
Deep Learning: When Conventional Wisdom Fails to be Wise

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A major tenet of conventional wisdom dictates that models should not be over-
2 parameterized: the number of free parameters should not exceed the number of
3 training data points. This tenet originates from centuries of shallow learning, pri-
4 marily in the form of linear or logistic regression. It is routinely applied to all kinds
5 of data analyses and modeling and even to infer properties of the brain. However,
6 through a variety of precise mathematical examples, we show that this conventional
7 wisdom is completely wrong as soon as one moves from shallow to deep learning.
8 In particular, we construct sequences of both linear and non-linear deep learning
9 models whose number of parameters can grow to arbitrarily large values, and which
10 remain well defined and trainable using a fixed, finite size, training set. In deep
11 models, the parameter space is partitioned into large equivalence classes. Learning
12 can be viewed as a communication process where information is communicated
13 from the data to the synaptic weights. The information in the training data only can,
14 and needs to, specify an equivalence class of the parameters. It cannot, and does
15 not need to, specify individual parameter values. As such, the number of training
16 examples can be smaller than the number of free parameters.

17 1 Introduction

18 A long held form of conventional wisdom is that in order to train a model with n parameters one
19 should have at least n training examples, and preferably more. The origin of this statistical “dogma”
20 stems from linear regression and other forms of shallow learning¹. The soundness of this dogma
21 appears to be obvious from our experiences with linear regression: in general n examples are
22 necessary and sufficient in order to solve a system of n linear equations in n unknown variables. As a
23 result, the dogma is routinely repeated and used in myriads applications of statistics to modeling data
24 across all areas of human inquiry, often well beyond shallow learning, and to inspire a fear, if not a
25 disgust, for the so-called over-parameterized models. The dogma is also routinely used in a variety
26 of “back-of-the-envelope” calculations, for instance to infer properties or processing strategies for
27 the human brain. Here we show, through a variety of examples, that this central dogma is valid only
28 for shallow learning and that it is completely wrong when it comes to deep learning. Hence, in deep
29 learning it may not be unwise to get rid of the conventional wisdom entirely.

30 1.1 The Origin of the Dogma

31 For the past three centuries, since the discovery of least square linear regression by Gauss and
32 Legendre in the late 1700s (e.g. [9]), one of the most central dogma of statistics has been that a model

¹The mathematically correct distinction between shallow and deep learning is whether there are hidden units/layers or not

33 should not have more parameters than data points. There is little doubt that the origin of this dogma
34 lies in linear regression, or equivalently in linear systems of equations where in general if there are n
35 unknown variables one needs n linear equations (or training examples) to uniquely solve the system.
36 However, this is not a characteristic of linear systems alone. The same holds true immediately for
37 logistic regression. Since the logistic function is monotone increasing, it has a unique inverse and
38 by inverting the targets one can reduce logistic regression to a linear system. While this is true for
39 single linear or logistic neurons, the same result holds for a shallow layer of linear or logistic neurons,
40 since in this case each neuron operates and learns independently of all the other neurons. Similar
41 observations can be made for single-variable polynomial regression. Thus, in short, the origin of the
42 conventional wisdom can easily be traced back to shallow learning and basic results in linear algebra.

43 Not only the soundness of the dogma seems obvious from basic linear algebra considerations, but
44 its violation in shallow learning leads to two kinds of problems: (1) an over-parameterized shallow
45 model is not well defined, in the sense that its parameters are not uniquely determined by the data;
46 and, as a result, (2) such a model can overfit the data by achieving low error on the training data,
47 while performing poorly on held out data. Finally, the widespread aversion for over-parameterized
48 models stems also from our sense of elegance and simplicity, as embodied in the principle of Occam's
49 razor.

50 1.2 Applications of the Dogma

51 While the dogma makes sense for shallow learning situations, it is often applied to deep learning
52 situations. For instance, many articles have been published in the literature recommending that deep
53 learning models ought to have training sets that are 10 times [1] or 50 times [2] bigger than the
54 number of free parameters. Obviously these arbitrary, constant, and widely discarding prescriptive
55 multiplicative factors should be viewed with a grain of suspicion.

