# Deep-Learning-Based Multivariate Pattern Analysis (dMVPA): A Tutorial and a Toolbox

1  **Karl M. Kuntzelman[1,2], Jacob M. Williams[3], Phui Cheng Lim[1,4], Ashok Samal[3], Prahalada K.**
2  **Rao[5], & Matthew R. Johnson[1,4]\***

3  [1]Center for Brain, Biology and Behavior, University of Nebraska-Lincoln, Lincoln, NE, USA

4  [2]Office of Technology Development and Coordination, National Institute of Mental Health, National
5  Institute of Health, Bethesda, MD, USA

6  [3]Department of Computer Science and Engineering, University of Nebraska-Lincoln, Lincoln, NE,
7  USA

8  [4]Department of Psychology, University of Nebraska-Lincoln, Lincoln, NE, USA

9  [5]Department of Mechanical and Materials Engineering, University of Nebraska-Lincoln, Lincoln,
10  NE, USA

11  **\* Correspondence:**
12  Matthew R. Johnson
13  matthew.r.johnson@gmail.com

14  **Keywords: deep learning, cognitive neuroscience, machine learning, EEG, fMRI, neural**
15  **networks, Python, MVPA**

16  **Abstract**

17  In recent years, multivariate pattern analysis (MVPA) has been hugely beneficial for cognitive
18  neuroscience by making new experiment designs possible and by increasing the inferential power of
19  functional magnetic resonance imaging (fMRI), electroencephalography (EEG), and other
20  neuroimaging methodologies. In a similar time frame, "deep learning" (a term for the use of artificial
21  neural networks with convolutional, recurrent, or similarly sophisticated architectures) has produced
22  a parallel revolution in the field of machine learning and has been employed across a wide variety of
23  applications. Traditional MVPA also uses a form of machine learning, but most commonly with
24  much simpler techniques based on linear calculations; a number of studies have applied deep learning
25  techniques to neuroimaging data, but we believe that those have barely scratched the surface of the
26  potential deep learning holds for the field. In this paper, we provide a brief introduction to deep
27  learning for those new to the technique, explore the logistical pros and cons of using deep learning to
28  analyze neuroimaging data – which we term "deep MVPA," or dMVPA – and introduce a new
29  software toolbox (the "Deep Learning In Neuroimaging: Exploration, Analysis, Tools, and
30  Education" package, DeLINEATE for short) intended to facilitate dMVPA for neuroscientists (and
31  indeed, scientists more broadly) everywhere.

32  **Word Count:** 11929

33  **Figures:** 3

34  **Tables:** 1

## 1    Introduction

35

36    Although the roots of cognitive neuroscience date to the 1920s (the advent of
37    electroencephalography, EEG; Berger, 1929), the modern neuroimaging era began in the mid-1990s,
38    with the development of functional magnetic resonance imaging (fMRI) methodology and the
39    increasingly widespread availability of (affordable) desktop computing workstations powerful
40    enough to process fMRI datasets. In those days, data analysis was primarily limited to univariate
41    investigations such as event-related potentials (ERPs) in EEG and univariate general linear model
42    (GLM) analyses aimed at detecting "blobs" of activation with fMRI (as well as differences in
43    activity, e.g. between experimental conditions, within such blobs)[1]. However, the march of progress
44    towards ever-more sophisticated models of brain function and the testing of ever-more refined
45    hypotheses has created a demand for corresponding improvements in analysis techniques.

46    Thus, somewhat more recently (beginning in the early-to-mid-2000s), a second age in neuroimaging
47    analysis arose with the advent of multivariate pattern analysis (MVPA; Haxby et al., 2001; Haxby et
48    al., 2014). Rather than focusing on whether a certain cognitive event elicits activity in a particular
49    cluster of fMRI voxels (or a voltage peak at a particular temporal latency with ERP), MVPA is
50    instead concerned with how a neural pattern or multivariate "brain state" comprising multiple voxels
51    (fMRI) or electrode/timepoint combinations (EEG) might collectively correspond to a certain
52    cognitive event or state. Numerous MVPA variations exist, including those based on correlation
53    (either Pearson or rank-based; Haxby et al., 2001), support vector machines (SVMs; De Martino et
54    al., 2008; Dosenbach et al., 2010), logistic regression (Akama et al., 2012), sparse multinomial
55    logistic regression (SMLR; Kohler et al., 2013; Krishnapuram et al., 2005), naïve Bayes classifiers
56    (Kassam et al., 2013), and more. Many of these techniques concern classification of brain patterns
57    into discrete cognitive states, whereas others examine different aspects of the data (e.g., overall
58    similarity between brain patterns; Xue et al., 2010; Lim et al., 2019) without explicit categorization,
59    but all of them represent increases in mathematical and conceptual sophistication over univariate
60    techniques. Importantly, when compared to earlier univariate techniques, MVPA has enabled us to
61    examine in a much more nuanced fashion how brain activity patterns encode mental states.

62    Although traditional MVPA techniques are substantially more advanced than univariate techniques,
63    they are nonetheless still fairly simple, both mathematically and conceptually. Traditional MVPA is a
64    form of machine learning (ML), but it is among the simplest forms; most MVPA approaches use
65    straightforward linear mathematical models. This comparative simplicity certainly confers
66    advantages – for example, faster computation times than more complex techniques (with some
67    caveats[2]), and a generally lower risk of "overfitting"[3]. However, simpler mathematical formulations

2

---

[1] Although most of our discussion focuses on fMRI and EEG, as those are the most common techniques in our field of cognitive neuroscience, most points should translate well to related technologies like structural MRI, magnetoencephalography (MEG), or electrocorticography (ECoG), and even to less closely related methods such as extracellular recordings (e.g., from rodents or nonhuman primates).

[2] For example, SVMs may take inordinately long to converge on extremely high-dimensional datasets that are handled more easily by deep neural networks. As discussed later, deep networks also have better support for GPU-based parallelization than simpler linear techniques, which can offset their computational costs.

[3] The creation of a predictive model that is highly customized to the data used to train the model, but generalizes poorly to new datasets that do not perfectly match the idiosyncrasies of the training data; a significant concern in ML. A good analogy is a bespoke garment perfectly tailored to the contours of a specific individual, which would fit him/her perfectly

68   are necessarily limited in what we call "informational resolution" – the specificity of the neural
69   patterns and cognitive states that they are able to capture.

70   How much informational resolution is required to glean as much about brain function as is possible
71   using current neuroimaging technology? The answer is hard to pin down, partly because it is difficult
72   to establish firm estimates of the "noise ceiling"[4] for these techniques. As neuroimagers, we often
73   complain that our techniques are "noisy," but with proper usage, the signal-to-noise ratios of EEG
74   and fMRI are really rather high, when considering only measurement noise from the instruments
75   themselves and the surrounding physical environment. Of significantly greater concern are "noise"
76   sources such as subject head/body motion, physiological artifacts (cardiac, respiratory, muscular,
77   etc.), and cognitive artifacts (distraction, poor understanding of instructions, falling asleep). Noise
78   ceilings for certain analytic techniques and datasets can be estimated (Kay et al., 2008; Nili et al.,
79   2014), but ultimately they will depend on which data components are considered "noise"; aside from
80   the noise that arises from the physics of the measurement itself, other biological and subject-driven
81   artifacts have some hope of being detected, modeled, and/or removed. And, much like the signal
82   components we actually care about (i.e., those related to our experimental questions), our ability to
83   detect and account for noise depends largely on the sophistication of our analytic techniques.

84   What we do know is that the brain is a highly complex, highly nonlinear system (Koch & Laurent,
85   1999; Sporns et al., 2000; Buzsaki & Mizuseki, 2014), and the addition of noise sources that are also
86   complex and nonlinear makes brain data no easier to analyze and interpret. Although the limits of the
87   usefulness of traditional MVPA, with its relatively low informational resolution, have not yet been
88   reached, those limits do loom on the horizon. As the size of neuroscience data continues to grow[5],
89   traditional MVPA's limitations become ever more apparent. It is a statistical truism that more
90   complex analytic models, with more parameters to fit, allow us to account for a greater proportion of
91   a dataset's variance, but they also require larger input data to estimate their parameters reliably. Yet
92   the sizes of many contemporary datasets are now such that they can potentially accommodate
93   significantly more sophisticated statistical models than traditional MVPA, with greater power to
94   identify, extract, and distinguish noise sources and signals of interest. Thus, we believe it is time for
95   cognitive neuroscience and related fields to place increased emphasis on developing, exploring, and
96   using more sophisticated techniques, and on producing tools that can be used to perform that
97   exploration more effectively and efficiently.

## 1.1 The case for deep learning

99    There are numerous potential analytic methods of greater complexity and sophistication than
100   traditional MVPA. One class of ML techniques that has been gaining popularity, and the one we
101   endorse in this paper, is "deep learning." Deep learning, briefly defined, refers to the use of artificial

---

but look terrible on most others. Conversely, an off-the-rack outfit with a simpler design would fit many individuals of roughly similar proportions reasonably well.

[4] Informally defined, the best we might be expected to do in using statistics to explain variance in the data, accounting for the fact that a certain amount of unexplainable variance, aka noise, will always exist.

[5] E.g., from better spatiotemporal resolution due to technological improvements; from increasingly large sample sizes, particularly from big-data initiatives such as the Human Connectome Project (Van Essen et al., 2013) and OpenNeuro (formerly OpenfMRI; Poldrack et al., 2013); and simply from the ongoing accumulation of data stockpiles from many years' worth of research studies.

102 neural networks (ANNs), typically with recurrent and/or convolutional architectures, that are more
103 complex, flexible, and powerful than both earlier generations of ANN architectures and the
104 techniques used for traditional MVPA. In the last few years, such deep neural networks (DNNs) have
105 been used increasingly heavily in a number of fields that employ ML for all kinds of purposes. Such
106 usage includes an ever-growing collection of studies in human neuroscience and related disciplines,
107 although a relatively small proportion have been devoted to neuroimaging analysis, and fewer still
108 devoted to decoding cognitive states from functional measurements of brain activity, which is a topic
109 of great interest to many. We believe the studies so far represent only the tip of the proverbial iceberg
110 in terms of what is achievable by using DNNs to analyze neuroscience datasets. In fact, we believe
111 deep learning has the potential to perform most of the tasks for which traditional MVPA is typically
112 employed, but with greater speed, flexibility, and power, and thus we advocate for the more
113 widespread use of what we call "deep MVPA," or dMVPA for short.

114 To achieve more widespread adoption of deep learning in the neurosciences, notable challenges to
115 confront include 1) a relatively low level of knowledge/awareness of these techniques, and 2)
116 insufficient availability of software tools to make dMVPA as approachable as traditional MVPA. In
117 this paper we address the first challenge by providing a brief review of deep learning techniques,
118 including how they can be used in neuroscience investigations, and the pros and cons of dMVPA
119 versus traditional MVPA. We address the second challenge by introducing a new Python-based
120 software toolbox (the "Deep Learning In Neuroimaging: Exploration, Analysis, Tools, and
121 Education" package; DeLINEATE for short) that builds upon previous DNN and MVPA tools and
122 aims to make dMVPA more approachable and efficient for other researchers.

