths of the biological
276 solution, others think there are useful clues in its limitations as well (e.g. crowding).

² The idea of “deep learning” is not exclusive to machine learning and neural networks (e.g. Dechter, 1986)

277 **Is deep learning the best solution for all visual tasks?** Deep learning is not the only thing in
278 the vision scientist's toolbox. Object recognition as a visual task has been very useful in vision
279 research because it is an objective task that is easily scored as right or wrong, is essential in
280 daily life, and captures some of the magic of seeing. Deep neural nets solve it, albeit with a
281 million parameters. Another interesting task is detection of image distortion. Currently a simple
282 model implementing gain-control normalization performs this much better than deep networks
283 do (Berardino et al. 2017). Scientists, like the brain, use whatever tool works best.

284 **Unproven convexity.** A problem is convex if there are no local minima other than the global
285 minimum (or minima if there are several equally good solutions). This guarantees that gradient-
286 descent will converge to a global minimum. As far as we know, classifiers that give inconsistent
287 results are not useful. Conservation of a solution across seeds and algorithms is evidence for
288 convexity. For some combinations of stimuli, categories, and classifiers, convexity can be
289 proved. For others, empirical tests can provide qualified assurance that the solution is a global
290 minimum. Many widely used networks are not convex, but still give mostly consistent answers
291 (LeCun, Bengio, & Hinton, 2015). In machine learning, kernel methods, including learning by
292 SVMs, have the advantage of easy-to-prove convexity, at the cost of limited generalization. In
293 the 1990s, SVMs were popular because they guaranteed fast convergence even with a large
294 number of training samples (Cortes & Vapnik, 1995). Thus, when the problem is convex, the
295 quality of solution is assured, and one can rate implementations by their demands for size of
296 network and training sample. Deep neural networks, on the other hand, generalize well, but are
297 not convex.

298 **Shallow vs. deep networks.** The field's imagination has focused alternately on shallow and
299 deep networks, beginning with the Perceptron in which only one layer learned, to backprop,
300 which allowed multiple layers and cleared the hurdles that doomed the Perceptron. Then SVM,
301 with its single layer, sidelined the multilayer backprop. Today multilayer deep learning reigns;

302 Krizhevsky et al. (2012) attributed the success of their network to its 8-layer depth; it performed
303 worse with fewer layers.

304 **Supervised vs. unsupervised.** Learning algorithms for a classifier can be supervised or not,
305 i.e. need labels for training, or don't. Today most machine learning is *supervised* (LeCun,
306 Bengio, & Hinton, 2015). The images are labeled (e.g. "car" or "face"), or the network receives
307 feedback on each trial from a cost function that assesses how well its answer matches the
308 image's category. In *unsupervised* learning, no labels are given. The algorithm processes
309 images, typically to minimize error in reconstruction, with no extra information about what is in
310 the (unlabeled) image. A cost function can also reward decorrelation and sparseness (e.g.
311 Olshausen and Field, 1996). This allows learning of image statistics and has been used to train
312 early layers in deep neural networks. Human learning of categorization is sometimes done with
313 explicitly named objects — "Look at the tree!" — but more commonly the feedback is implicit.
314 Consider reaching your hand to raise a glass of water. Contact informs vision. On specific
315 benchmarks, where the task is well-defined and labeled examples are available, supervised
316 learning can excel (e.g. AlexNet), but unsupervised learning may be more useful when few
317 labels are available.

318 **CURRENT DIRECTIONS**

319 **What does deep learning add to the vision-science toolbox?** Deep learning is more than
320 just a souped-up regression (Marblestone et al., 2016). Like Signal Detection Theory (SDT), it
321 allows us to see more in our behavioral and neural data. In the 1940's, Norbert Wiener and
322 others developed algorithms to automate and optimize signal detection and classification. A lot
323 of it was engineering. The whole picture changed with the SDT theorems, mainly the proof that
324 the maximum-likelihood receiver is optimal for a wide range of simple tasks (Peterson et al.,
325 1954). In white noise a traditional receptive field computes the likelihood of the presence of a