56 Another standard application of the dogma is to infer properties of complex, non-shallow systems,
57 like the brain. For instance, Geoff Hinton and others like to point out that the human brain has on
58 the order of say 10^{15} synapses, while human lives last on the order of 3×10^9 seconds. Assuming
59 one training example per second, or even 1000 training examples per second, the brain does not have
60 enough training examples to train its army of synapses. From this false premise, one may draw all
61 kinds of conclusion from "the brain must be doing something special" to "the majority of synapses
62 must be hardwired". However, as we shall see, all these conclusions are worthless: they may be false
63 or true, since they are derived from a false premise. The false premise is obtained by applying a
64 statistical principle, correctly observed in shallow learning situations, to deep learning situations.

65 2 Preliminary Evidence against the Dogma

66 Preliminary evidence that something may be wrong with the dogma comes from at least three
67 directions: Bayesian statistical theory, statistical ensembles, and deep learning practice.

68 From a purely Bayesian perspective, selecting the complexity of a model based on the amount of
69 training data makes no sense at all, as there is in general no relationship between the two. Using a
70 prior that favors simple models may be convenient, or satisfy tradition, however there is no intrinsic
71 epistemological reason for selecting such a prior. If anything, a situation with few data points may be
72 the sign that data are hard or expensive to acquire. In turn, this is possibly the sign of an underlying
73 complex phenomena, which may call for a complex model rather than a simple one. Using a prior
74 that favors models with few parameters is analogous to the paradigm of searching for one's car keys
75 at night under the only lamp present in a dark parking lot: there is no epistemological reason for the
76 keys to be under the lamp. But what about Occam's razor? As noted in [10, 11], such a prior is
77 not needed to implement Occam's razor which naturally emerges from the Bayesian framework. To
78 see this in a simple way, imagine having an overall class of models comprising two sub-classes of
79 models: simple models (S) and complex models (C). Imagine that a priori one has no preference
80 between the two classes S and C , and likewise that within each class one has no preference among
81 the models in that class. Let s and c denote the value of the constant prior probability shared by all
82 the models in class S and in class C respectively. Thus the overall prior distribution must satisfy

$$s|S| + c|C| = 1$$

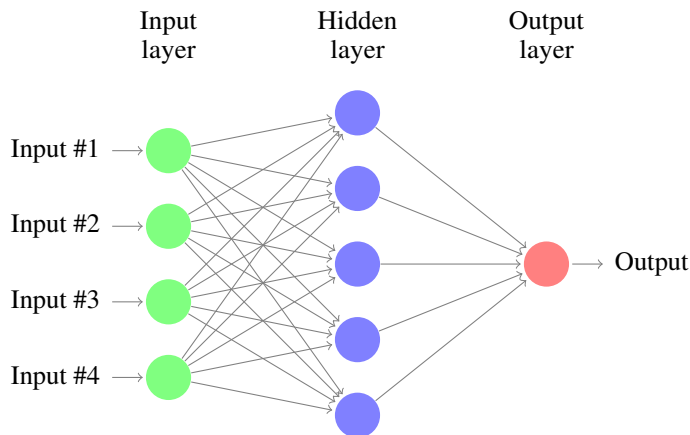


Figure 1: A $A(4, 5, 1)$ architecture.

83 where $|S|$ and $|C|$ represent the volumes of the corresponding classes. Because the complex models
 84 have more parameters, in general $|S| \ll |C|$. As a result, we must have: $s \gg c$. In short, simple
 85 models will automatically have a much higher prior probability, and this effect will tend to be reflected
 86 also in the posterior probabilities.

87 A second line of evidence against the soundness of the dogma comes from the widespread use,
 88 and recognized effectiveness, of statistical ensembles, where many different models are combined
 89 together, for instance through a simple weighted averaging operation. This combination alone
 90 generally results in a deep overall model, even if the individual models are shallow. And even if
 91 the number of parameters of each individual model satisfies the dogma, obviously as the number of
 92 models in the ensembles is increased, there is a point where the overall model starts to violate the
 93 dogma. Perhaps surprisingly, the over-parameterization aspect of ensembles does not seem to have
 94 systematically worried statisticians.