123 **2    dMVPA: A tutorial**

124 **2.1    A brief history of neural networks**

125 The techniques we now collectively call "deep learning" are generally extensions of older "shallow"
126 ANNs, which are significantly less complex and powerful than DNNs but not much different in their
127 basic principles. The concept behind all ANNs originates from a highly abstracted view of non-
128 artificial neural networks, i.e., the biological nervous system (Figure 1A). In this framework, most
129 implementation details are stripped away, and what remains is the basic idea of a network of simple
130 computational units ("neurons") that receive input (which can typically be excitatory or inhibitory),
131 perform an operation on their inputs (typically some variation on summation), and produce an output
132 (typically a single value analogous to an action potential or a firing rate), which might then serve as
133 input to one or more downstream neurons[6]. The original and simplest case is the McCulloch-Pitts
134 neuron (McCulloch & Pitts, 1943; Figure 1B), a processing unit whose input and output values are
135 exclusively binary (0 or 1). The McCulloch-Pitts neuron sums its inputs, compares the sum to some
136 threshold value, and outputs a 1 ("action potential") or 0 according to whether the sum exceeds the
137 threshold. Although a pioneering idea and an interesting (if highly simplified) early model of neural

4

[6] It should be noted up front that the artificial "neurons" used in ML applications bear about as much resemblance to real neurons as a paper airplane bears to a commercial airliner. In both cases, the barest core principles are similar between the pared-down model and the real thing, but little else. However, despite the low resemblance, ANNs can still be extremely useful tools for ML and data processing. Readers are nonetheless cautioned to be as circumspect about over-aggressive comparisons between artificial and real neural networks as they would about buying transatlantic tickets on paper airplanes.

138 information processing, McCulloch-Pitts neurons can only implement a limited set of functions and
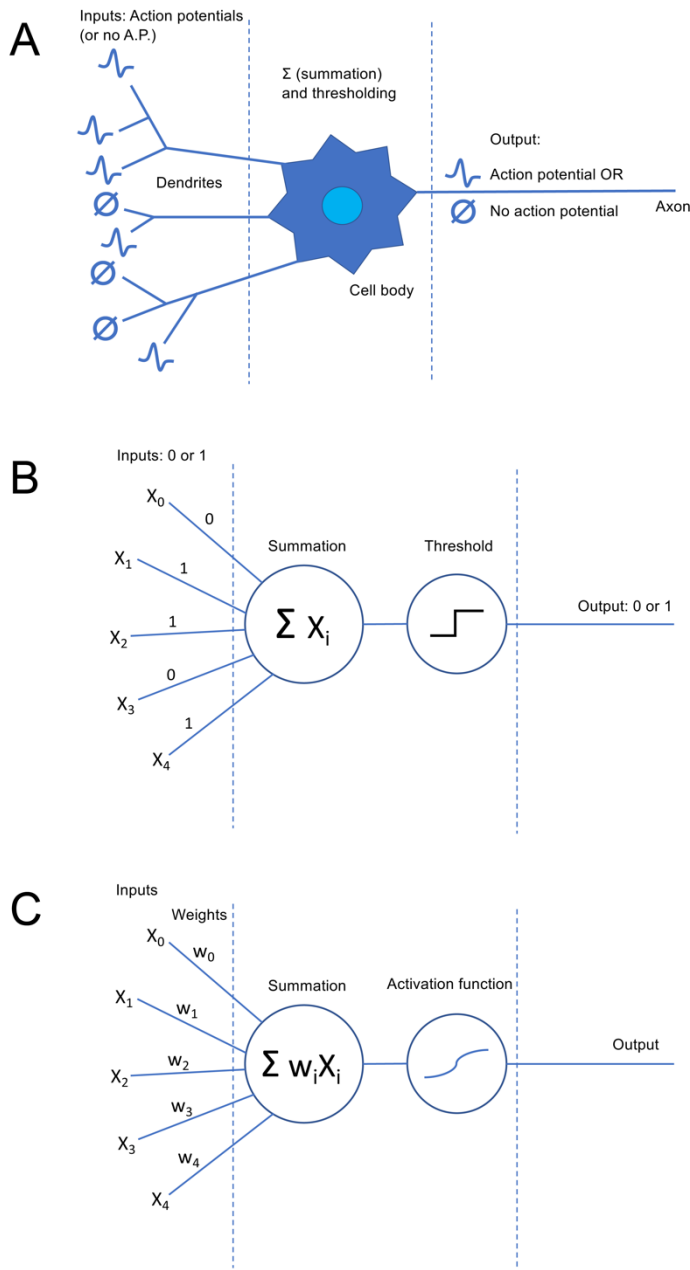139 are thus not considered very useful for modern ML applications.



**Figure 1. Comparison of biological and artificial neural models.**

(A) A simplified "textbook" model of a biological neuron. Inputs come in via the dendrites in the form of action potentials (or the lack thereof). The inputs are summed in the cell body (soma) and, if the threshold voltage is reached, the cell produces an action potential as output that is delivered via the cell's axon. (B) The original and simplest version of an artificial neuron model, the McCulloch-Pitts neuron. Similar to the biological neuron, inputs (xi) and outputs are binary (although we now know this to be an oversimplified view of biological neurons). Inputs are summed and the result passed to a thresholding function; if the threshold is met, an output of 1 is produced, and otherwise the output is 0. (C) A perceptron, a more sophisticated revision of the McCulloch-Pitts neuron that has an important place in modern artificial neural networks. The concept of trainable weights (wi; mimicking biological potentiation at synapses) is introduced, and inputs are now multiplied by their corresponding weight before summation. In addition, in contemporary perceptron models, the threshold function can be replaced by any arbitrary function, called the "activation function." Popular activation functions like the hyperbolic tangent may still act largely like thresholding functions, but with the ability to deliver graded rather than strictly binary output values.

140 A few years later, though, ANNs took a significant step forward when Rosenblatt (1958)
141 incorporated Hebb's theoretical views on the strengthening and weakening of synaptic connections
142 (Hebb, 1949) into a McCulloch-Pitts-like unit that came to be called the *perceptron*[7]. In its simplest
143 form, a perceptron (Figure 1C) is largely identical to a McCulloch-Pitts neuron with one critical

5

[7] As originally conceived, "perceptron" referred to a more complex network of units that could be implemented in a physical machine to produce artificial vision, hence the name. However, the most salient feature that researchers latched onto was the structure of the neural units, and via synecdoche "perceptron" came to be the name of such a unit, so there is some degree of fuzziness around nomenclature and definitions. Here, we use the contemporary sense of "perceptron" to refer to the architecture of an artificial neural unit, rather than the original plan for the physical perceptron machine.

144 addition: Each input is now associated with a "synaptic weight" (often denoted $w_0$, $w_1$, etc.) that
145 determines whether it is excitatory or inhibitory and how strongly it influences the output.
146 Summation is then performed on the inputs after they have been multiplied by their respective
147 weights. Statistically inclined readers may recognize this as not-dissimilar-to a regression model,
148 particularly logistic regression; to conceptually convert between a multiple regression model and a
149 perceptron, simply rename the weights from the $\beta_i$ typically used in regression equations to $w_i$ and
150 pass the regression output through a thresholding function, logistic function (to essentially replicate
151 logistic regression), or other function as desired[8]. This function is known as the artificial neuron's
152 *activation function*; activation functions are a key feature of contemporary ANN designs, and there
153 are many options to choose from.

154 In most important ways, the perceptron-like artificial neural units used in some DNNs today are not
155 substantially different than the classical perceptrons discussed by Minsky and Papert (1969) in their
156 seminal book some 50 years ago. Yet the original perceptron architectures retained many of the
157 McCulloch-Pitts neurons' limitations and still had significant constraints on the classes of problems
158 they could solve. The key developments that distinguish the powerful deep learning techniques of
159 today from the toy models of the past are 1) improved methods for establishing what the proper
160 synaptic weights should be for a given dataset/problem, i.e., *training* the neural network, and 2) new
161 and ultimately better ways of digitally connecting groups of artificial neurons together into more
162 complex structures, i.e., improved ANN *architectures*[9].

### 2.2 Training algorithms and neural network architectures

164 The earliest ANN architectures were very simple indeed; either a single artificial neuron or, in the
165 next major architectural advance after that, a *layer* of such units. In this latter (still very simple)
166 architecture, the units are *fully connected*, meaning that each unit receives a copy of each possible
167 input value (see Figure 2A). Note that in this figure, as in many neural network diagrams, inputs and
168 outputs are represented as "layers" of a sort, but there is only one true layer of computational units[10].
169 If the ANN is meant to calculate a classification problem (a common application), the outputs are
170 typically assumed to each correspond to one of the possible classes, and are interpreted in a winner-
171 take-all fashion (i.e., for a given set of input data, whichever output value is highest is interpreted as
172 the network's prediction of the class that the input data belong to). Although the transition from
173 single-neuron to single-layer architectures laid a critical foundation for later work, single layers of
174 perceptrons were soon shown not to be terribly useful as artificial intelligence agents, no matter what
175 their synaptic weights were or how those weights were determined. As Minsky and Papert
176 demonstrated in *Perceptrons* (1969), it is mathematically impossible for any single-layer perceptron

6

[8] The concept of the perceptron is also somewhat looser than the McCulloch-Pitts neuron regarding whether inputs and outputs are constrained to be binary or can be continuously valued, and regarding what kind of thresholding or other function the summed inputs are passed through in order to create the output.

[9] Not to mention the ~billion-fold increase in computational power (IBM 704 at 12,000 flops versus a recent desktop GPU at ~11 *tera*flops, for an NVIDIA GeForce GTX 1080 Ti) that helps to make such sophisticated architectures viable.

[10] In a biological neural network, one might relate these to a layer of dendrites, a layer of cell bodies, and a layer of axons, but all of those together would comprise a single layer of neurons.

177 network – no matter how many units are in it – to perform certain fundamental computational
178 operations[11].

179 This revelation may not seem surprising in retrospect; after all, a single layer of neurons, all receiving
180 the same inputs, is not a very viable architecture for a biological neural network either. Still, it was
181 enough to significantly dampen enthusiasm for ANN research for over a decade. Although adding
182 another layer of computational units (known as a *hidden layer*) would allow the network to maintain
183 an intermediate representation of the input and enable more complex operations[12], the algorithms
184 available for training single-layer perceptron networks could not be readily extended to multi-layer
185 architectures. In the 1980s, however, interest was reignited with the popular (re-)discovery of the
186 *backpropagation* algorithm (or simply *backprop*, to its friends). This algorithm was known and even
187 applied to ANNs previously (Linnainmaa, 1970; Werbos, 1974), but it did not reach mainstream
188 awareness until the publication of Rumelhart and colleagues' (Rumelhart et al., 1986a; Rumelhart et
189 al., 1986b) seminal formulations of it. Backprop proved to be a highly robust method for training
190 ANNs across many applications, and is still the dominant training algorithm in use today.

191 The main principle behind backprop is to take any errors made by the network during training and
192 propagate responsibility for them from the output layer (where the error is assessed, by comparing
193 the network's decision to the known correct decision[13]) backwards through the network towards the
194 input layer, penalizing the synaptic weights most responsible for the error along the way. It is
195 analogous to the human behavior encapsulated by the vernacular phrase, "Shit rolls downhill." For
196 example, imagine that a CEO – the final decision-maker in her company's chain of command –
197 makes a decision that loses the company money. She turns to her immediate inferiors and doles out
198 punishment to them proportional to how influential they were in guiding that decision, and decides to
199 trust those influential individuals less in the future. In turn, each of those upper-level managers
200 passes along the punishment and distrust they have received to their immediate inferiors, again
201 proportional to their influence on the upper managers' actions, and so on down the corporate
202 hierarchy. In this way, one hopes that the next time a similar decision is faced, the shift in influences
203 and communication channels throughout the hierarchy will produce a better outcome.

204 The advent of effective backprop-based training for ANNs reignited interest in them for a time, and
205 backprop-trained ANNs were found to perform admirably in a number of ML domains. Still, before
206 long, interest waned again, as neural nets with many hidden layers were found to present
207 mathematical difficulties for backpropagation algorithms, and complex networks also took a long
208 time to train on the CPUs of the era. Concurrently, the 1990s also saw the development of promising
209 alternative ML algorithms, most notably the modern incarnation of support vector machines (SVMs;
210 Cortes & Vapnik, 1995; Boser et al., 1992). SVMs were easier to work with than ANNs and
211 performed nearly equivalently (or even better) in many problem domains of the time. Thus, when
212 traditional MVPA techniques arose in neuroimaging in the 2000s, it is unsurprising that SVMs and

[11] Put more formally, single-layer networks cannot solve problems that are not linearly separable, which famously includes the relatively simple XOR function. (For binary inputs A and B, respond "yes" if A is true and B is false, or if A is false and B is true, but respond "no" if A and B have the same value.)

