

An Improved Bound on the VC-Dimension of Neural Networks with Polynomial Activation Functions

M. Vidyasagar
Advanced Technology Centre
Tata Consultancy Services
1-2-10 S. P. Road
INDIA 500 003
sagar@atc.tcs.co.in

J. Maurice Rojas
Department of Mathematics
Texas A&M University
College Station, TX 77843-3368
U.S.A.
rojas@math.tamu.edu

November 21, 2018

Abstract

We derive an improved upper bound for the VC-dimension of neural networks with polynomial activation functions. This improved bound is based on a result of Rojas [Roj00] on the number of connected components of a semi-algebraic set.

1 Introduction

We examine neural networks with polynomial activation functions. The specific architecture of the neural networks is described in detail in the next section. Such neural networks have been the subject of active investigation for several years, since powerful tools from algebraic geometry can be brought into play in analyzing the VC-dimension of such networks. Perhaps [GJ93] was the first paper to connect these two subjects. For several years (see for example [KM97]) it has been known that every bound on the number of connected components of a semi-algebraic set can be readily translated into a corresponding bound on the VC-dimension of a neural network architecture. Practically all of the known bounds on the VC-dimension of neural networks with polynomial activation bounds make use of a classical result discovered by Milnor, Oleinik, Petrovsky, and Thom [OP49, Mil64, Tho65].¹ This bound, while easy to use, is usually much larger than necessary, since it only uses coarse information about the underlying set such as the number of variables and the maximum degree of the input polynomials. More recently, sharper bounds using more refined data from the input polynomials have been discovered. In the present note, we use a result due to Rojas [Roj00] that is particularly well-suited to neural networks with polynomial activation functions. The present bound is, *in all cases*, sharper than the earlier bound of Goldberg and Jerrum [GJ93]. Moreover, it is intuitively appealing, as the improvement can be quantified as the relative entropy of two probability vectors, whose dimension equals the number of layers in the neural network. This shows that the problem of bounding the VC-dimension of a neural network architecture continues to be interesting, and that we should strive to derive even tighter upper bounds.

Our main result is stated in theorem 4 of section 4. The recent semi-algebraic bound it is based on is stated as theorem 3 of section 3. However, let us first review a bit of background and some of the earlier bounds.

2 Known Results

The following definition of the VC-dimension is standard; see for example the books by Vapnik [Vap95] or Vidyasagar [MV97].

Definition 1 *Suppose X is a set and \mathcal{F} is a collection of $\{0, 1\}$ -valued functions on X . A set $S = \{x_1, \dots, x_n\} \subseteq X$ is said to be **shattered** by \mathcal{F} if each of the 2^n functions mapping*

¹Actually, this result bounds the sum of the Betti numbers of a semi-algebraic set, and this quantity is always at least as large as the number of connected components. In practice, one usually only needs an upper bound on the number of connected components.

S into $\{0, 1\}$ is the restriction to S of some function in \mathcal{F} . The **Vapnik-Chervonenkis (VC)-dimension** of \mathcal{F} is the largest integer n such that there exists a set of cardinality n that is shattered by \mathcal{F} . \diamond

By identifying a $\{0, 1\}$ -valued function with its support set, it is also possible to speak of the VC-dimension of a collection of sets. In the sequel, we shall use both notions interchangeably.

Following by now familiar approaches, we view a *neural network* as a verifier of formulas. Specifically, let $\mathbf{w} \in \mathbb{R}^k$ denote the “weight vector” or the vector adjustable parameters in a neural networks. A neural network with input space $X \subseteq \mathbb{R}^N$ and weight vector \mathbf{w} evaluates a logical proposition $\phi(\mathbf{x}, \mathbf{w})$ which is a Boolean combination of s atomic expressions of the form $\tau_i(\mathbf{x}, \mathbf{w}) = 0$ or $\tau_j(\mathbf{x}, \mathbf{w}) > 0$. Letting 1 (resp. 0) denote “true” (resp. “false”), we can thus think of ϕ as a function from \mathbb{R}^{k+N} to $\{0, 1\}$. So for each weight vector \mathbf{w} , define

$$A_{\mathbf{w}} := \{\mathbf{x} \in \mathbb{R}^N : \phi(\mathbf{w}, \mathbf{x}) = 1\}.$$