95 Finally, and perhaps most importantly, it has been observed several times that in deep learning practice
 96 that over-parameterized models can work well, with no significant sign of overfitting. However, this
 97 phenomena has been used either to criticize deep learning, or is regarded as some kind of oddity or a
 98 mystery (e.g. [12, 13, 8]), possibly requiring novel strategies for combating the over-fitting curse.

99 Here we set out to prove why the conventional wisdom is simply wrong when it comes to deep
 100 learning. In particular we give several examples of large networks with many parameters that can be
 101 trained with far fewer examples in both the linear and non-linear cases. We consider primarily the
 102 supervised learning framework, but through the use of autoencoder architectures we show that the
 103 same basic ideas can be applied to the unsupervised, or semi-supervised, learning frameworks. At
 104 the linear end of the spectrum of models, we look at deep, fully-connected, linear networks. At the
 105 other extreme non-linear end of the spectrum, we look at deep, fully-connected, unrestricted Boolean
 106 networks. And in the middle of the spectrum, we look at deep fully-connected networks of linear
 107 threshold gates.

108 **Notation:** We use the notation $A(n_0, n_1, \dots, n_L)$ to denote a deep feedforward architectures with
 109 n_i units in layer i , where the input layer is layer 0 and the output layer is layer L (Figure 1).

110 3 The Linear Regime

111 Deep feed-forward linear networks have been studied for quite some time (e.g. [4, 7, 5, 6]) in the
 112 context of least square linear regression. One of the main theoretical results is that, in the fully-
 113 connected case, the error functions of these networks does not have any spurious local minima. All
 114 the critical points where the gradient of the error function is zero are either global minima or saddle
 115 points. As a result, properly applied stochastic gradient descent will tend to converge to a global
 116 minimum. The structure of the global minima and the saddle points can be understood in terms of
 117 Principal Component Analysis (CS). Clearly, as the depth of these models is increased the number



Figure 2: An $A(1, \dots, 1)$ architecture with L single-layer neurons. There is a single synaptic weight w_i connecting neuron $i - 1$ to neuron i . In the linear case, with no biases, the input-output function is given by $y = Px$ where P is the product of all the synaptic weights. While the number of parameters L can be arbitrarily large, a single training example is sufficient to constrain the value of the multiplier. Gradient descent rapidly converges onto an optimal solution where the product of the synaptic weight has the optimal value: $P = \alpha/\beta$ where $\alpha = E(xt)$ and $\beta = E(x^2)$ (see text).

118 of parameters can grow to infinity. But what are the requirements on the size of the corresponding
 119 training sets?

120 3.1 The Simplest Deep Linear Model

121 To begin with, we consider an architecture $A(1, 1, \dots, 1)$, with a single linear neuron in each layer
 122 (Figure 2). For simplicity we assume that there are no biases, but the same analysis can easily be
 123 extended to the case with biases. The weights are w_1, \dots, w_L and the neural network behaves as a
 124 multiplier, in the sense that given an input x the output is simply:

$$y = Px \quad \text{with} \quad P = \prod_i w_i$$

125 This is a deep linear regression architecture with L parameters. The supervised training data consists
 126 of input-target pairs of the form (x, t) that provide information about what the overall multiplier P
 127 should be. Taking expectations over the training data, let $E(tx) = \alpha$ and $E(x^2) = \beta$. The error
 128 \mathcal{E} is the standard least square error. It is easy to check that the error is convex in P and that at the
 129 optimum one must have $\alpha - \beta P = 0$ or $P = \alpha/\beta$. It can be shown (see [3]) that, except for trivial
 130 cases, given any initial starting point, gradient descent, or even random backpropagation (feedback
 131 alignment), will converge to a global minimum satisfying $P = \alpha/\beta$.

132 While the architecture has an arbitrary large number of parameters L , in principle a single training
 133 example is sufficient to determine the value of the correct multiplier. The value of the overall product
 134 P partitions the space of synaptic weights into equivalence classes: all the architectures which
 135 produce the same value P are equivalent. The training data need only to provide enough information
 136 for selecting one equivalence class, but not the value of the individual weights within the equivalence
 137 class. Thus there is a manifold of equivalent solutions satisfying the optimal relationship $P = \alpha/\beta$
 138 and the volume of this manifold grows with the number L of parameters. However the training set
 139 can remain as small as a single training example, a clear violation of the dogma.