[12] Including XOR and many others.

[13] As backprop is performed by comparing the performance of a network on a training dataset against an already-known ground truth for that dataset, it is thus considered a form of *supervised learning*, in ML parlance.

213  other similarly robust linear classification algorithms, well-suited to the mid-sized datasets of the
214  time, dominated within that emerging field.

215  ## 2.3   The deep learning Renaissance

216  Research interest in ANNs experienced another upswing, which has continued to the present,
217  beginning around 2006. This rebirth happened for several reasons, including: 1) solutions to some of
218  the technical and mathematical problems that had plagued networks with complex, many-layered
219  architectures (Hinton et al., 2006); 2) methods for training ANNs on desktop workstations using the
220  GPU instead of the CPU, producing speed improvements of up to ~70x (Raina et al., 2009); 3) the
221  advent of the so-called "Big Data" era, which provided the larger datasets required to adequately
222  train more complex neural architectures; and 4) the re-branding of neural net research as "Deep
223  Learning," which, despite being more public relations than true substance, still likely helped ignite
224  new interest in a field formerly seen as relatively tired and unpopular. Since this Renaissance began,
225  there have naturally been several key architectural and methodological developments[14]. However,
226  these newer architectures are still trained and used similarly to the older, simpler networks described
227  above, and the variations are not too difficult to comprehend once one understands the fundamental
228  concepts and terminology behind ANNs.

229  During the early days of this revival, deep learning research had a number of notable successes,
230  including advances in speech recognition, natural language processing, computer vision, financial
231  fraud detection, and more. Large technology companies, who had access to Big Data and financial
232  motivations for finding better ways to process it, also had their interest piqued. Thus they began to
233  invest in deep learning research themselves, including developing improved software tools (for
234  example, the TensorFlow toolbox, developed primarily at Google, and PyTorch, developed primarily
235  at Facebook). These tools typically rely on lower-level driver and software library support for GPU-
236  based computation, most notably NVIDIA's CUDA libraries for general GPU-accelerated computing
237  and their cuDNN framework, built atop CUDA, specifically for DNN applications[15]. Although the
238  use of such tools has exploded in the technology sphere and in basic computer science research,
239  adoption in other areas, such as cognitive neuroscience, has been slower. This lag can partly be
240  attributed to fundamental limitations and difficulties of DNN-based data analysis (e.g., potential for
241  overfitting), but another large factor is the lack of higher-level software tools that make it convenient
242  for neuroscience researchers to implement dMVPA without needing to write large amounts of their
243  own code. And, when better software tools exist, it will be more efficient to explore the space of
244  possibilities and limitations of dMVPA. In short, neuroscience and related fields need more software
245  tools that match, or exceed, the versatility and ease-of-use of existing traditional MVPA tools. This is
246  the goal of the DeLINEATE toolbox (Deep Learning In Neuroimaging: Exploration, Analysis, Tools,
247  and Education), which we introduce below.

---

[14] E.g., the use of ReLU activation functions (Maas, Hannun, and Ng, 2013); new approaches to regularization (Zeiler & Fergus, 2013); and other architectural elements that were available earlier became more prominently used, once sufficient data and computing power existed to use them more effectively (e.g., convolutional network layers).

[15] However, alternatives for other GPU architectures do exist, such as the CoreML library used in Apple devices, which use primarily non-NVIDIA GPUs.

248    ## 2.4    Pros, cons, and caveats of dMVPA

249    *Pro: Potentially greater suitability for complex, many-featured datasets*. As discussed earlier, one
250    great promise of dMVPA is the potential to unearth more fine-grained patterns in neuroscience data
251    than the simpler (and commonly linear) techniques of traditional MVPA. However, a fundamental
252    principle of statistics is that more powerful (i.e., more complex) models require more parameters[16],
253    and reliably estimating more parameters requires larger input datasets. Hence, why deep learning and
254    Big Data are commonly associated with each other. Unlike, say, the Google Images team, most
255    neuroscientists are unfortunately not swimming in training data for sophisticated machine learning;
256    neuroscience data are frequently "Big," but more from features[17] than from number of examples[18]. Of
257    course, in deep learning (and most statistical analyses), the inverse situation is usually more
258    desirable: A relatively large ratio of examples to features.

259    Potential solutions to the too-many-features problem include finding ways to intelligently select
260    (*feature selection*) or algorithmically condense[19] the feature set. However, beyond those options,
261    most traditional MVPA techniques do not many choices for constraining the feature set, and in
262    particular lack any built-in ability to take the structure of the input data into account. This is
263    unfortunate because neuroscience data[20] tend to be highly structured (temporally and spatially) in
264    ways that could be informative for MVPA[21]. DNNs, on the other hand, have numerous potential
265    architectural configurations that can be optimized to take advantage of known structure in the input
266    data. Most notably, certain types of ANN layers (e.g., convolutional layers) can handle multi-
267    dimensional input data, whereas traditional MVPA's linear classifiers typically just vectorize multi-
268    dimensional inputs. Thus, dMVPA makes it possible to design customized classifiers that are more
269    suited to a particular shape/dimensionality of input data.

270    *Caveat*. Having more architectural options for structuring and condensing complex input data also
271    leads to a paradox of choice; how can one possibly decide on the best DNN architecture for a given
272    dataset? Unfortunately, dMVPA is still a young field, and we are still working on establishing good
273    heuristics for network architectures to handle many-featured datasets. Also unfortunately, this is not
274    one of those methodological choices where differences between options can be chalked up to

9

[16] "Parameters" used in the statistical sense, i.e., numeric values that need to be estimated.

[17] In the machine learning sense; for example, the number of voxels in a trial of fMRI data or the number of (electrodes ×
timepoints) in a trial of EEG data.

[18] Also in the machine learning sense, i.e., instances of a set of features that can be assigned a category label. In
psychology and neuroscience, such "examples" are generally called "trials" (e.g., of a cognitive task), although in some
cases examples may correspond to experimental subjects – an even more limited resource.

[19] For example, in techniques like elastic nets (Zou & Hastie, 2005) or SMLR, which use regularization or similar tricks
to reduce the number of predictor features.

[20] Again, our discussion focuses on neuroscience data, but these techniques, lessons, and software tools can readily be
translated to related (or even not-so-related) research fields with similarly-structured datasets and classification problems.

[21] For example, it may be useful to condense several spatially adjacent EEG electrodes with similar waveforms into a
single data channel. Or, if trying to classify whether a subject is viewing faces or houses, to construct a feature detector
that is sensitive to a certain voltage peak (say, the N170; Bentin 1996) but time-invariant within a ~20ms window, to
account for trial-to-trial latency variability.

275  rounding error; the wrong dMVPA architecture may completely fail to perform above chance in
276  situations where a superior architecture classifies the data fairly accurately.

277  *Con: Many potential types of analysis architecture; many of these carry an increased danger of*
278  *overfitting.* Most conventional MVPA techniques (SVM, SMLR, etc.) have a relatively small number
279  of hyperparameters[22] to adjust, and those hyperparameters can often either be left at default values or
280  automatically estimated by the algorithm without serious adverse effects on performance. In contrast,
281  the number of possible hyperparameters to adjust in dMVPA is effectively infinite. These
282  hyperparameters include the number of layers in the network, the number of units in each layer, the
283  type of each layer[23], and any number of additional layer-type-specific hyperparameters that can be
284  separately specified for each layer. Thus, even choosing a starting point for how to construct a
285  dMVPA model can be daunting for inexperienced researchers (and experienced ones, too).
286  Furthermore, thanks to the No Free Lunch (NFL) theorem(s) (Wolpert & Macready, 1997; Shalev-
287  Shwartz & Ben-David, 2014), we know that no estimation- or optimization-based analysis technique
288  will be optimal for every dataset or problem domain, and therefore it is impossible to know *a priori*
289  whether a given analysis technique will be optimal for a particular problem. Put another way, if we
290  knew in advance that a particular analysis technique *were* optimal for our problem, then that
291  technique would necessarily be exquisitely tailored to the problem – which means we would
292  essentially already know the structure of the data perfectly, which obviates the need to conduct the
293  analysis.

294  Compounding the problem, there is no real upper limit, other than available computing power, to
295  how complex dMVPA models can be allowed to grow[24]. For the current status quo of neuroscience
296  data, most possible dMVPA models would be far too complex; many would even contain more
297  parameters to estimate than there are data points in the input set! It would be inaccurate to say these
298  models would fit the data poorly; rather, they would fit the training data *too* well. It is not uncommon
299  to see a complex dMVPA model effectively memorize its training data, producing perfect
300  classification of the training dataset but extremely poor generalization to a test dataset – the classic
301  problem of overfitting.

302  *Caveat.* Much as SVMs provide a fairly robust method for classification across a surprisingly wide
303  range of data types and problem domains (though they are rarely truly optimal due to NFL), there is
304  some hope that such "pretty good, most of the time" dMVPA architectures might exist as well.
305  Again, the field is young, but during development of the DeLINEATE toolbox, we have often found
306  that relatively simple dMVPA models, consisting of just 1–2 convolutional layers and 1–2 dense

---

[22] This term is less commonly used in the MVPA literature than the ANN literature, but it refers essentially to a parameter of the algorithm set by the user before running the analysis (for example, the amount of regularization), to distinguish those values from plain (non-hyper) parameters, which are the values estimated by the statistical process or model-fitting algorithm.

[23] A full rundown of layer types is beyond this article's scope and better-suited to a general introduction to deep learning, but common types include perceptron-style "dense" layers, "convolutional" layers, "recurrent" layers, and supporting utility layers that calculate simpler mathematical functions; discussed in more detail below.

[24] Complexity could be defined many ways, but for now, we will use it mainly to refer to how many parameters (not hyperparameters) need to be estimated for a given model.

307  layers[25], perform comparably to (or better than) the industry workhorse of SVMs. A bit of
308  customization is often required to fit the size and shape of the input dataset, and it can be useful to
309  test out different variations of dMVPA architecture on one portion of the dataset before applying the
310  best-performing architecture to the remaining held-out data, but a satisfactory architecture is typically
311  not too difficult to find without excessive trial-and-error. We have found that after some experience
312  using dMVPA, one begins to develop fairly good intuitions about what kinds of architecture might be
313  best suited to a specific problem, but it is still far from an exact science.

314  As the field progresses, we hope that it will converge on more heuristics for designing dMVPA
315  architectures that perform as robustly as SVMs across datasets, while still retaining the flexibility and
316  other advantages of dMVPA. Still, for many practical applications, it is less important to identify an
317  optimal model than it is to determine if the data can be reliably classified above chance (Hebart &
318  Baker, 2018). With properly implemented cross-validation, this can often be achieved by a wide
319  variety of architectures (assuming the data do contain enough meaningful signal for reliable
320  decoding), with the accuracy difference between sets of hyperparameters being only a few percentage
321  points. Conversely, if the input data contain only noise with respect to the classification problem, any
322  sane architecture should perform at chance on the test set. Thus, while some trial and error may be
323  necessary before deciding that data cannot be classified, exhaustive model search is seldom required.
324  When possible, it is often helpful to conduct a traditional MVPA to get a ballpark estimate of how a
325  reasonably well-configured dMVPA should be expected to perform.