The objective is to obtain an upper bound on the VC-dimension of the collection of sets $\mathcal{A} := \{A_{\mathbf{w}} \mid \mathbf{w} \in \mathbb{R}^k\}$ or, equivalently, the VC-dimension of the collection of $\{0, 1\}$ -valued functions $\Phi := \{\phi(\mathbf{w}, \cdot) \mid \mathbf{w} \in \mathbb{R}^k\}$.

To state the result, we need one final bit of notation: Let $\mathbf{x}_1, \dots, \mathbf{x}_v \in \mathbb{R}^N$, and suppose $sv \geq k$. From the sv polynomials $\tau_j(\cdot, \mathbf{x}_i)$ determined by all $(i, j) \in \{1, \dots, v\} \times \{1, \dots, s\}$, choose $r \leq k$ polynomials, and label them $\theta_1(\cdot), \dots, \theta_r(\cdot) : \mathbb{R}^k \rightarrow \mathbb{R}$. Define

$$\mathbf{f}(\mathbf{w}) := [\theta_1(\mathbf{w}) \dots \theta_r(\mathbf{w})] \in \mathbb{R}^r.$$

Finally, let B denote the maximum number of connected components of any pre-image $f^{-1}(y)$ with $y \in \mathbf{R}^r$, for any choice of r and $\theta_1, \dots, \theta_r$ as above. With the above set-up, the following result is proved in [KM97]. For further background on our setting below see [KM97] or [MV97], p. 329.

Theorem 1 *Following the notation above, assume further that restrict to those $y \in \mathbf{R}^r$ that are regular values of f . Then*

$$\text{VCDIM}(\Phi) \leq 2 \lg B + 2k \lg(2es).$$

That B is in fact finite and admits an explicit upper bound is obtained by appealing to the aforementioned classical result of Milnor, Oleinik, Petrovsky, and Thom [OP49, Mil64, Tho65], which we now state as follows:

Lemma 1 *Suppose $\theta_1, \dots, \theta_r$ are polynomials in k variables, with degree no larger than d . Then whenever \mathbf{y} is a regular value of \mathbf{f} as defined above, the preimage $\mathbf{f}^{-1}(\mathbf{y})$ contains no more than $d(2d)^{k-1}$ connected components.*

Note that Milnor actually proves the theorem in the case where $\mathbf{y} = \mathbf{0}$, but we can clearly perturb the constant terms of the θ_i to enforce this assumption. If we replace the quantity $d(2d)^{k-1}$ by the larger number $(2d)^k$ and substitute $B = (2d)^k$ into the upper bound (1), we get the following result.

Theorem 2 *Suppose ϕ is a Boolean formula involving a total of s polynomial equalities and inequalities, where each polynomial has degree no larger than d with respect to \mathbf{w} . Then $\text{VCDIM}(\Phi) \leq 2k \lg(4eds)$.*

The above result is the same as that derived in [GJ93]. It should be noted, however, that Goldberg and Jerrum actually consider neural networks with *piecewise polynomial* activation functions. With more elaborate notation, their results can be derived as special cases of Theorem 1.

Theorem 1 shows the importance of deriving *tight* upper bounds on the number of connected components of a semi-algebraic set. This is a long-standing problem in real algebraic geometry that has received considerable attention from the research community. It is obvious from the bound (1) that any improvement over Milnor's upper bound translates *directly* into a corresponding improvement in the estimate of the VC-dimension of a neural network architecture with polynomial activation functions. This leads us to the next topic.