140 Of course, here and everywhere else in the following examples, one may wonder what could be the
 141 purpose of having L layers, when a single layer could be sufficient to implement the same overall
 142 input-output function. There could be multiple purposes. The most obvious one is that the volume
 143 of the solutions grows with the depth of the architectures and this may facilitate learning. But in
 144 addition, one must also think about the possible constraints that may be associated with physical
 145 neural systems, as opposed to the virtualized simulations of neural systems we routinely carry on
 146 our digital computers using the likes of Keras, PyTorch, and TensorFlow. For example, even in the
 147 simplest linear case described above, imagine that the overall desired multiplier is $P = 2^{10} = 1024$

148 but that the individual synaptic weights connecting one neuron to the next are bounded in the $[-2, +2]$
149 range. Then no architecture with less than 10 layers is capable of implementing the optimal input-
150 output function. Deeper architectures are needed to implement the overall optimal function and to
151 robustly distribute the load across multiple synapses.

152 3.2 Deep Linear Models with No Bottlenecks

153 At first sight, one may be tempted to think that the example above is due to the fact that there is a
154 single neuron per layer. However, this is not the case and exactly the same phenomena is observed
155 for a linear regression architecture of the form $A(n, n, \dots, n)$ where all the layers have size n
156 and the weights are given by matrices W_1, \dots, W_L . Again, in vector-matrix form, the input-output
157 relationship is given by:

$$y = Px \quad \text{with} \quad P = W_L W_{L-1} \dots W_1$$

158 Again it is easy to see that this architecture has Ln^2 parameters. The overall input-output function
159 corresponds to a single $n \times n$ matrix P . But in order to specify such a linear map, we only need
160 to specify the images of the canonical basis of \mathbb{R}^n , in other words, n training examples in general
161 position are sufficient, again violating the dogma.

162 Note that this property remains true if the architecture also contains expansive hidden layers of size
163 greater than n , or if the input and output layers have different sizes and all the hidden layers have size
164 greater than the input and output layer (i.e. the hidden layers do not affect the rank of the optimal
165 overall input-output function).

166 3.3 Deep Linear Models with Bottlenecks

167 In the previous two examples, all the layers have the same size, or are expansive. However it is easy to
168 relax this assumption and consider compressive architectures. To begin with, consider a purely linear
169 compressive autoencoder architecture of the form $A(n, m, n)$, with $m < n$ (Figure 3). In this case,
170 the bottleneck layer imposes a rank restriction on the overall transformation. It is well known [4] that
171 not only the quadratic error function of such an autoencoder has no spurious local minima, but all
172 its critical points correspond, up to changes of coordinates in the hidden layer, to projections onto
173 subspaces spanned by eigenvectors of the data covariance matrix. The global minima is associated
174 with Principal Component Analysis using projections onto a subspace of dimension m . Obviously
175 one can include additional linear layers of size greater or equal to m between the input layer and
176 the bottleneck layer, or between the bottleneck layer and the output layer, arbitrarily increasing the
177 total number of parameters, but without affecting the essence of the optimal solution. The minimal
178 training set to specify the optimal solution consists of m vectors of size n to specify the project
179 hyperplane, providing another egregious violation of the dogma. Again there are large equivalence
180 classes of parameters associated with the same overall performance (e.g. in the linear case with a
181 single bottleneck, we have $P = AB = ACC^{-1}B$; thus the overall map P is defined up to invertible
182 transformations applied to the hidden layer). The results in [4, 7] show that the same observations
183 can be made for arbitrary fully connected deep linear architectures (i.e. beyond autoencoders) and
184 not only in the real-valued case, but also in the complex-valued case [6].

185 All the previous examples correspond to linear networks. Thus one may be misled to think that
186 the analyses apply only to linear networks. Next we show that exactly the same phenomena can be
187 observed in non-linear deep architectures. Among the non-linear model to be discussed, we will
188 examine first the most non-linear model of all which is the unrestricted Boolean model, where each
189 neuron implements a Boolean function, with no restrictions on the kinds of Boolean functions. An
190 unrestricted Boolean neuron with n inputs implements a function f with 2^n parameters, since one
191 must specify one binary value for each of the 2^n possible entries of the truth table of f . Then we will
192 consider also the case of Boolean neurons implemented by linear threshold functions, or perceptrons.