326  *Pro: Intrinsically multiclass classification.* One advantage of dMVPA whose value is likely
327  underestimated is that it is straightforward to design a "true" multiclass classifier, whereas most
328  traditional MVPA methods are intrinsically binary. Thus, in traditional MVPA, multiclass decisions
329  must generally be built from a combination of binary classifiers[26]. While there is nothing
330  methodologically wrong *per se* with building multiclass decisions from binary ones, the implications
331  are slightly different than those of a true multi-way decision, which should be taken into account
332  when interpreting results. Furthermore, in some commonly-used MVPA tools (e.g., PyMVPA), the
333  multiclass decision procedure is not always transparent to the end user, which can be a point of
334  confusion. Conversely, dMVPA classifiers are able to consider all classification options
335  simultaneously; as a consequence, it is also trivially easy to obtain meaningful prediction scores
336  across all classes for each example in the testing set, which can then be used in analyses that go
337  beyond simple winner-take-all accuracy measures.

338  *Pro/Con: Performance.* Performance, in the sense of speed, can be either an advantage or a
339  disadvantage of dMVPA. Although dMVPA network architectures can vary so widely that it is
340  difficult to generalize, *prima facie* dMVPA should typically run slower than traditional MVPA,
341  because the calculations involved in training a dMVPA network are more complex. However, for
342  larger datasets (in terms of numbers of features and/or examples), the performance of traditional
343  MVPA techniques may scale more poorly than dMVPA. (See "Benchmarks" below and Table 1 for
344  details.) Thus, beyond a certain dataset size, dMVPA may be the only feasible choice. Also, because

---

[25] Technically, these "deep" MVPA networks would not be very deep in terms of how many layers they contain. Still, a fair portion of "deep" learning these days does not use particularly complex network structures; the term now seems to refer more to the contemporary era of ANN-based data analysis than any particular network structure.

[26] Typically, if we have classes ABC, the multiclass decision would be made either by training up classifiers "A vs not-A," "B vs not-B," and "C vs not-C," or by training up classifiers "A vs B," "A vs C," and "B vs C," and then summing up the scores in favor of each category across classifiers in order to obtain an overall score for that category.

345    the network architecture of dMVPA can be adjusted, researchers have more options; e.g., whether to
346    employ a simpler network that may not achieve maximum accuracy but runs quickly, versus a more
347    complex network that runs slower.

348    *Caveat.* As alluded earlier, dMVPA's computational costs can be somewhat offset by parallelization,
349    which is better supported by deep-learning software tools than most traditional MVPA tools. This is
350    true even if parallelizing across CPUs/cores, but especially true if using the computer's GPU. Results
351    vary widely depending on dataset size, network architecture, and the specific hardware involved, but
352    users might roughly anticipate anywhere from a 5x–100x or more speedup for running dMVPA on a
353    GPU versus a CPU. On one hand, these benefits make dMVPA a more competitive option, speed-
354    wise. On the other hand, GPU-accelerated dMVPA does require more specialized hardware and more
355    human effort setting up the relevant drivers and software packages. While we have striven in our
356    toolbox and documentation to keep this process as painless as possible, it is still more effort than is
357    required to run non-GPU-accelerated analyses; whether that effort is well-spent will heavily depend
358    on individual users and what tasks they are trying to accomplish.

359    *Pro: Flexibility of applications.* Although our focus has been on dMVPA, we should note that
360    modern neural networks have an ever-increasing number of uses beyond simple classification. For
361    example, one currently popular strategy is to train a model for categorization within some domain
362    (e.g., the contents of a photograph) and then interrogate the model's intermediate layers, in an
363    attempt to understand what strategy the model is using (Zeiler & Fergus, 2014). Autoencoder-style
364    architectures allow for, e.g., unsupervised learning of feature structure (Xie et al., 2016), feature-
365    sharpening for degraded inputs (Lore et al., 2017), and principled fusion of multimodal data (Ngiam
366    et al., 2011). Deep networks can also be used to implement classification techniques that are not well-
367    suited to traditional MVPA – for example, "transfer learning," in which a network is initially trained
368    on one dataset, and then refined by training it further on a different dataset. As another example, we
369    have recently explored using deep networks to create "smarter" similarity/distance metrics tailored to
370    particular datasets/applications, unlike traditional formula-based metrics (e.g., Pearson correlation,
371    Euclidean distance), which do not afford such flexibility (Williams et al., 2020). The DeLINEATE
372    toolbox can, with varying degrees of effort, support many of these advanced applications.

373    *Con: Field and dependencies are in active development.* While the software tools for traditional
374    MVPA will presumably keep receiving periodic updates, the field overall is fairly mature and not
375    changing particularly rapidly. However, deep learning and dMVPA are newer; as such, the
376    techniques and their underlying software tools are continually being updated. This means that
377    documentation can rapidly go out of date, and incompatibilities can arise easily if developers are not
378    careful. We have aspired to make our own toolbox as robust as possible to the changing software
379    landscape, but it is still worth being aware of. Of course, there are mitigating strategies: Users can
380    find one version that works and refuse to update anything, but this deprives them of future
381    enhancements. Alternately, they can continually update, but this makes it harder to exactly replicate
382    earlier work run with previous software versions. If only Python toolboxes (our DeLINEATE
383    toolbox, and the Keras/PyMVPA backends it relies on) are updated, Python's "virtual environment"
384    feature can be helpful for maintaining different software setups, each in their own containers. But, if
385    later updates require newer hardware drivers, *and* users wish to maintain backward compatibility
386    with their earlier work, they may wish to do what our lab has done: Purchase several small hard
387    drives for each machine, set up a fresh operating system for each new major driver version, and
388    simply reboot from a different boot drive when one wishes to work with current vs. legacy versions
389    of the software.

## 2.5 A brief introduction to network architecture

390

391 In an abstract sense, all feedforward[27] neural networks may be viewed as a collection of
392 mathematical operations to be applied in sequence to an input of some fixed size, along with rules for
393 updating the parameters of those operations during training. In a classic perceptron, the core
394 operations are multiplication (input data times weight values), summation, and then activation (a
395 thresholding operation, traditionally). In a *multi-layer* perceptron network (Figure 2A), this complete
396 multiplication-summation-activation sequence is repeated, with each layer's outputs becoming the
397 next layer's inputs. A typical, slightly simplified mental model for such networks treats those
398 multiplication-summation-activation operations as all occurring within a self-contained unit or node,
399 like in a biological neuron; a number of such units in parallel constitutes a layer of the network, and
400 the main free parameter chosen by the designer of the network is the number of units in each layer.
401 However, unlike a biological neuron, in an ANN this set of operations is not immutable – one might
402 opt to omit activation, invert values after every step, or do any other sort of mathematical
403 transformation, at any step of the sequence. One could also adopt a different mental framework in
404 which every individual operation is a layer of the network, such that each layer of a perceptron
405 network expands into three sequential computational layers: a multiplication layer, a summation
406 layer, and an activation layer. In Keras, the Python framework upon which the DeLINEATE
407 toolbox's deep-learning functionality rests, it is possible to work with either of these
408 conceptualizations – e.g., there are individual layer types that can perform thresholding/activation,
409 but the activation operation can also be specified as an argument of other layer types, with the
410 understanding that activation is applied last, after that layer's primary operation.

411 In lay terms, when sufficiently tortured and beaten into submission, contemporary deep learning
412 frameworks can be mangled into performing virtually any kind of mathematical operation or
413 transformation on the input data. A full discussion of all the possibilities could fill several books, and
414 is thus beyond the introductory scope of this paper. However, there are a few broadly useful kinds of
415 operation/layer that are particularly worth understanding; novices to deep learning should focus on
416 understanding the basic gist of these fundamental tropes before getting lost in the details. Here, they
417 are described briefly in broad categories; Keras has several subtypes of each depending on details of
418 the desired implementation.

### 2.5.1 Classic

419

420 Called "Dense" layers in Keras, these are layers made of perceptrons (Figure 2A). They compute
421 weighted sums and apply an activation function. Varying the number of computational units in such a
422 layer allows one to increase (e.g., consider more potential weightings) or decrease (e.g., prune less
423 informative features) the dimensionality of the data as it passes through the layer. By default, these
424 layers are fully-connected, meaning that all outputs from one layer are used as inputs for each

13

[27] "Feedforward" meaning that all outputs from earlier (closer to the input) layers are fed "forward" into later (closer to
the output) layers; outputs are never fed back into earlier layers. Feedforward networks are generally easier to work with
and design. Our toolbox currently supports only networks with a broadly feedforward design (implemented via the
"Sequential" model class in Keras) when using the graphical interface or text-based job files; however, when using it as a
collection of Python functions, other network types are possible. One exception is recurrent layers, which feed their
output back into themselves; thus networks containing recurrent layers are not strictly feedforward. However, as
implemented in our toolbox and the Keras backend we rely on, the recurrency can be viewed as something that recurrent
layers handle within themselves; the user does not have to think about this recurrency in terms of their network
architecture. From the user's point of view, the *layers* of the network still follow a feedforward/sequential structure, even
if the individual *units* within some layers have recurrency built-in.

425    computational unit in the next layer of the network. As noted earlier, a neural network made entirely
426    of dense layers is sometimes called a "multi-layer perceptron" network architecture.
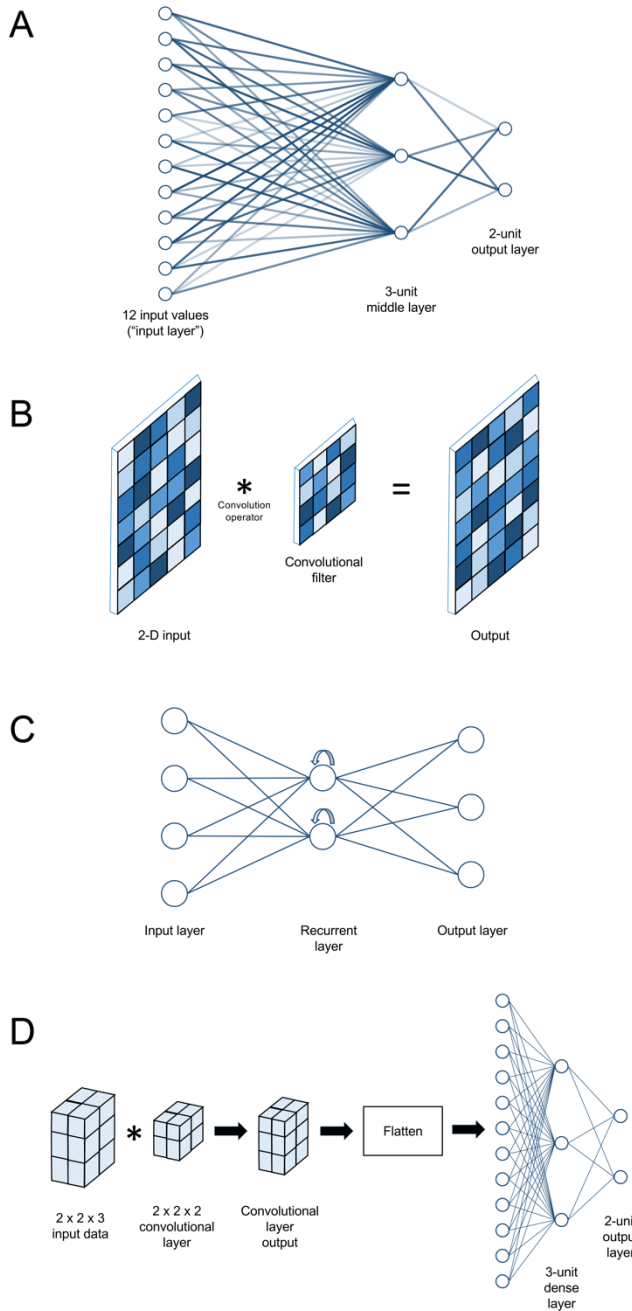


**Figure 2. Examples of artificial neural network architectures.**

(A) A simple fully-connected multi-layer perceptron model with 12 input values, a middle layer comprising three perceptrons, and an output layer with two perceptrons. Line lightness is used to represent synaptic weight strength. (B) An example of a convolutional neural network layer that might be used to analyze 2-D input. Here, the layer looks less like a set of artificial neurons and more like a digital filter used in image processing. Two-dimensional input is convolved with a 2-D filter to yield a 2-D output, sized similarly to the input. During training, it is the values in the convolutional filter that get adjusted. Square lightness represents the numeric values in the cells of each 2-D matrix. (C) An example of a simple recurrent neural network. There are many types of recurrent neural network structures with varying degrees of complexity, but all share the property that recurrent units' output gets passed back into them (represented here by curved arrows), giving them some form of "memory" for previous input values. (D) An example of a complete neural network architecture that might be used to analyze 3-D input such as MRI data for a two-class classification problem. In this simple example, 12 input values in a $2{\times}2{\times}3$ array are first passed through a $2{\times}2{\times}2$ convolutional filter, yielding another $2{\times}2{\times}3$ array as output. This is then passed through a "flattening" layer to convert it to a $12{\times}1$ vector, which then passes through a 3-unit dense layer to a 2-unit output layer (as shown in panel A).