3 Improved Upper Bound on the Number of Connected Components

In [Roj00], an improvement is provided over Milnor's bound. To state this improved result, a bit of notation is introduced.

Let Δ_n denote the standard n -simplex in \mathbf{R}^n , with vertices the standard basis vectors and the origin. Note that

$$d\Delta_n := \left\{ (x_1, \dots, x_n) \in \mathbf{R}^n \mid x_i \geq 0 \text{ for all } i \text{ and } \sum_{i=1}^n x_i \leq d \right\}.$$

Let $\text{Vol}_n(\cdot)$ denote the renormalization of the usual volume in \mathbf{R}^n satisfying $\text{Vol}_n(\Delta_n) = 1$. (Since the usual n -dimensional volume is multiplicative for orthogonal subspaces, it is easy to prove by induction that Vol_n is just $n!$ times the usual n -dimensional volume.)

Theorem 3 *Suppose τ_1, \dots, τ_r are polynomials in (w_1, \dots, w_k) , and let $\mathbf{e}_1, \dots, \mathbf{e}_k$ and \mathbf{O} denote the standard basis vectors and the origin of \mathbf{R}^k . Also, let V denote the convex hull of*

the union of $\{\mathbf{0}, e_1, \dots, e_k\}$ with the set of all (i_1, \dots, i_k) such that $w_1^{i_1} \dots w_k^{i_k}$ is a monomial of some $\theta_j(\cdot)$. Then

$$B \leq 2^k \text{Vol}_k(V).$$

In the special case where every k -tuple with $\sum_{j=1}^k i_j \leq d$ occurs in V , we recover the (adjusted) Milnor bound $(2d)^k$. However, the whole point of the preceding refined bound is that there are many instances where the input polynomial are far more sparse, and this can be exploited.

4 Improved Upper Bound on the VC-Dimension

In this section, we derive an improved upper bound on the VC-dimension of neural networks with polynomial activation functions. The improved bound is a direct consequence of coupling Theorems 1 and 3.

Let us begin by describing the class of neural networks under study. It is assumed that the network has N real inputs denoted by x_1, \dots, x_N . There are l levels in the network, and at level i there are q_i output neurons; however, at the output layer (level l) there is only a single neuron (see below). Let k_i denote the number of adjustable parameters, or “weights,” at level i , and let $k = \sum_{i=1}^l k_i$ denote the total number of adjustable parameters. Let $\mathbf{w}_i := (w_{i,1}, \dots, w_{i,k_i})$ denote the weight vector at level i , and $\mathbf{w} = (\mathbf{w}_1 \dots \mathbf{w}_l)$ denote the total weight vector. The input-output relationship of each neuron at level i is of the form

$$y_{i,j} = \tau_{i,j}(\mathbf{w}_i, y_{i-1,1}, \dots, y_{i-1,q_{i-1}}), \quad j = 1, \dots, q_i.$$

where $y_{i,j}$ is the output of neuron j at level i , and $\tau_{i,j}$ is a polynomial of degree no larger than α_i in the components of the weight vector \mathbf{w}_i , and no larger β_i in the components of the vectors $y_{i-1,j}$. At the final layer, there is a simple perceptron device following the polynomial activation function.

With this class of neural networks, it is clear that the output will equal one if and only if a polynomial inequality of the form

$$y_l(\mathbf{w}, \mathbf{x}) \geq 0,$$

is satisfied, where \mathbf{w} is the weight vector and $\mathbf{x} = (x_1 \dots x_N)$ is the input vector. Thus we can apply Theorem 1 with $s = 1$. The issue now is to determine the number of connected components B of the semi-algebraic set defined by $y_l(\mathbf{w}, \mathbf{x}) = \mathbf{y}$.