326 signal matching the receptive field weights. It was exciting to realize that the brain contains 10^{11}
327 likelihood computers. Later work added prior probability, for a Bayesian approach. Tanner &
328 Birdsall (1958) noted that, when figuring out how a biological system does a task, it is very
329 helpful to know the optimal algorithm and to rate observed performance by its *efficiency* relative
330 to the optimum. SDT solved detection and classification mathematically, as maximum likelihood.
331 It was the classification math of the sixties. Machine learning is the classification math of today.
332 Both enable deeper insight into how biological systems classify. In the old days we used to
333 compare human and ideal classification performance (Pelli et al. 2006). Today, we also
334 compare human and machine learning. Deep learning is the best model we have today for how
335 complex systems of simple units can recognize objects as well as the brain does. Deep
336 learning, i.e. learning by multi-layered neural networks using backprop, is not just AlexNet but
337 also includes ConvNets and other architectures of trained artificial neural networks. Several labs
338 are currently comparing patterns of activity of particular layers to neural responses in various
339 cortical areas of the mammalian visual brain (Yamins et al. 2014; Khaligh-Razavi &
340 Kriegeskorte, 2014).

341 **What computer scientists can learn from psychophysics.** Computer scientists build
342 classifiers to recognize objects. Vision scientists, including psychologists and neuroscientists,
343 study how people and animals classify in order to understand how the brain works. So, what do
344 computer and vision scientists have to say to each other? Machine learning accepts a set of
345 labelled stimuli to produce a classifier. Much progress has been made in physiology and
346 psychophysics by characterizing how well biological systems can classify stimuli. The
347 psychophysical tools (e.g. threshold and signal detection theory) developed to characterize
348 behavioral classification performance are immediately applicable to characterize classifiers
349 produced by machine learning (e.g. Ziskind, Hénaff, LeCun, & Pelli, 2014; Testolin, Stoianov, &
350 Zorzi, 2017).

351 **Psychophysics.** “Adversarial” examples have been presented as a major flaw in deep neural
352 networks. These slightly doctored images of objects are misclassified by a trained network,
353 even though the doctored images have little effect on human observers. The same doctored images are
354 similarly misclassified by several different networks trained with the same stimuli (Szegedy, et
355 al., 2013). Humans too have adversarial examples. Illusions are robust classification errors. The
356 blindspot-filling-in illusion is a dramatic adversarial example in human vision. While viewing with
357 one eye, two finger tips touching in the blindspot are perceived as one long finger. If the image
358 is shifted a bit so that the fingertips emerge from the blindspot the viewer sees two fingers.
359 Neural networks lacking the anatomical blindspot of human vision are hardly affected by the
360 shift. The existence of adversarial examples is intrinsic to classifiers trained with finite data,
361 whether biological or not. In the absence of information, neural networks interpolate and so do
362 biological brains. Psychophysics, the scientific study of perception, has achieved its greatest
363 advances by studying classification errors (Fechner, 1860). Such errors can reveal “blindspots”.
364 Stimuli that are physically different yet indistinguishable are called *metamers*. The systematic
365 understanding of color metamers revealed the three dimensions of human color vision (Palmer,
366 1777; Young, 1802; Helmholtz, 1860). In recent work, many classifiers have been trained solely
367 with the objects they are meant to classify, and thus will classify everything as one of those
368 categories, even doctored noise that is very different from all of the images. It is important to
369 train with sample images that represent the entire test set.

370 **CONCLUSION**

371 Machine learning is here to stay. Deep learning is better than the “neural” networks of the
372 eighties. Machine learning is useful both as a model for perceptual processing, and as a
373 decoder of neural processing, to see what information the neurons are carrying. The large size
374 of the human cortex is a distinctive feature of our species and crucial for learning. It is
375 anatomically homogenous yet solves diverse sensory, motor, and cognitive problems. Key

376 biological details of cortical learning remain obscure, even if they ultimately preclude backprop,
377 the performance of current machine learning algorithms is a useful benchmark.