### 2.5.2 Convolutional

428    Convolutional layers (Figure 2B) may be conceptualized as collections of filters that are swept across
429    (in mathematical terms, convolved with) their input. When used to process 2-D photographic data,
430    their function is often likened to visual neurons, which take input from a spatially restricted receptive
431    field, extract some feature if present, and pass along the result to the next layer of the visual
432    processing hierarchy. For readers familiar with digital image processing, they are essentially like
433    other kinds of digital filters (e.g., a blur filter, an edge detector), except that convolutional layers can
434    work with any dimensionality of data (not just 2-D images) and their parameters are learned over the

435 course of training, rather than being pre-defined. The combination of filter shape and input data
436 structure will determine what kinds of feature may be selected for and passed along as output. For
437 example, if each example of input data is a $32 \times 1000$ array of EEG voltages (e.g., 2 seconds of 32-
438 channel data sampled at 500 Hz), a set of $1 \times 10$ filters would be capable of detecting high-frequency
439 patterns within individual channels (in this example, patterns that fit inside a 20 ms time window),
440 but insensitive to lower-frequency or purely spatial patterns. Conversely, a set of $10 \times 1$ filters could
441 detect patterns distributed across multiple channels, but only those that occur instantaneously.
442 However, one could instead employ, for example, a set of $8 \times 20$ filters, which would be capable of
443 detecting patterns spread across up to eight adjacent channels over a 40 ms time window. Choices
444 about data structure are consequently more important for this class of layers than for a multi-layer
445 perceptron; the input examples would contain identical information if flattened from 32 channels $\times$
446 1000 timepoints to a single $1 \times 32,000$ vector, but the meaning of a $1 \times 10$ filter bank's outputs would
447 be very different.

### 2.5.3 Recurrent

449 Recurrent layers (Figure 2C) are named for their property of having their outputs fed back into
450 themselves as inputs. By maintaining an internal state determined by previous inputs, recurrent units
451 develop a form of memory for sequential data. For example, a $1 \times 10$ vector input to a classic dense
452 unit would be combined to a single value in only two steps – multiplying each element of the vector
453 by its weight and then summing the results. If the same vector were fed into a recurrent unit
454 (typically called a cell), the first element would be handled in isolation, but evaluation of the second
455 element would include the output of the cell's operation on the first element. The result of this would,
456 in turn, update the unit's state to influence its response to the third element, and so on until each
457 element of the input is consumed. Recurrent networks are frequently used to process natural language
458 data (both audio and text) and in general are considered good choices for timeseries data. In our own
459 work, we have not observed any significant benefit over convolutional layers when working with
460 human neuroscience data, and have found recurrent-based networks to take longer to train than
461 convolutional-based networks; however, these findings are likely highly dependent on details of the
462 dataset and research question. As alluded earlier, for common types of recurrent cells, the recurrency
463 is handled within the cell as a form of internal "memory" that is not visible to the rest of the network,
464 so network architectures using recurrent layers can still be considered broadly "sequential" or
465 feedforward, and are thus supported by our toolbox.

### 2.5.4 Supporting

467 This is a broad category of operations that, for various reasons, are generally thought of as secondary
468 or historically baked-in to more interesting operations. In Keras, this includes activation layers,
469 various purely utilitarian data-reshaping or simple mathematical operations, dropout (an operation in
470 which some percentage of a layer's units are ignored; thought to mitigate overfitting), etc. Some of
471 these operations (e.g., activation functions) can be specified either as distinct layers or as parameters
472 to a primary layer, whereas others (e.g., a layer that downsamples the output of the previous layer via
473 averaging) can only be specified as distinct layers.

### 2.5.5 Practical advice

475 The following is a combination of our experience and advice we have received from other
476 colleagues. We hope it is helpful as a starting point, but readers should not feel overly constrained by
477 it. While the modern leaders in image recognition involve dozens of layers (Szegedy et al., 2016), in
478 our experience the aim of dMVPA can typically be accomplished with much smaller networks. When
479 working with minimally-processed fMRI/EEG/eye-tracking data, we have found that a good starting

480    point often consists of 1–2 convolutional layers followed by 2–3 dense layers; based on preliminary
481    results from that architecture, one could add or remove layers, adjust the layers' sizes, or tweak other
482    hyperparameters. See Figure 2D for an example. For maximal effectiveness and interpretability,
483    consideration should be given to the match between the shape of the per-example input data and
484    shape of convolutional filters (e.g., should the filters look across EEG channels, or only within? If
485    across, are channels arranged to be spatially adjacent in the data?). *Leaky ReLU* is usually our
486    preferred activation function, and we have often found dropout values of ~0.3 in dense layers to be
487    beneficial. We have found *Stochastic Gradient Descent (SGD)* with the momentum parameter
488    (classical or Nesterov) set to something on the order of 0.9 to be a generally successful optimizer,
489    although the *Adam* optimizer (Kingma & Ba, 2014) also performs well in some situations[28]. New
490    users are encouraged to experiment with everything and keep track of the results; soon, you will
491    likely develop your own favorite architectures and hyperparameters. Do not be afraid to experiment
492    broadly; dMVPA has some powerful advantages, but we are also in a more exploratory phase for this
493    kind of research, and designing a sufficiently performant dMVPA architecture can take significant
494    trial-and-error. Of course, the extent to which that exploration might constitute *p*-hacking depends on
495    your research aims; if that is a potential concern, you may want to design your analysis based on an
496    independent dataset (e.g., one of the sample datasets included in our toolbox), or consider a split-half
497    design in which one half of your data is used to explore analysis architectures and the other half is
498    used for confirmatory purposes.

499    ## 3    dMVPA: A toolbox

500    ### 3.1    The DeLINEATE Toolbox

501    One major purpose of the DeLINEATE toolbox is to enable rapid exploration of model
502    architectures/hyperparameters while maintaining an accurate record of what was done and how it
503    turned out. These are conflicting goals in common practice – a researcher attempting to iterate on an
504    analysis is often tweaking a script or working directly with a command-line interpreter, perhaps in a
505    Notebook type environment (Grus, 2018), and discarding fruitless branches of exploration along the
506    way. Maintaining an accurate record of each tweak and its results during such rapid prototyping is
507    not easy, and can take more time and coding discipline than many of us have.

508    Our solution to this problem was a processing pipeline in which a single JSON (JavaScript Object
509    Notation) format[29] job configuration file fully specifies an analysis: the input data, how it will be
510    divided for cross-validation and rescaled, the model architecture to be trained and evaluated, and the
511    outputs to be saved (Figure 3A). The toolbox translates this JSON file into Python code to execute
512    the specified analysis (or analyses), and saves all desired outputs into .tsv (tab-separated values) files
513    with names that include a user-defined prefix linking them to the original JSON file. A copy of that
514    original JSON file can also be saved alongside the other output, so that even if the original is

16

---

[28] We realize that all this terminology can be overwhelming at first, but readers unfamiliar with deep learning should try not to feel discouraged by the sheer number of architecture/hyperparameter choices available. Rest assured that it does become more familiar and accessible after some hands-on experience.

[29] JSON is a format that allows data structures to be written to plain text files with human-readable syntax. Although not as intuitive as a graphical interface, editing JSON-formatted job files is certainly easier for beginners than writing their own Python code. There are also JSON modules available for many popular text editors and a handful of standalone JSON editing programs to make the task even easier.

515    subsequently overwritten during the exploration process, the "output" copy remains a pristine record
516    of what was run to create a particular set of results.
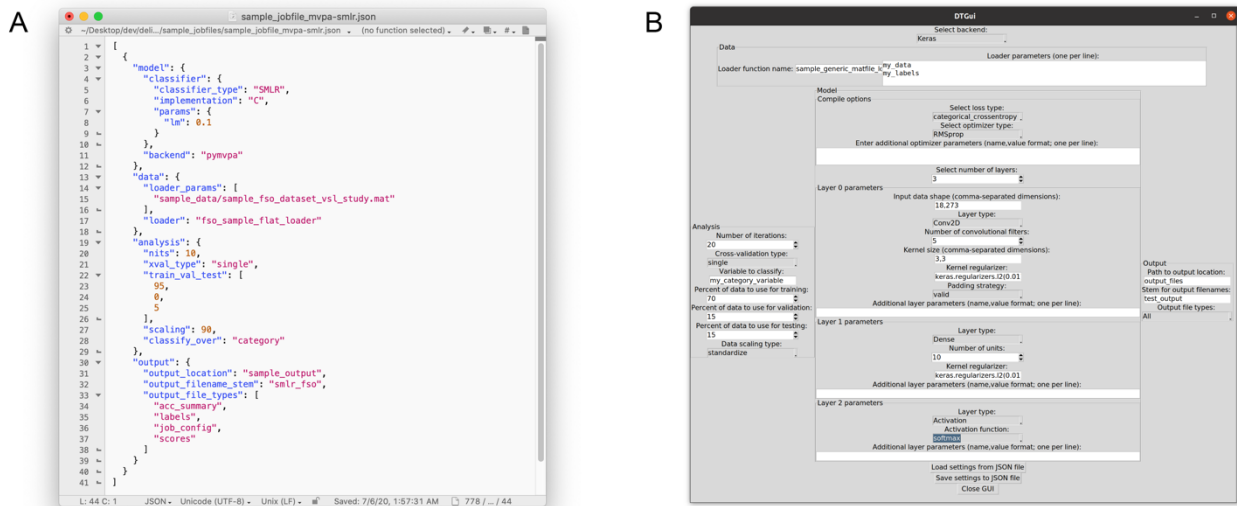


**Figure 3. Ways that users can configure an analysis in the DeLINEATE toolbox.**

(A) Most users will likely configure analyses using a text-based JSON (JavaScript Object Notation) format job file. In this example, the file is open in a generic text-editor program, but JSON-format-specific editing software also exists. Each job has four main sections: "model," "data," "analysis," and "output," corresponding to the major object types in the toolbox. The file shown is configured to run 10 iterations of a PyMVPA-based SMLR analysis using a sample face-scene-object-viewing EEG dataset, using a randomly selected 95% of trials as training data and 5% as test data on each iteration. (B) A basic graphical user interface (GUI) that allows users to configure a job file without having to edit the text directly. The most frequently used options for several common analysis types are available (although editing the text file directly will always allow more flexibility than is possible to express in a GUI). The GUI also contains sections for data, analysis, model, and output, as well as buttons for loading in an existing job file and saving the settings configured in the dialog box to a new JSON file. The settings shown are configured to run 20 iterations of a Keras-based deep learning analysis, using 70% of trials as training data, 15% as validation data, and 15% as test data on each iteration.