Now we are in a position to state the main result. To facilitate the statement, we introduce a bit of notation. Define

$$d_l = \alpha_l, d_{l-1} = \alpha_{l-1}\beta_l, \dots, d_i = \alpha_i \prod_{j=i+1}^l \beta_j, \quad i = 1, \dots, l-1.$$

Recall that k_i denotes the number of adjustable parameters at level i , and that k denotes the total number of adjustable parameters. Define the probability vectors

$$\mathbf{v} := (k_1/k \dots k_l/k), \quad \mathbf{u} := (d_1/d \dots, d_l/d),$$

and define the “binary” relative entropy $H(\mathbf{v}|\mathbf{u})$ as

$$H(\mathbf{v}|\mathbf{u}) := \sum_{i=1}^l v_i \lg \left(\frac{v_i}{u_i} \right) = \frac{1}{k} \sum_{i=1}^l k_i \lg \left(\frac{dk_i}{kd_i} \right).$$

Note that the above is the same as the conventional relative entropy of two probability vectors, except that we use base-2 logarithms instead of natural logarithms. Following standard convention, we take $0 \lg(0/0) = 0$.

Theorem 4 *With the above notation, we have*

$$B \leq 2^k k! \prod_{i=1}^l \frac{d_i^{k_i}}{k_i!}. \quad (4.1)$$

$$\leq \left(\frac{2d}{e^{7/8}} \right)^k 2^{-kH(\mathbf{v}|\mathbf{u})}. \quad (4.2)$$

where $d := \sum_{i=1}^l d_i$ and we assume $k_1, \dots, k_l \geq 2$ in the last inequality. Consequently, when $k_1, \dots, k_n \geq 2$, the VC-dimension of the neural network architecture is bounded above by

$$2k(\lg(4ed) - H(\mathbf{v}|\mathbf{u})).$$

Remark The above theorem shows that the reduction in the VC-dimension estimate over that of Theorem 2 is precisely $2k$ times the (binary) relative entropy of the two probability vectors v and u defined above. Thus if $k_i/k = d_i/d$ for all i , there will not be any reduction at all. In general, the fraction by which the older VC-dimension estimate is reduced is precisely the ratio $H(\mathbf{v}|\mathbf{u})/(\lg(4ed))$. Note also that the assumption that there are at least 2 adjustable parameters at each levels is a reasonably mild assumption. \diamond

Proof of Theorem 4: The proof depends on a careful book-keeping of the degree of $y_l(\mathbf{w}, \mathbf{x})$ with respect to the various components of \mathbf{w} . From the architecture of the neural network, it is clear that at the first level, each of the $y_{1,j}$ is a polynomial in the components

of \mathbf{w}_1 of degree no larger than α_1 . At the second level, each of the $y_{2,j}$ is a polynomial, whose monomials are of (combined) degree no larger than α_2 in the components of \mathbf{w}_2 , and of (combined) degree no larger than $\beta_2\alpha_1$ in the components of \mathbf{w}_1 . Thus, while each $y_{2,j}$ could have a total degree of $\alpha_2 + \beta_2\alpha_1$ in the components of \mathbf{w}_1 and \mathbf{w}_2 , the total degree of the monomial terms involving the components of \mathbf{w}_1 does not exceed $\beta_2\alpha_1$, while the total degree of the monomial terms involving the components of \mathbf{w}_2 does not exceed α_2 . A simple argument by induction then tells us that at the output layer (level l), the single output y_l is a polynomial whose monomials have total degree no larger than $d_l = \alpha_l$ in the components of \mathbf{w}_l , no larger than $d_{l-1} = \beta_l\alpha_{l-1}$ in the components of \mathbf{w}_{l-1} , and so on. With the d_i 's defined as above, the components of each \mathbf{w}_i appear with total degree no larger than d_i . Thus the total degree of y could be as large as d , but the monomial terms involving the components of \mathbf{w}_i have total degree no larger than d_i . So the set V defined in Theorem 3 satisfies the following containment:

$$V \subseteq \prod_{i=1}^l S_{d_i}^{k_i}.$$

Because of this containment, it follows that

$$\text{Vol}_n(V) \leq k! \prod_{i=1}^l \frac{d_i^{k_i}}{k_i!}.$$

Combining this with the bound (3) establishes the first estimate (4.1).