378 **ACKNOWLEDGEMENTS**

379 Thanks to Yann LeCun for helpful conversations. Thanks to Aenne Brielmann, Kaitie Holman,
380 Laura Suci, and Avi Ziskind for helpful comments on the manuscript. We thank both reviewers,
381 Nikolaus Kriegeskorte and anonymous, for many constructive suggestions. DGP was supported
382 by NIH grant R01 EY027964.

383 **REFERENCES**

- 384 Bengio, Y. (2009). Learning deep architectures for AI. *Foundations and Trends in Machine*
385 *Learning*, 2(1), 1–127.
- 386 Bengio, Y., Lee, D. H., Bornschein, J., Mesnard, T., & Lin, Z. (2015). Towards biologically
387 plausible deep learning. *arXiv preprint arXiv:1502.04156*.
- 388 Berardino, A., Laparra, V., Ballé, J., & Simoncelli, E. (2017). Eigen-Distortions of Hierarchical
389 Representations. In *Advances in Neural Information Processing Systems* (pp. 3533-3542).
- 390 Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by
391 the author). *Statistical science*, 16(3), 199-231.
- 392 Bryson, A. E., Denham, W. F., & Dreyfus, S. E. (1963). Optimal programming problems with
393 inequality constraints. *AIAA journal*, 1(11), 2544-2550.
- 394 Cadieu, C. F., Hong, H., Yamins, D. L., Pinto, N., Ardila, D., Solomon, E. A., Majaj, N. J. &
395 DiCarlo, J. J. (2014). Deep neural networks rival the representation of primate IT cortex for core
396 visual object recognition. *PLoS computational biology*, 10(12), e1003963.
- 397 Caporale, N., & Dan, Y. (2008). Spike timing-dependent plasticity: a Hebbian learning rule.
398 *Annu. Rev. Neurosci.*, 31, 25-46.
- 399 Cox, D. D., & Dean, T. (2014). Neural networks and neuroscience-inspired computer vision.
400 *Current Biology*, 24(18), R921-R929.
- 401 Crick, F. (1989). The recent excitement about neural networks. *Nature*, 337(6203), 129.
- 402 Dechter, R. (1986). Learning while searching in constraint-satisfaction-problems. In *Proceedings*
403 *of the Fifth AAAI National Conference on Artificial Intelligence* (pp. 178-183). AAAI Press.
- 404 Deng, L., & Yu, D. (2014). Deep learning: methods and applications. *Foundations and Trends®*
405 *in Signal Processing*, 7(3–4), 197-387.
- 406 Efron, B., & Hastie, T. (2016). *Computer Age Statistical Inference: Algorithms, Evidence, and*
407 *Data Science* (Vol. 5). Cambridge University Press.
- 408 Fechner, GT (1860). Elements of psychophysics. Breitkopf & Härtel, Leipzig (reprinted in 1964
409 by Bonset, Amsterdam); German translation by HE Adler (1966): Elements of psychophysics.
- 410 Fisher, R. A. (1922). The goodness of fit of regression formulae, and the distribution of
411 regression coefficients. *Journal of the Royal Statistical Society*, 85(4), 597-612.
412 doi:10.2307/2341124.
- 413 Friedman, J. H. (1998). Data mining and statistics: What's the connection? *Computing Science*
414 *and Statistics*, 29(1), 3-9