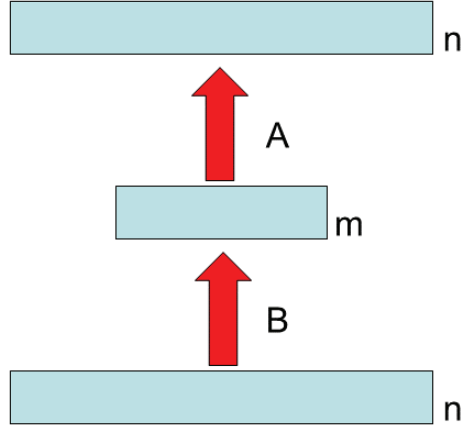


Figure 3: An $A(n, m, n)$ compressive ($m < n$) autoencoder architecture. In the linear case, the transformations A and B correspond to matrices and the overall linear transformation P is given by: $y = Px = ABx$.

193 4 The Non-Linear Regime: Unrestricted Boolean Model

194 4.1 The Simplest Deep Non-Linear Model

195 We can use the same architecture $A(1, \dots, 1)$ as in the first example above. In the Boolean unrestricted
 196 model, each Boolean function from one neuron to the next is either the identity or the negation
 197 (Boolean NOT function). So there is one binary degree of freedom associated with each layer and
 198 again the number of degrees of freedom grows linearly with the depth. The overall input-output
 199 function is either the identity, or the negation, and a single training example is sufficient to establish
 200 whether the overall function ought to be the identity or the negation of the identity. If the architecture
 201 contains an even number of negations the overall input-output function is the identity, and if the
 202 architecture contains an odd number of negations, the overall input-output function is the negation.
 203 Thus again the dogma is violated.

204 To get a slightly more interesting non-linear example, we can use the same architecture $A(1, \dots, 1)$
 205 as in the first example above, with L weights w_1, \dots, w_L . The difference is that all the neurons have
 206 a non-linear activation function $g(x) = x^2$ (more generally we could use for instance $g(x) = x^k$).
 207 Thus the overall input-output function is given by:

$$y = (w_L \dots ((w_2 w_1 x)^2) \dots)^2 = w_L^2 w_{L-1}^4 \dots w_1^{2L} x^{2^L}$$

208 or

$$y = Px^{2^L} \quad \text{with} \quad P = \prod w_i^{2^{L-2i+2}}$$

209 Thus in this case the multiplier P realized by the architecture is positive. Again the number of
 210 parameters is L and it can be arbitrarily large. As in the linear case, a single training example of
 211 the form (x, t) is sufficient to determine the multiplier P , with a manifold of equivalent solutions
 212 corresponding to parameters satisfying $P = \prod w_i^{2^{L-2i+2}} = \alpha/\beta$, with this time $\alpha = E(tx^2)$ and
 213 $\beta = E(x^4)$, when $\alpha > 0$. If $\alpha < 0$, the optimum is obtained for $P = 0$ which can be achieved by
 214 having at least one of the weights of the architectures equal to zero. In short, in both examples treated
 215 in this subsection, the dogma is again violated.

216 4.2 Deep Non-Linear Models with No-Bottlenecks (Unrestricted Boolean)

217 Consider an architecture $A(n_0, \dots, n_L)$ where each neuron can implements any Boolean function
218 of the neurons in the previous layer. The error function is the Hamming distance between target
219 and output vectors. For simplicity, let us first assume that all the layers have the same size n . The
220 overall input-output function is a Boolean map from \mathbb{H}^n to \mathbb{H}^n , where \mathbb{H}^n denotes the n -dimensional
221 hypercube. This architecture has $Ln2^n$ parameters, since each unrestricted Boolean neuron with
222 n inputs has 2^n free parameters. The overall input-output map can be specified using only $n2^n$
223 examples. It can easily be implemented with 0 error through a large class of equivalent networks. As
224 the number of layers L goes to infinity the number of parameters goes to infinity, while the number
225 of required training examples remains fixed and is determined entirely by the size of the input and
226 output layers. This can easily be generalized to a Boolean unrestricted architecture of the form
227 $A(n_0, \dots, n_L)$, as long as there are no bottleneck layers. In such an architecture, the total number of
228 parameters is given by: $\sum_{i=1}^L n_i 2^{n_i-1}$. The number of necessary and sufficient training examples
229 needed to specify the overall input-output function is given by: $n_L 2^{n_0}$, and thus again the dogma is
230 violated. The case with bottle-neck layers is treated below.