517    A secondary goal was to facilitate comparison of dMVPA approaches to traditional MVPA while, as
518    much as possible, maintaining parity in data handling. To this end, classic MVPA is also supported
519    alongside the dMVPAs that are our primary focus. This is currently implemented with a PyMVPA
520    backend. Traditional MVPA uses the same JSON job file format as dMVPA, as well as similar
521    output file formats, cross-validation/rescaling options, etc., making it a simple task to conduct
522    parallel MVPA and dMVPA on the same data. Currently we support SVM (Support Vector Machine)
523    and SMLR (Sparse Multinomial Logistic Regression) classifiers for traditional MVPA, although our
524    framework is readily extensible to most other classifiers in the PyMVPA toolbox.

525    For a typical user, the primary entry point to the toolbox is *delineate.py*, a simple script that accepts
526    one or more JSON-format configuration files as arguments, validates their contents, and uses them to
527    create and run one or more analysis job(s). This allows users to run analyses without requiring them
528    to write any code of their own. To further increase accessibility, we have recently developed a simple
529    graphical user interface (GUI) that some find more approachable than a text editor (Figure 3B). GUI
530    users can click on a collection of interactive menus to create properly-formatted job configuration
531    files, which can then be used as input to the main *delineate.py* script. The GUI can also auto-populate

532  selections based on an existing job configuration file for users who have a starting point (such as one
533  of the included sample job files) they wish to modify for future analyses.

534  For Python-proficient users who want more complex or flexible analysis options, the toolbox can
535  also be used as a Python programming library, and users can write their own code instead of creating
536  JSON files. JSON functionality and code-library functionality can also be mixed-and-matched (e.g.,
537  JSON files can be used to create a template analysis, which can then be tweaked and iterated upon
538  with custom code). For users who wish to write their own Python code as well as JSON users who
539  simply want some familiarity with the toolbox's underlying functionality, we next present a brief
540  overview of the code structure; more detail is available in the toolbox documentation.

541  ### 3.2  DeLINEATE Toolbox structure

542  The DeLINEATE Toolbox is an object-oriented collection of Python modules, each responsible for a
543  different aspect of the (d)MVPA process. It comprises five main object classes and a small number of
544  supporting files that contain utility functions or facilitate batch analysis. Each main class is housed in
545  a .py file named for that class. In typical usage, the toolbox follows a minimum-import philosophy;
546  to use it as a code library, one simply needs to navigate to its main directory and directly import the
547  desired class file(s). The primary classes are:

548  1.  *DTJob*, responsible for parsing JSON files that define DeLINEATE jobs and passing the
549      appropriate information to constructors for the other object types. In typical usage, a DTJob is
550      responsible for creating one of each other object type and then triggering the DTAnalysis
551      object to actually run the analysis. However, users can also eschew DTJob entirely if they
552      prefer to instantiate the other objects manually in their own Python code.

553  2.  *DTAnalysis*, a parent class that contains one instance each of DTModel, DTData, and
554      DTOutput; it is responsible for coordinating the operations of those other objects. This
555      includes dividing data into training/validation/testing sets, iterating through portions of the
556      data when desired (e.g., to loop through individual subjects), and initiating the model
557      training/testing procedures.

558  3.  *DTModel*, responsible for constructing the model in the appropriate machine learning
559      backend (currently, either Keras or PyMVPA). The "model" in this sense refers either to the
560      artificial neural network (Keras) or an object representing a simpler classifier, e.g., a support
561      vector machine with a linear kernel and parameter C=1 (PyMVPA).

562  4.  *DTData*, responsible for loading the dataset from a data file, storing it, and performing certain
563      operations on it (such as scaling/normalization or slicing it up into smaller training,
564      validation, and/or test subsets).

565  5.  *DTOutput*, responsible for writing analysis results to output files.

566  The four main sections of a JSON-format job file are the *analysis*, *model*, *data*, and *output* sections,
567  which map directly onto the corresponding Python classes; each section contains the parameters

568  necessary to instantiate an object of the appropriate class[30]. Another (purely optional) class, *DTGui*,
569  implements the aforementioned GUI.

### 3.3  Current functionality

### 3.3.1 Model types and backends

572  At present, the DeLINEATE toolbox has been used in-house for approximately two years to conduct
573  analyses across a number of studies. It is a high-level toolbox with a flexible, extensible architecture
574  that potentially allows it to sit atop multiple underlying machine-learning libraries. Currently, we
575  support a subset of functionality for two backends: Keras (Chollet et al., 2015) for dMVPA and
576  PyMVPA (Hanke et al., 2009) for traditional MVPA. With our heavy focus on providing a flexible
577  architecture, it is relatively easy to add support for additional backends in the future, as well as
578  enhancing the breadth of support for features of Keras and PyMVPA, enabling new data types to be
579  imported, etc. The relative prioritization of such extensions will be guided by user demand.

### 3.3.2 Cross-validation

581  We currently support two approaches to cross-validation. The first is a "universal" approach
582  (specified in configuration files with the name "single") in which all data are treated as belonging to a
583  single pool, which is randomly divided into training/validation/test sets according to percentages
584  specified in the configuration file. The second divides the data according to some attribute of the
585  samples[31] and iterates through each value of this property, dividing the data within each iteration into
586  training/validation/test sets (specified in configuration files as "loop_over_sa"). Regardless of which
587  scheme is used, because classification performance can be influenced by a model's initial
588  conditions[32], it is common practice to run multiple complete cross-validation iterations in order to
589  ensure a stable estimate of the architecture's performance. With properly configured input data (see
590  below), these two cross-validation schemes can cover most common MVPA use cases; however,
591  additional schemes can be added in the future according to demand.

### 3.3.3 Rescaling

593  Although some MVPA methods are invariant to the scaling of the input data, others, such as many
594  dMVPA applications, require data to be on a certain scale for good classification. The issue is
595  slightly complicated by the need to prevent features of the test data from influencing the training
596  data. We support several methods for rescaling data that avoid this issue by calculating necessary
597  parameters solely on the training data, and using those parameters to adjust validation/test data as

19

[30] Although a non-Python-savvy user does not need to know these implementation details, the parity between job file sections and Python classes makes it easy for more experienced coders to switch back and forth between job files and their own Python scripts. As noted above, it is also possible to mix-and-match the two approaches.

[31] A "sample attribute," if you will, which is the terminology used by other MVPA toolboxes for a tag or property associated with each data sample/example. For instance, a subject ID or session ID.

[32] Especially for dMVPA; for a given architecture, a classification might sometimes perform well and sometimes at chance depending on the random values assigned to weights at the beginning of training, which is generally a sign that the architecture needs adjusting. Other classification techniques, such as SVMs, are deterministic; as such, they may or may not benefit from multiple cross-validation iterations, depending on dataset and cross-validation scheme.

598    well. Again, these methods are readily extensible with additional options, or users can always pre-
599    scale their own data however they like. Currently supported methods are:

600      1. "percentile", which identifies the value at a specified percentile of the data and divides all
601        data by that value,

602      2. "standardize", which mean-centers and divides all values by the standard deviation of the
603        data,

604      3. "mean_center", which subtracts the mean of the data from all values,

605      4. "map_range", which translates values into the range between a user-specified minimum and
606        maximum (0 and 1, by default).

### 3.3.4 Input data and loaders

608    By far, the most common question we have received from potential users concerns the necessary
609    format for input data. The toolbox operates on, at minimum, one NumPy array and one Python
610    dictionary. The former contains the actual data to be analyzed in a two-or-more-dimensional array,
611    where one dimension represents examples (e.g., trials) and the other dimension(s) are feature
612    dimensions. For instance, an fMRI dataset might be shaped as (examples × voxels), whereas an EEG
613    dataset might be (examples × electrode × timepoint). Higher-dimensional structure is ignored in
614    traditional MVPA and simply collapsed into a 2-D (examples × features) array, as those simple
615    classifiers can only operate on vectors of data. However, dMVPA, when run with an appropriate
616    network architecture, can operate on any dimensionality of data and can potentially take that
617    information into account for classification. If the spatiotemporal structure of the data is meaningful,
618    this may produce superior performance. The Python dictionary contains the metadata needed to
619    interpret the data array, in the form of one or more "sample attributes" (defined earlier; e.g.,
620    experimental condition, participant identity) for each sample. These sample attributes may be used as
621    targets for classification (i.e., the class labels to be predicted) or as grouping variables in cross-
622    validation (e.g., for leave-one-subject-out cross-validation).

623    Data are read into the toolbox by a "loader" Python function specified in the job configuration file.
624    Loaders can reside in a specific subdirectory of the toolbox or in an arbitrary user-specified location.
625    We include several example datasets and corresponding loader functions that should be easily
626    modifiable by researchers to fit their own needs. This is the main place where a typical user *might*
627    need to write their own Python code; because of the many idiosyncratic formats used to store
628    experimental data, some users may need to write a short function to read their files in and reshape
629    them into the expected format. However, if the format is well-supported by NumPy or other Python
630    libraries, these functions can typically be quite short (on the order of 10 lines of code). We also
631    provide generic functions included for data in the NumPy and MATLAB native file formats, which
632    will accept any .mat or .npy file containing one array variable of data examples and at least one
633    variable of sample attributes. Thus, if users are able to save their data in one of those formats
634    beforehand, there may be no need for a custom loader function.

635    Because neuroscience data vary widely in format, we recognize that a need for additional loader
636    options could still present a barrier to some researchers. We encourage such individuals to reach out
637    to us so that we can offer assistance and expand the range of formats we are able to support natively.
638    On the other hand, the overall flexibility in format means that with just a few lines of code, any
639    dataset that can be represented as a multi-dimensional array is a candidate for analysis with our

640  toolbox, not limited to neuroscience data; for instance, we have used the toolbox to analyze eye-
641  tracking data (Cole et al., under review), photographic images, and more.

### 3.3.5 Graphical User Interface

643  As described earlier, the GUI currently allows users to generate a job configuration structure via
644  menu selections and free-entry fields (Figure 3B) that can be auto-populated by loading an existing
645  job file. For frequently used Keras layer types, some reasonable default hyperparameters are
646  provided; however, there are minimal defaults available for less common layer types, and in general
647  it is still recommended for users to have some baseline knowledge of Keras's workings and
648  hyperparameter options, even when using the GUI. As the number of potential analysis
649  configurations is effectively limitless and this module is a relatively recent addition, error checking is
650  currently somewhat limited. Still, we recognize that a usable GUI is a critical feature for some users,
651  and we expect this to be a primary target for expansion and refinement in upcoming releases.

### 3.4   Availability

653  All toolbox code is currently hosted at https://bitbucket.org/delineate/delineate and is freely
654  accessible and open-source under the MIT License. There is also a project website at
655  http://www.delineate.it/ that hosts older releases, documentation, links to video tutorials, and more.

### 3.5   Hardware/software requirements

657  The DeLINEATE toolbox has few software dependencies of its own. However, as noted earlier, it
658  requires either a Keras or PyMVPA backend to perform dMVPA or traditional MVPA, respectively,
659  and those packages have their own corresponding dependencies. Fortunately, both Keras and
660  PyMVPA are well-documented and readily available; we also provide start-to-finish setup guides on
661  the toolbox website. In brief, DeLINEATE is compatible with any recent version of either backend,
662  and in principle can be run on any Python version from 2.7 onward, including all versions of Python
663  3; however, specific Python version compatibility may depend on which version of Keras/PyMVPA
664  the user is running, and which Python versions *those* libraries are compatible with. The only
665  additional dependency of DeLINEATE is Python support for Tcl/Tk (a graphical interface toolkit) if
666  one wishes to use DTGui; most Python installations include Tcl/Tk libraries, but some might require
667  a separate installation. As Python is available on all major operating systems (Windows, macOS, and
668  Linux), DeLINEATE will also run on any of them, although hardware choices may constrain
669  operating system options.