- 415 Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism
416 of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4), 193-202.
- 417 Galton, F. (1877). Typical laws of heredity. *Nature*, 15(389), 512-514.
- 418 Gauss, C.F. (1809). *Theoria motus corporum coelestium in sectionibus conicis solem*
419 *ambientium*. Hamburg: Friedrich Perthes und I. H. Besser.
- 420 Geisler, W. S. (1989). Sequential ideal-observer analysis of visual discriminations.
421 *Psychological review*, 96(2), 267.
- 422 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning* (Vol. 1). Cambridge: MIT press.
- 423 Helmholtz, H. von (1860/1925). *Handbuch der physiologischen Optik*, volume II. Leopold Voss,
424 Leipzig, third edition. Translated as *Treatise on Physiological Optics*, volume II. The Optical
425 Society of America, 1925. Edited by James P. C. Southall.
- 426 Hinton, G. E., Osindero, S., & Teh, Y. W. (2006). A fast learning algorithm for deep belief nets.
427 *Neural computation*, 18(7), 1527-1554.
- 428 Hubel, D. H., & Wiesel, T. N. (1962). Receptive fields, binocular interaction and functional
429 architecture in the cat's visual cortex. *The Journal of physiology*, 160(1), 106-154.
- 430 Hong, H., Yamins, D. L., Majaj, N. J., & DiCarlo, J. J. (2016). Explicit information for category-
431 orthogonal object properties increases along the ventral stream. *Nature neuroscience*, 19(4),
432 613.
- 433 Hung, C. P., Kreiman, G., Poggio, T., & DiCarlo, J. J. (2005). Fast readout of object identity from
434 macaque inferior temporal cortex. *Science*, 310(5749), 863-866.
- 435 Ivakhnenko, A. G. (1971). Polynomial theory of complex systems. *IEEE Transactions on*
436 *Systems, Man and Cybernetics*, 4, 364-378.
- 437 Ivakhnenko, A. G. & Lapa, V. G. (1965). *Cybernetic Predicting Devices*. CCM Information
438 Corporation.
- 439 Jarrett, K., Kavukcuoglu, K., & LeCun, Y. (2009). What is the best multi-stage architecture for
440 object recognition? In *IEEE 12th International Conference on Computer Vision*. pp. 2146-2153.
- 441 Khaligh-Razavi, S. M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised,
442 models may explain IT cortical representation. *PLoS computational biology*, 10(11), e1003915.
- 443 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep
444 convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097-
445 1105.
- 446 LeCun, Y. (1985). Une procedure d'apprentissage pour reseau a seuil asymmetrique (A
447 learning scheme for asymmetric threshold networks). In *Proceedings of Cognitiva 85*, Paris,
448 France. pp. 599-604.

- 449 LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., & Jackel, L. D.
450 (1989). Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4),
451 541-551.
- 452 LeCun, Y., Boser, B. E., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W. E., & Jackel,
453 L. D. (1990). Handwritten digit recognition with a back-propagation network. In *Advances in*
454 *neural information processing systems* (pp. 396-404).
- 455 LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to
456 document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- 457 LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.
- 458 Legendre, A. M. (1805). *Nouvelles méthodes pour la détermination des orbites des comètes*
459 (No. 1). F. Didot.
- 460 Lin, C. J. (2001). On the convergence of the decomposition method for support vector
461 machines. *IEEE Transactions on Neural Networks*, 12(6), 1288-1298.
- 462 Lowel, S., & Singer, W. (1992). Selection of intrinsic horizontal connections in the visual cortex
463 by correlated neuronal activity. *Science*, 255(5041), 209.
- 464 Majaj, N. J., Hong, H., Solomon, E. A., & DiCarlo, J. J. (2015). Simple learned weighted sums of
465 inferior temporal neuronal firing rates accurately predict human core object recognition
466 performance. *Journal of Neuroscience*, 35(39), 13402-13418.
- 467 Marcus, G. F., Pinker, S., Ullman, M., Hollander, M., Rosen, T. J., Xu, F., & Clahsen, H. (1992).
468 Overregularization in language acquisition. *Monographs of the society for research in child*
469 *development*, i-178.
- 470 Marr, D. (1982). *Vision: A computational investigation into the human representation and*
471 *processing of visual information*. San Francisco, CA: Freeman and Company.
- 472 Matloff, N. (2014). Statistics: Losing Ground to CS, Losing Image Among Students. *Revolutions*.
473 August 26, 2014. [http://blog.revolutionanalytics.com/2014/08/statistics-losing-ground-to-cs-losing-image-among-](http://blog.revolutionanalytics.com/2014/08/statistics-losing-ground-to-cs-losing-image-among-students.html)
474 [students.html](http://blog.revolutionanalytics.com/2014/08/statistics-losing-ground-to-cs-losing-image-among-students.html)
- 475 Marblestone, A. H., Wayne, G., & Kording, K. P. (2016). Toward an integration of deep learning
476 and neuroscience. *Frontiers in computational neuroscience*, 10.
- 477 Mazzone, P., Andersen, R. A., & Jordan, M. I. (1991). A more biologically plausible learning rule
478 for neural networks. *Proceedings of the National Academy of Sciences*, 88(10), 4433-4437.
- 479 Minsky, M., & Papert, S. (1988). *Perceptrons: An Introduction to Computational Geometry*.
480 Cambridge, MA: MIT press.
- 481 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by
482 learning a sparse code for natural images. *Nature*, 381(6583), 607.