231 4.3 Deep Non-Linear Models with Bottlenecks (Unrestricted Boolean)

232 For simplicity, consider first an unrestricted Boolean compressive autoencoder with architecture
233 $A(n, m, n)$ and $m < n$. The error function is the Hamming distance between the input vector and
234 the output vector. The hidden layer can have 2^m states. Thus if the number of training examples is
235 at most 2^m , it can be realized by the architecture with 0 Hamming distortion, since every input can
236 be mapped to a unique hidden representation and the corresponding representation can be mapped
237 back to the same input using unrestricted Boolean gates. Obviously if additional layers of size at
238 least m are added between the input layer and the hidden layer, or between the hidden layer and the
239 output layer, the number of parameters can be arbitrarily increased, while maintaining the same fixed
240 training set and the ability to implement it exactly with no Hamming distortion. Thus in this regime
241 the dogma is again violated.

242 In the more interesting regime where the number of training examples exceeds 2^m , then there must
243 be clusters of training examples that are mapped to the same hidden representation. It is easy to see
244 that for optimality purposes the corresponding representation must be mapped to the binary vector
245 closest to the center of gravity of the cluster, essentially the majority vector, in order to minimize
246 the Hamming distortion. Thus, in short, in this regime the optimal solution corresponds to a form of
247 optimal clustering with respect to the Hamming distance with, in general, 2^m clusters. As a back of
248 the envelope calculation, assuming the clusters are spherical, these can be described by providing
249 two points corresponding to a diameter. Thus in principle a training set of size $2 \times 2^m = 2^{m+1}$ could
250 suffice. The number of parameters of the architecture is given by: $m2^n + n2^m$ which far exceeds
251 the number of training examples. And even without the assumption of spherical clusters, it is clear
252 that the number of parameters far exceeds the number of training examples, and that the gap can be
253 made as large as possible, just by adding additional layers of size at least m between the input and
254 the hidden layer, or the hidden layer and the output layer. Thus again the dogma is grossly violated.

255 Finally, we turn to deep non-linear architecture where the neurons are linear or polynomial threshold
256 gates. Linear threshold neurons, or perceptrons, are very similar to sigmoidal (e.g. logistic) neurons.

257 5 The Non-Linear Regime: Linear or Polynomial Threshold Gates

258 Here each neuron in the architecture is a linear or polynomial threshold function of degree d . In
259 the linear threshold case ($d = 1$), any neuron with n inputs $x = (x_1, \dots, x_n)$ produces an output
260 equal to $\text{sign}(\sum_i w_i x_i)$ in the $-/+$ case; or $H(\sum_i w_i x_i)$ in the 0/1 case, where H denotes the
261 Heaviside function. Such a neuron has n synaptic parameters. In the polynomial case of degree d ,
262 the output of a neuron has the form $\text{sign}(p(x))$ in the $-/+$ case; or $H(p(x))$ in the 0/1 case, where
263 $p(x) = p(x_1, \dots, x_n)$ is a polynomial of degree d . The number of parameters of a polynomial
264 threshold neuron increases accordingly. As usual a bias can also be added or, equivalently, one of the
265 input variables is considered to be constant and equal to 1.

266 5.1 The Simplest Deep Non-Linear Model with Linear or Polynomial Threshold Gates

267 We can use the same architecture $A(1, \dots, 1)$ as in the first example above. Linear or polynomial
268 threshold neurons can realize the identity and the negation, depending on whether the corresponding
269 incoming weight is positive or negative. So the result here is similar to the Boolean unrestricted
270 case. For instance for linear threshold gates, without the bias, the number of parameters is equal to L .
271 The number of negative weights determines how many negations are present in the chain. A single
272 input-output example determines whether the overall chain should be the identity or the negation.
273 Thus again the dogma is violated.