670  In terms of hardware, a bare-bones DeLINEATE installation will run on any computer with enough
671  RAM to hold the user's dataset in memory, as long as the user only wishes to run analyses on the
672  CPU. Traditional MVPA via PyMVPA does not presently employ GPU acceleration, but most
673  dMVPA users will want to enable GPU acceleration for a dramatic increase in speed (see
674  "Benchmarks" below). As Keras relies on the TensorFlow library for its own backend (or the older
675  Theano library; now deprecated in recent Keras versions but still supported by DeLINEATE), which
676  in turn relies on the CUDA (Compute Unified Device Architecture) and cuDNN (CUDA deep neural
677  network) libraries from NVIDIA, effectively this means that an NVIDIA-compatible GPU is required
678  for accelerated dMVPA. Different GPUs will have different compatibility with various versions of
679  CUDA, cuDNN, TensorFlow/Theano, and Keras; however, as long as compatible versions of those
680  tools are installed, DeLINEATE should work with any of them. At the time of writing, we
681  recommend midrange to high-end GPUs from the GeForce 10 series or higher; our lab's workstations
682  mostly use GeForce GTX 1070 through GeForce GTX 1080 Ti cards, but other users may have

683    higher or lower requirements. Currently, a reasonably powerful workstation for many dMVPA
684    applications could be built from parts for $1500–2000 US[33], although prices can vary widely
685    depending on users' specific requirements and budgets. Since no current Apple computers support
686    compatible NVIDIA GPUs, GPU-accelerated dMVPA is currently unavailable on macOS. Generally,
687    for scientific computing, we recommend Linux-based operating systems for their widespread
688    compatibility and open-source nature; however, GPU-accelerated dMVPA will work on Windows as
689    well. In the future, if the macOS/NVIDIA compatibility situation changes, or if DeLINEATE adds
690    support for additional backends, GPU-accelerated dMVPA may become available on macOS.

691    It has historically been difficult to implement large neural networks without setting up dedicated
692    hardware, largely because the virtualization approaches favored for cloud-based computing do not
693    provide sufficient access to GPUs. However, we have recently seen the emergence of an option that
694    may be useful to those who lack either the budget or the technical confidence to set up their own deep
695    learning environments. Google Colab (https://colab.research.google.com) is a browser-based Python
696    environment akin to Jupyter Notebooks with some access to GPUs. Because the provided
697    environment includes Keras/TensorFlow and allows interaction with files stored on Google Drive, it
698    is relatively straightforward to execute DeLINEATE-based analyses by importing some of the classes
699    and manually calling the method that begins an analysis. An example IPython notebook is provided
700    in the Colab subfolder of the DeLINEATE repository. This approach requires some proficiency in
701    Python and is subject to fluctuating resource limitations, so no promises can be made about speed or
702    stability; however, it may be a good jumping-off point for beginning users wishing to explore the
703    toolbox before investing in their own equipment.

### 3.6    Benchmarks

704

705    For both traditional MVPA and dMVPA, performance (both accuracy and computation time) will
706    vary drastically across datasets, hardware, and choice of MVPA classifier or neural network
707    architecture. Thus, the generalizability of any benchmarks is limited. However, to give readers a
708    rough sense of the computational advantages of dMVPA and how running times scale for different
709    dataset sizes, we prepared several datasets and analyzed them with both traditional MVPA and
710    dMVPA. These benchmark datasets emulate the format of an fMRI dataset, but are entirely synthetic.
711    The code to generate them is included in the toolbox.

712    We simulated datasets with three conditions (classes). Datasets ranged from 200 features (e.g.,
713    voxels) to 25,600 features in a doubling progression (200, 400, 800, …). The number of examples
714    (trials) per condition ranged from 100 to 10,000 in the progression: $10^2$, $20^2$, $30^2$, …. Full details
715    are given in the code. Briefly, for each condition, a random signal with the appropriate number of
716    features was generated. Then, supposing for this example that we are generating 900 trials/condition,
717    30 variations on the "canonical" signal for that condition would be generated by blending the
718    canonical signal with a certain proportion of random noise. Then, for each of those 30 variations, 30
719    sub-variations were generated by the same process. Although we did not particularly strive for
720    biological verisimilitude, the intent was to somewhat mimic a circumstance where brain patterns had
721    a small number of "true" variations (e.g., if the condition were "faces," subjects might have slightly
722    different voxel response patterns for different genders/races) as well as trial-to-trial variations due to

22

[33] Based on market prices for parts to build a system similar to ours at the time they were built, with an eight-core Intel i7-9700K CPU, GeForce GTX 1070 GPU, 32GB RAM, 1TB SSD primary storage, 4TB HDD secondary storage, and a compatible CPU cooler, motherboard, case, and power supply, for a total of $1750 US. Newer GPUs and other parts have been released since those were built, but pricing for current parts is in a similar range.

723   stimulus exemplar effects and/or measurement noise. To make the classification more challenging,
724   each trial's signal was also blended with a proportion of the signal of a trial from each of the other
725   two conditions.

726   The datasets were analyzed with three classifier models: a simple CNN, SMLR, and SVM. The CNN
727   used GPU acceleration (NVIDIA GeForce GTX 1080 Ti), whereas the other models used only the
728   CPU (Intel Xeon X5650 @ 2.67GHz). Each analysis was typically run for 10 iterations (cycles of
729   training/test with different randomly-selected training/test sets) except when running times became
730   prohibitive, in which case the analysis was terminated after as few as five iterations.

731   Mean running times (Table 1) ranged drastically, from less than one second to several days. As
732   expected, running times for all model types generally increased with greater numbers of features and
733   trials. SVMs had both the shortest and longest running times. Compared to SVMs, SMLR had both a
734   longer shortest running time and a shorter longest running time (i.e., the range was compressed on
735   both ends), and CNNs continued this trend with an even longer shortest running time and a still
736   shorter longest running time (i.e., the range was even more compressed). Notably, the CNN never
737   took less than 10 seconds (largely due to a relatively fixed start-up time for Keras models) but its
738   longest running times, for the most complex datasets, were still under 15 minutes. By comparison,
739   SMLR's longest running times were over four hours, and SVMs' were multiple days. (And a few
740   SVM models never converged in any reasonable amount of time.) Thus, as expected, deep learning
741   models were less time-efficient than traditional MVPA for simpler datasets but were vastly more
742   scalable for large datasets.

743   Benchmark datasets were intended to be classifiable at moderate accuracies but not particularly
744   designed to be benchmarks *of* accuracy, so we do not report comprehensive accuracy results, which
745   could invite misleading extrapolations to real data. However, generally all methods performed above
746   chance, in a comparable range. Typically, the CNN had the lowest accuracy of all three models on
747   datasets with few trials but usually had the highest accuracy with large trial counts, especially when
748   feature counts were low. Conversely, SVM had the highest accuracy when trial counts were low or
749   with very high feature counts, although in those high-feature-count analyses, the SVM running time
750   was long enough to be unusable in many real-world scenarios. SMLR accuracy almost always fell
751   between CNN and SMLR. Again, we do not expect these accuracies on synthetic data to perfectly
752   reflect performance on real-world data, but they do fit general expectations of how models of varying
753   complexity might be expected to overfit or underfit datasets of varying sizes.

754   ## 4    Discussion

755   ### 4.1    Future development

756   Toolbox development is ongoing and will largely be steered by community feedback. Current goals
757   include adding support for non-sequential Keras models (e.g., those including feedback connections),
758   transfer learning, model introspection, Generative Adversarial Networks (GANs), and additional
759   built-in data loaders and cross-validation schemes. We also plan to make the GUI more informative
760   and intuitive for users who are less familiar with Keras, and to include some tools for visualization
761   and potentially analysis of results (although this remains an unsettled topic; see Hebart & Baker,
762   2018, for relevant discussion). Although we have kept discussion in this paper fairly general,
763   information is still liable to go out-of-date quickly due to the rapid pace of deep learning methods
764   development; users are encouraged to consult our website for the most updated details.

765 **4.2   Summary**

766   Deep learning continues to grow and offer new possibilities for computation in many areas of
767   research and private industry. While it is being increasingly used in neuroimaging and other
768   neuroscience applications, adoption has been hampered by the complexity of the topic and the lack of
769   approachable software tools. We hope that this tutorial review will help researchers new to deep
770   learning address the former, and that the DeLINEATE software toolbox will help address the latter.
771   In years to come, we expect dMVPA to enable a forward leap in neuroscience discoveries
772   comparable to, or exceeding, that of traditional MVPA over older analyses.

773 **5      Tables**

774 **5.1    Table 1. Running times on synthetic benchmark datasets, in minutes**

| CNN | Trials/condition | 100 | 400 | 900 | 1600 | 2500 | 3600 | 4900 | 6400 | 8100 | 10000 |
|-----|------------------|-----|-----|-----|------|------|------|------|------|------|-------|
| Features | | | | | | | | | | | |
| 200 | | .207 | .231 | .244 | .250 | .258 | .295 | .365 | .363 | .438 | .475 |
| 400 | | .216 | .227 | .238 | .240 | .249 | .278 | .297 | .322 | .428 | .520 |
| 800 | | .213 | .231 | .242 | .248 | .262 | .279 | .305 | .341 | .421 | .437 |
| 1600 | | .202 | .258 | .270 | .282 | .298 | .337 | .367 | .414 | .499 | .532 |
| 3200 | | .205 | .322 | .357 | .372 | .376 | .430 | .499 | .564 | .683 | .709 |
| 6400 | | .195 | .326 | .505 | .542 | .657 | .678 | .724 | .900 | 1.18 | 1.15 |
| 12800 | | .292 | .503 | .914 | 1.20 | 1.61 | 1.99 | 2.35 | 2.97 | 2.17 | 2.28 |
| 25600 | | .360 | .920 | 1.32 | 2.49 | 3.71 | 7.41 | 7.73 | 9.43 | 13.5 | 12.1 |

| SMLR | Trials/condition | 100 | 400 | 900 | 1600 | 2500 | 3600 | 4900 | 6400 | 8100 | 10000 |
|------|------------------|-----|-----|-----|------|------|------|------|------|------|-------|
| Features | | | | | | | | | | | |
| 200 | | .024 | .049 | .053 | .086 | .081 | .112 | .127 | .137 | .178 | .219 |
| 400 | | .047 | .326 | .321 | .382 | .292 | .358 | .447 | .484 | .541 | .676 |
| 800 | | .087 | .329 | 1.49 | 1.74 | 1.70 | 1.77 | 1.89 | 1.94 | 2.19 | 2.57 |
| 1600 | | .137 | .563 | 1.70 | 5.40 | 7.78 | 9.09 | 10.1 | 10.3 | 10.8 | 10.8 |
| 3200 | | .194 | 1.37 | 2.53 | 6.06 | 12.6 | 24.4 | 30.0 | 41.3 | 50.8 | 54.2 |
| 6400 | | .258 | 2.59 | 5.59 | 9.57 | 16.6 | 29.4 | 46.9 | 70.7 | 112 | 140 |
| 12800 | | .354 | 4.19 | 13.7 | 20.5 | 28.3 | 42.9 | 64.6 | 91.2 | 128 | 170 |
| 25600 | | .456 | 5.90 | 23.6 | 52.8 | 63.8 | 76.6 | 104 | 144 | 192 | 253 |

| SVM | Trials/condition | 100 | 400 | 900 | 1600 | 2500 | 3600 | 4900 | 6400 | 8100 | 10000 |
|-----|------------------|-----|-----|-----|------|------|------|------|------|------|-------|
| Features | | | | | | | | | | | |
| 200 | | .003 | .103 | .516 | 1.09 | 2.28 | 6.37 | 20.6 | 47.5 | 88.1 | 148 |
| 400 | | .004 | .067 | 1.25 | 3.71 | 5.67 | 16.7 | 65.1 | 155 | 304 | 502 |
| 800 | | .008 | .090 | .710 | 5.32 | 10.4 | 32.0 | 159 | 439 | 920 | 1663 |
| 1600 | | .014 | .216 | .950 | 2.76 | 9.70 | 36.9 | 240 | 855 | 2302 | 4794 |
| 3200 | | .031 | .438 | 2.03 | 5.64 | 12.1 | 22.7 | 109 | 1098 | 3168 | $\infty$ |
| 6400 | | .068 | .892 | 4.32 | 13.2 | 28.3 | 55.2 | 134 | 392 | 2065 | $\infty$ |
| 12800 | | .134 | 1.85 | 9.11 | 28.0 | 65.5 | 133 | 357 | 1241 | $\infty$ | $\infty$ |
| 25600 | | .263 | 3.72 | 18.5 | 57.1 | 138 | 314 | 1048 | 3699 | $\infty$ | $\infty$ |

**Table 1. Running times on synthetic benchmark datasets, in minutes**

We processed a synthetic benchmark dataset with three models: a convolutional neural network (CNN), Sparse Multinomial Logistic Regression (SMLR), and Support Vector Machines (SVM). Average running time is listed in minutes. A few SVM models never converged in any reasonable amount of time and are represented in the table with the infinity symbol $\infty$. See text for further details.