- 483 Palmer, G. (1777). *Theory of Colour and Vision*, London: Leacroft.
- 484 Pearson, K., Yule, G. U., Blanchard, N., & Lee, A. (1903). The law of ancestral heredity.
485 *Biometrika*, 2(2), 211-236.
- 486 Pelli, D. G., Burns, C. W., Farell, B., & Moore-Page, D. C. (2006). Feature detection and letter
487 identification. *Vision research*, 46(28), 4646-4674.
- 488 Peterson, W. W. T. G., Birdsall, T., & Fox, W. (1954). The theory of signal detectability.
489 *Transactions of the IRE professional group on information theory*, 4(4), 171-212.
- 490 Ranzato, M. A., Huang, F. J., Boureau, Y. L., & LeCun, Y. (2007). Unsupervised learning of
491 invariant feature hierarchies with applications to object recognition. In *2007 IEEE conference on*
492 *computer vision and pattern recognition*, pp. 1-8.
- 493 Ranzato, M. A., Poultney, C., Chopra, S., & LeCun, Y. (2007). Efficient learning of sparse
494 representations with an energy-based model. In *Proceedings of NIPS*.
- 495 Riesenhuber, M., & Poggio, T. (1999). Hierarchical models of object recognition in
496 cortex. *Nature neuroscience*, 2(11), 1019-1025.
- 497 Rojas, R. (1996). A short proof of the posterior probability property of classifier neural networks.
498 *Neural Computation*, 8(1), 41-43.
- 499 Rosenblatt, F. (1958), The Perceptron: A Probabilistic Model for Information Storage and
500 Organization in the Brain, *Psychological Review*, 65, 6, pp. 386–408.
- 501 Rubinov, M. (2015). Neural networks in the future of neuroscience research. *Nature Reviews*
502 *Neuroscience*.
- 503 Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-
504 propagating errors. *Nature*, 323, 533-536. doi:10.1038/323533a0.
- 505 Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. *IBM*
506 *Journal of research and development*, 3(3), 210-229.
- 507 Samuel, A. L. (1967). Some studies in machine learning using the game of checkers. II—recent
508 progress. *IBM Journal of research and development*, 11(6), 601-617.
- 509 Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural networks*, 61,
510 85-117.
- 511 Sejnowski, T. J., & Rosenberg, C. R. (1987). Parallel networks that learn to pronounce English
512 text. *Complex systems*, 1(1), 145-168.
- 513 Seung, H. S., & Sompolinsky, H. (1993). Simple models for reading neuronal population codes.
514 *Proceedings of the National Academy of Sciences*, 90(22), 10749-10753.