274 5.2 Deep Non-Linear Models with Bottlenecks (Linear or Polynomial Threshold Gates)

275 We can again start with a compressive autoencoder architecture with shape $A(n, m, n)$ and $m < n$
276 and linear threshold neurons with the Hamming error function. In the most interesting case where the
277 number of examples exceeds 2^m , then the optimal solution corresponds to the optimal approximation
278 to the optimal Hamming clustering that can be achieved using linear threshold gates. The number of
279 parameters of this architecture is $2nm$ which is not necessarily less than the number 2^{m+1} of required
280 training examples, under the spherical cluster assumption. However, as in the similar previous
281 examples, the number of parameters can be increased arbitrarily by adding additional layers of size
282 at least m between the input and the hidden layer, or between the hidden layer and the output layer.
283 Thus once again there are large equivalence classes in parameter space (e.g. applying permutations to
284 the neurons in a given layer) and the dogma is grossly violated.

285 6 Discussion

286 The conventional dogma that models ought to have less parameters than the number of training
287 examples is a mere product of shallow learning. It arises, and should be applied, only in shallow
288 learning situations. As soon as one moves to deep learning situations, the dogma becomes non-sense
289 and all the expectations it creates are simply wrong, *even in the linear case*. It is simply time to
290 think about deep models in a different way, without the expectation that over-parameterization must
291 necessarily lead to over-fitting. This is not to say, of course, that over-parameterized deep learning
292 models cannot overfit, but expecting them to do so just because they are over-parameterized is unwise
293 and unnecessary.

294 Over-parameterized models tend to partition the parameter space into large equivalence classes.
295 All the parameter settings within one class are equivalent in terms of overall performance. Neural
296 learning can be viewed as a communication process where information is communicated from the
297 training data to the synaptic weights. The training data needs to contain enough information to select
298 one of the equivalence classes, but not any particular setting of the weights within that class. Thus
299 the information needed to specify one equivalence class is much less than the information required
300 to specify a particular setting of the weights. And this explains why the number of data points can
301 be much less than the number of parameters. Furthermore, the structure of the deep models and the
302 partitioning into equivalence classes is such that it is not even possible for the training data to be able
303 to specify each individual weight of the architecture. This is because the system cannot distinguish
304 between two different settings of the parameters within the same equivalence class. For instance,
305 once the optimal class is achieved with a particular setting of the weights, the gradient of the error
306 is zero and there is no way of exploring or distinguishing other optimal architectures in the same
307 equivalence class.

308 Shallow learning, in particular linear regression, already contains many of the central themes of
309 machine learning: from the use of a parameterized family of models, to model fitting by error
310 minimization, to prediction and so forth. However, when transitioning to deep learning, linear
311 regression is misleading in three major aspects. First, it has an analytic closed-form solution. Second,
312 it is interpretable (or visualizable, at least in low dimensions). Third, it requires that the number of
313 training examples be equal or even exceed the number of parameters in order to completely determine
314 the solution. The first two points are now well established and accepted. We use stochastic gradient
315 descent for deep learning model fitting and almost no one cares about not having a closed-form
316 analytic solution. Likewise, no one expects to be able to easily visualize complex non-linear surfaces
317 in high-dimensional spaces, although many are still working on various other issues related to

318 interpretability. However, we are still struggling with the third point. It is time to move on this front
319 too.

320 It should be clear from the examples presented that one of the emergent characteristics of over-
321 parameterized regimes is the existence of large equivalence classes in parameter space, all associated
322 with roughly the same level of overall performance. The training data needs only to provide enough
323 information to select one of the equivalence classes (at the relevant quantization level), and not
324 to specify the value of each one of the parameters. Reflecting back on the human brain, most
325 mature human brains can pass the Turing test and achieve some form of general intelligence using
326 architectures that are similar, at least at the macroscopic level, and at the level of the basic hardware
327 components (e.g. pyramidal cells), but presumably with significant differences at the level of
328 individual synapses.