## 6     Nomenclature

ANN: artificial neural network

CNN: convolutional neural network

CUDA®: NVIDIA Compute Unified Device Architecture

cuDNN: NVIDIA CUDA® Deep Neural Network library

DeLINEATE: Deep Learning In Neuroimaging: Exploration, Analysis, Tools, and Education

dMVPA: deep multivariate pattern analysis

DNN: deep neural network

GAN: generative adversarial network

JSON: Javascript object notation

ML: machine learning

MVPA: multivariate pattern analysis

SMLR: sparse multinomial logistic regression

SVM: support vector machine

## 7     Conflict of Interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## 8     Author Contributions

KMK, JMW, PCL, and MRJ worked on toolbox code and co-wrote the manuscript. AS and PKR consulted on the analyses and related projects intertwined with toolbox development, and contributed to the writing of the manuscript.

## 9     Funding

## 10    Acknowledgments

## 11    References

Akama, H., Murphy, B., Na, L., Shimizu, Y., and Poesio, M. (2012). Decoding semantics across fMRI sessions with different stimulus modalities: a practical MVPA study. Frontiers in neuroinformatics, 6:24. doi: https://doi.org/10.3389/fninf.2012.00024

Bentin, S., Allison, T., Puce, A., Perez, E., and McCarthy, G. (1996). Electrophysiological studies of face perception in humans. Journal of cognitive neuroscience, 8:6, 551-565.

Berger, H. (1929). Über das elektroenkephalogramm des menschen. Archiv für psychiatrie und nervenkrankheiten, 87:1, 527-570.

Boser, B. E., Guyon, I. M., and Vapnik, V. N. (1992). "A training algorithm for optimal margin classifiers" in Proceedings of the Fifth Annual Workshop on Computational Learning Theory. (New York, NY, USA: ACM Press), 144–152.

Buzsaki, G., and Mizuseki, K. (2014). The log-dynamic brain: how skewed distributions affect network operations. Nature Reviews Neuroscience, 15:4, 264-278.

Chollet, F. (2015). Keras. https://keras.io [Accessed September 17, 2020]

Cortes, C., and Vapnik, V. (1995). Support-vector networks. Machine Learning, 20, 273–297.

De Martino, F., Valente, G., Staeren, N., Ashburner, J., Goebel, R., and Formisano, E. (2008). Combining multivariate voxel selection and support vector machines for mapping and classification of fMRI spatial patterns. Neuroimage, 43:1, 44-58.

Dosenbach, N. U., Nardos, B., Cohen, A. L., Fair, D. A., Power, J. D., Church, J. A., ... and Barnes, K. A. (2010). Prediction of individual brain maturity using fMRI. Science, 329:5997, 1358-1361.

Grus, J. (2018). I Don't Like Notebooks. Talk given at Jupytercon, New York, NY. Video available at: https://www.youtube.com/watch?v=7jiPeIFXb6U [Accessed November 13, 2019]

Hanke, M., Halchenko, Y. O., Sederberg, P. B., Hanson, S. J., Haxby, J. V., and Pollmann, S. (2009). PyMVPA: a python toolbox for multivariate pattern analysis of fMRI data. Neuroinformatics, 7:1, 37-53.

Haxby, J. V., Connolly, A. C., and Guntupalli, J. S. (2014). Decoding neural representational spaces using multivariate pattern analysis. Annual Review of Neuroscience, 37, 435-436. doi: 10.1146/annurev-neuro-062012-170325

Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P. (2001). Distributed and overlapping representations of faces and objects in ventral temporal cortex. Science, 293:5539, 2425-2430.

838    Hebart, M. N., and Baker, C. I. (2018). Deconstructing multivariate decoding for the study of brain
839    function. Neuroimage, 180, 4-18.

840    Hebb, D. O. (1949). The organization of behavior: a neuropsychological theory. New York: John
841    Wiley.

842    Hinton, G. E., Osindero, S., & Teh, Y. W. (2006). A fast learning algorithm for deep belief nets.
843    Neural computation, 18:7, 1527-1554.

844    Kassam, K. S., Markey, A. R., Cherkassky, V. L., Loewenstein, G., and Just, M. A. (2013).
845    Identifying emotions on the basis of neural activation. PloS one, 8:6. doi:
846    https://doi.org/10.1371/journal.pone.0066032

847    Kay, K.N., David, S.V., Prenger, R.J., Hansen, K.A. and Gallant, J.L. (2008). Modeling low-
848    frequency fluctuation and hemodynamic response timecourse in event-related fMRI. Human Brain
849    Mapping, 29:2, 142-156.

850    Kingma, D. P., and Ba, J. (2014). Adam: A method for stochastic optimization. arXiv [Preprint].
851    Available at: https://arxiv.org/abs/1412.6980 (Accessed September 17, 2020).

852    Koch, C., and Laurent, G. (1999). Complexity and the nervous system. Science, 284:5411, 96-98.

853    Kohler, P. J., Fogelson, S. V., Reavis, E. A., Meng, M., Guntupalli, J. S., Hanke, M., ... and Peter, U.
854    T. (2013). Pattern classification precedes region-average hemodynamic response in early visual
855    cortex. NeuroImage, 78, 249-260.

856    Krishnapuram, B., Figueiredo, M., Carin, L., and Hartemink, A. (2005). Sparse Multinomial Logistic
857    Regression: Fast Algorithms and Generalization Bounds. IEEE Transactions on Pattern Analysis and
858    Machine Intelligence (PAMI), 27, 957–968.

859    Lim, P. C., Ward, E. J., Vickery, T. J., and Johnson, M. R. (2019). Not-so-working memory: Drift in
860    functional magnetic resonance imaging pattern representations during maintenance predicts errors in
861    a visual working memory task. Journal of Cognitive Neuroscience, 31:10, 1520-1534.

862    Linnainmaa, S. (1970). The representation of the cumulative rounding error of an algorithm as a
863    Taylor expansion of the local rounding errors. [Master's thesis]. [Helsinki, Finland]: University of
864    Helsinki.

865    Lore, K. G., Akintayo, A., and Sarkar, S. (2017). LLNet: A deep autoencoder approach to natural
866    low-light image enhancement. Pattern Recognition, 61, 650-662.

867    Maas, A. L., Hannun, A. Y., and Ng, A. Y. (2013). Rectifier nonlinearities improve neural network
868    acoustic models. Proc. icml, 30:1, 3.

869    McCulloch, W.S., and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity.
870    Bulletin of Mathematical Biophysics 5, 115–133.

871    Minsky, M., and Papert, S. (1969). Perceptrons: An Introduction to Computational Geometry. MIT
872    Press.

873   Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., and Ng, A. Y. (2011). Multimodal deep learning.
874   ICML, 11, 689-696.

875   Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., and Kriegeskorte, N. (2014). A
876   toolbox for representational similarity analysis. PLoS computational biology, 10:4, e1003553.

877   Poldrack, R. A., Barch, D. M., Mitchell, J., Wager, T., Wagner, A. D., Devlin, J. T., ... and Milham,
878   M. (2013). Toward open sharing of task-based fMRI data: the OpenfMRI project. Frontiers in
879   neuroinformatics, 7, 12.

880   Raina, R., Madhavan, A., and Ng, A. Y. (2009). Large-scale deep unsupervised learning using
881   graphics processors. Proceedings of the 26th annual international conference on machine learning,
882   873-880.

883   Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and
884   organization in the brain. Psychological Review. 65:6, 386–408. doi:10.1037/h0042519.

885   Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986a). Learning representations by back-
886   propagating errors. Nature, 323, 533–536.

887   Rumelhart, D. E., Hinton, G. E.,  and Williams, R. J. (1986b). Learning internal representations by
888   error propagation. Parallel distributed processing: Explorations in the microstructure of cognition, 1.

889   Shalev-Shwartz, S., and Ben-David, S. (2014). Understanding machine learning: From theory to
890   algorithms. Cambridge University Press.

891   Sporns, O., Tononi, G., & Edelman, G. M. (2000). Connectivity and complexity: the relationship
892   between neuroanatomy and brain dynamics. Neural networks, 13(8-9), 909-922.

893   Szegedy, C., Ioffe, S., Vanhoucke, V., & Alemi, A. (2016). Inception-v4, inception-resnet and the
894   impact of residual connections on learning. arXiv preprint arXiv:1602.07261.

895   Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E., Yacoub, E., Ugurbil, K., & Wu-Minn
896   HCP Consortium. (2013). The WU-Minn human connectome project: an overview. Neuroimage, 80,
897   62-79.

898   Werbos, P. (1974). Beyond Regression: New Tools for Prediction and Analysis in the Behavioral
899   Sciences. [PhD thesis]. [Cambridge (MA)]: Harvard University.

900   Williams, J. M., Samal, A., Rao, P. K., and Johnson, M. R. (2020). Paired Trial Classification: A
901   Novel Deep Learning Technique for MVPA. Frontiers in Neuroscience, 14, 417.

902   Wolpert, D.H., Macready, W.G. (1997), "No Free Lunch Theorems for Optimization", IEEE
903   Transactions on Evolutionary Computation 1, 67

904   Xie, J., Girshick, R., and Farhadi, A. (2016). Unsupervised deep embedding for clustering analysis.
905   International conference on machine learning, 478-487.

906   Xue, G., Dong, Q., Chen, C., Lu, Z., Mumford, J. A., and Poldrack, R. A. (2010). Greater neural
907   pattern similarity across repetitions is associated with better memory. Science, 330:6000, 97-101.

908  Zeiler, M. D., and Fergus, R. (2014). Visualizing and understanding convolutional networks.
909  European conference on computer vision, 818-833.

910  Zeiler, M. D., and Fergus, R. (2013). Stochastic pooling for regularization of deep convolutional
911  neural networks. arXiv [Preprint]. Available at https://arxiv.org/abs/1301.3557 (Accessed September
912  17, 2020)

913  Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. Journal of the
914  royal statistical society: series B (statistical methodology), 67(2), 301-320.

915  **12  Data Availability Statement**

916  The toolbox code and sample data are available through a Git repository hosted at
917  https://bitbucket.org/delineate/delineate/src/master/. Release versions of the toolbox and additional
918  documentation, as well a link to the Git repository, can be found at http://delineate.it.