To prove the second estimate, we use Stirling's approximation. In particular, [Rud76, ex. 20, pg. 200] tells us that for all $t \in \{2, 3, 4, \dots\}$, we have

$$e^{7/8}(t/e)^t \sqrt{t} < t! < e(t/e)^t \sqrt{t}.$$

Consequently, we easily obtain

$$\frac{k!}{k_1! \cdots k_l!} < e^{1-\frac{7}{8}k} \frac{k^k}{k_1^{k_1} \cdots k_l^{k_l} \sqrt{k_1 \cdots k_l}}.$$

Dropping the square root term on the bottom can of course be done, and then an elementary calculation yields $2^k k! \prod_{i=1}^l \frac{d_i^{k_i}}{k_i!} \leq \left(\frac{2d}{e^{7/8}}\right)^k 2^{-kH(\mathbf{v}|\mathbf{u})}$, provided $k_1, \dots, k_n \geq 2$.

The VC-dimension estimate (4) now follows readily from Theorem 1. \square

5 Numerical Example

Consider a network with four inputs, five hidden-layer neurons at the first level and an output-layer neuron. As is common, let us suppose that $\alpha_i = 1$ for all i . This means

that all the adjustable parameters *enter linearly* into the corresponding activation function. Suppose $\beta_1 = 2, \beta_2 = 3$. This means that the hidden-layer neurons have quadratic activation functions, whereas the output-layer neuron has a cubic activation function. It remains to specify the integers k_1 and k_2 , representing the number of adjustable parameters. Let us assume that practically all of the monomial terms are present in each neural characteristic. Thus it is reasonable to assume $k_1 = 50, k_2 = 20$. Finally, $d_1 = 3, d + 2 = 1$. With these figures, one has

$$\mathbf{v} = (5/7, 2/7), \quad \mathbf{u} = (0.25, 0.75),$$

$$H(\mathbf{v}|\mathbf{u}) \approx 0.684033, \quad \lg(4ed) \approx 5.4427, \quad \frac{H(\mathbf{v}|\mathbf{u})}{\lg(4ed)} \approx 0.12567.$$

Thus, in this case, the improved bound is roughly 12.5% sharper.

6 Conclusions

References

- [GJ93] P. Goldberg and M. Jerrum, “Bounding the Vapnik-Chervonenkis dimension of concept classes parametrized by real numbers,” *Proc. 6th ACM Workshop on Computational Learning Theory*, pp. 361–369, 1993.
- [KM97] Marek Karpinski and Angus J. Macintyre, “Polynomial bounds for VC dimension of sigmoidal and general Pfaffian neural networks,” *J. Comp. Sys. Sci.*, 54, 169–176, 1997.
- [Mil64] John W. Milnor, “On the Betti numbers of real varieties,” *Proc. Amer. Math. Soc.*, 15, 275–280, 1964.
- [OP49] O. Oleinik and I. Petrovsky, “*On the Topology of Real Algebraic Hypersurfaces*,” *Izv. Akad. Nauk SSSR* 13, pp. 389–402, 1949.
- [Roj00] J. Maurice Rojas, “Some Speed-Ups and Speed Limits in Real Algebraic Geometry,” *Journal of Complexity*, FoCM 1999 special issue, vol. 16, no. 3 (sept. 2000), pp. 552–571.
- [Rud76] Rudin, Walter, *Principles of Mathematical Analysis*, 3rd edition, McGraw-Hill, 1976.

- [Tho65] René Thom “*Sur l’homologie des variétés algébriques réelles,*” In S. Cairns (Ed.), *Differential and Combinatorial Topology*, Princeton University Press, 1965.
- [Vap95] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, New York 1995.
- [MV97] M. Vidyasagar, *A Theory of Learning and Generalization*, Springer-Verlag, London, 1997.