329 Finally, there is the question of when deep architectures overfit the data. The results presented
330 here provide a clear answer. Consider an architecture with w parameters. At the proper level of
331 quantization of the weights and the error function, the architecture may partition the space of weights
332 into e equivalence classes. Thus $\log_2 e$ bits are needed to specify one of the equivalence classes. If the
333 training data provides less than $\log_2 e$ bits of information, then it does not contain enough information
334 to select a relevant equivalence class and overfitting may occur. If the training data provides $\log_2 e$
335 bits of information to select an equivalence class, then there is no overfitting and providing more data
336 is not necessary. In the case of a classification architecture with independent binary inputs of length
337 n , k training examples contain on the order of kn bits of information. Thus the important question is
338 not whether $k \approx w$ (conventional wisdom) but whether $kn \approx \log_2 e$.

339 References

- 340 [1] Yaser S Abu-Mostafa. Hints. *Neural computation*, 7(4):639–671, 1995.
- 341 [2] Ahmad Alwosheel, Sander van Cranenburgh, and Caspar G Chorus. Is your dataset big enough?
342 sample size requirements when using artificial neural networks for discrete choice analysis.
343 *Journal of choice modelling*, 28:167–182, 2018.
- 344 [3] P. Baldi. *Deep Learning in Science*. Cambridge University Press, Cambridge, UK, 2021.
- 345 [4] P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from
346 examples without local minima. *Neural Networks*, 2(1):53–58, 1989.
- 347 [5] P. Baldi and K. Hornik. Learning in linear networks: a survey. *IEEE Transactions on Neural
348 Networks*, 6(4):837–858, 1994. 1995.
- 349 [6] P. Baldi and Z. Lu. Complex-valued autoencoders. *Neural Networks*, 33:136–147, 2012.
- 350 [7] Pierre Baldi. Linear learning: Landscapes and algorithms. *Advances in neural information
351 processing systems*, 1, 1988.
- 352 [8] Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model
353 alignment explain generalization in kernel regression and infinitely wide neural networks.
354 *Nature communications*, 12(1):1–12, 2021.
- 355 [9] Adrien Marie Legendre. *Nouvelles méthodes pour la détermination des orbites des comètes*. F.
356 Didot, 1805.
- 357 [10] D. J. C. MacKay. Bayesian interpolation. *Neural Computation*, 4:415–447, 1992.
- 358 [11] D. J. C. MacKay. A practical Bayesian framework for backprop networks. *Neural Computation*,
359 4:448–472, 1992.
- 360 [12] Terrence J Sejnowski. The unreasonable effectiveness of deep learning in artificial intelligence.
361 *Proceedings of the National Academy of Sciences*, 117(48):30033–30038, 2020.
- 362 [13] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
363 deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–
364 115, 2021.

365 **Checklist**

- 366 1. For all authors...
- 367 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
368 contributions and scope? [Yes]
- 369 (b) Did you describe the limitations of your work? [Yes]
- 370 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 371 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
372 them? [Yes]
- 373 2. If you are including theoretical results...
- 374 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 375 (b) Did you include complete proofs of all theoretical results? [Yes]
- 376 3. If you ran experiments...
- 377 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
378 mental results (either in the supplemental material or as a URL)? [N/A]
- 379 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
380 were chosen)? [N/A]
- 381 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
382 ments multiple times)? [N/A]
- 383 (d) Did you include the total amount of compute and the type of resources used (e.g., type
384 of GPUs, internal cluster, or cloud provider)? [N/A]
- 385 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 386 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 387 (b) Did you mention the license of the assets? [N/A]
- 388 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 389
- 390 (d) Did you discuss whether and how consent was obtained from people whose data you're
391 using/curating? [N/A]
- 392 (e) Did you discuss whether the data you are using/curating contains personally identifiable
393 information or offensive content? [N/A]
- 394 5. If you used crowdsourcing or conducted research with human subjects...
- 395 (a) Did you include the full text of instructions given to participants and screenshots, if
396 applicable? [N/A]
- 397 (b) Did you describe any potential participant risks, with links to Institutional Review
398 Board (IRB) approvals, if applicable? [N/A]
- 399 (c) Did you include the estimated hourly wage paid to participants and the total amount
400 spent on participant compensation? [N/A]