- 515 Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., & Fergus, R.
516 (2013). Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- 517 Tanner Jr, W. P., & Birdsall, T. G. (1958). Definitions of d' and η as psychophysical measures.
518 *The Journal of the Acoustical society of America*, 30(10), 922-928.
- 519 Testolin, A., Stoianov, I., & Zorzi, M. (2017). Letter perception emerges from unsupervised deep
520 learning and recycling of natural image features. *Nature Human Behaviour*, 1(9), 657.
- 521 Turing, A.M. (1953). 'Digital computers applied to games'. in 'Faster than thought', ed. B.V.
522 Bowden, London 1953. Published by Pitman Publishing.
- 523 Ullman, S., Assif, L., Fetaya, E., & Harari, D. (2016). Atoms of recognition in human and
524 computer vision. *Proceedings of the National Academy of Sciences*, 113(10), 2744-2749.
- 525 Vapnik, V. (2013). *The nature of statistical learning theory*. Springer Science & Business Media.
- 526 Weber, E. H. (1834/1996). *EH Weber on the tactile senses*. Psychology Press. Translated by
527 Helen E. Ross from E.H Weber (1834) *De Tactu*.
- 528 Werbos, P. (1974). *Beyond regression: New tools for prediction and analysis in the behavioral*
529 *sciences*. PhD thesis, Harvard University.
- 530 Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014).
531 Performance-optimized hierarchical models predict neural responses in higher visual
532 cortex. *Proceedings of the National Academy of Sciences*, 111(23), 8619-8624.
- 533 Young, T. (1802). The Bakerian Lecture. On the theory of light and colours, *Philosophical*
534 *Transactions of the Royal Society of London* 92, 12-48. doi: 10.1098/rstl.1802.0004
- 535 Yule, G. U. (1897). On the theory of correlation. *Journal of the Royal Statistical Society*, 60(4),
536 812-854.
- 537 Zeiler, M. D., & Fergus, R. (2013). Visualizing and understanding convolutional networks. arXiv
538 preprint arXiv:1311.2901.
- 539 Ziskind, A.J., Hénaff, O., LeCun, Y., & Pelli, D.G. (2014) The bottleneck in human letter
540 recognition: A computational model. Vision Sciences Society, St. Pete Beach, Florida, May 16-
541 21, 2014, 56.583. <http://f1000.com/posters/browse/summary/1095738>

542

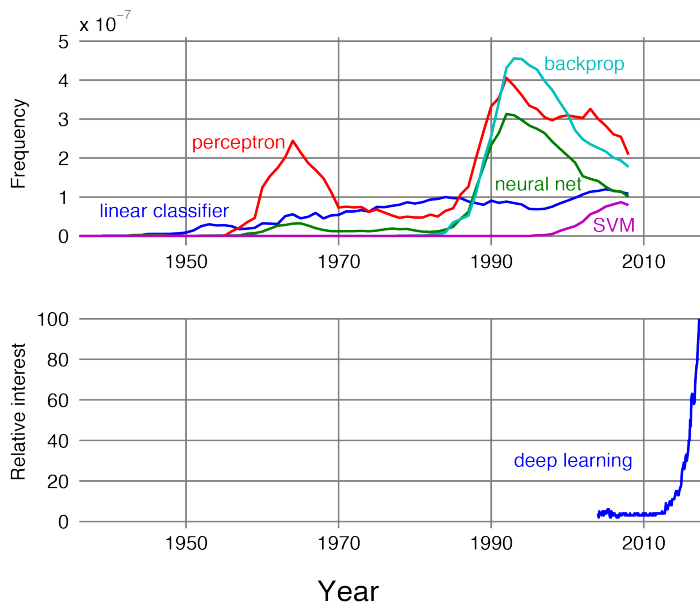


Figure 1. Top: The frequency of appearance of each of five terms — “linear classifier”, “perceptron”, “support vector machine”, “neural net” and “backprop” — in books indexed by Google in each year of publication. Frequency is reported as a fraction of all instances of that number of words (1,2, or 3) normalized by the number of books published that year (ngram / year / books published). The figure was created using Google’s n-gram viewer (<https://books.google.com/ngrams>), which contains a yearly count of n-grams found in sources printed between 1500 and 2008. **Bottom:** Numbers represent worldwide search interest relative to the highest point on the chart for the given year for the term “deep learning” (as reported by <https://trends.google.com/trends/>).

Figure 1. Top: The frequency of appearance of each of five terms — “linear classifier”, “perceptron”, “support vector machine”, “neural net” and “backprop” — in books indexed by Google in each year of publication. Frequency is reported as a fraction of all instances of that number of words (1,2, or 3) normalized by the number of books published that year (ngram / year / books published). The figure was created using Google’s n-gram viewer (<https://books.google.com/ngrams>), which contains a yearly count of n-grams found in sources printed between 1500 and 2008. **Bottom:** Numbers represent worldwide search interest relative to the highest point on the chart for the given year for the term “deep learning” (as reported by <https://trends.google.com/trends/>).

