

VC Dimension of Sigmoidal and General Pfaffian Neural Networks

Marek Karpinski¹ Angus Macintyre²

TR-95-065

November 1995

Abstract

We introduce a new method for proving explicit upper bounds on the VC Dimension of general functional basis networks, and prove as an application, for the first time, that the VC Dimension of analog neural networks with the sigmoidal activation function $\sigma(y) = 1/1 + e^{-y}$ is bounded by a quadratic polynomial $O((lm)^2)$ in both the number l of programmable parameters, and the number m of nodes. The proof method of this paper generalizes to much wider class of Pfaffian activation functions and formulas, and gives also for the first time polynomial bounds on their VC Dimension. We present also some other applications of our method.

¹Dept. of Computer Science, University of Bonn, 53117 Bonn. Research partially supported by the International Computer Science Institute, Berkeley, by the DFG Grant KA 673/4-1, and by the ESPRIT BR Grants 7097 and ECUS 030. Research partially done while visiting Dept. of Computer Science at Princeton University. Email: marek@cs.uni-bonn.de

²Mathematical Institute, University of Oxford, Oxford OX1 3LB. Research supported in part by a Senior Research Fellowship of the EPSRC. Email: ajm@maths.ox.ac.uk

0 Introduction

This paper studies the VC Dimension of general functional basis networks, and the resulting Boolean combinations of certain formulas. We develop a new method for proving explicit upper bounds for a wide class of analog neural networks with general Pfaffian activation functions.

The most commonly used activation function in various neural networks applications is the sigmoid $\sigma(y) = 1/(1 + e^{-y})$ (cf. [HKP91]). We refer to [AB92], [M93a], and [MS93] for all the necessary background on the computation by neural networks and the VC dimension (particularly, to the connection between their computational power, and the sample complexity).

In [MS93] the finiteness of VC Dimension of sigmoidal neural networks has been established for the first time using a deep result in model theory. It is perhaps worth nothing that slightly more general analytic increasing activation functions do not always have finite VC-dimension [S92].

In Maass's 1993 lecture notes [M93a] (see also [GJ93] and [MS93]), Open Problem 10 asks:

Is the VC Dimension of analog neural nets with the sigmoid activation function $\sigma(y) = 1/(1 + e^{-y})$ bounded by a polynomial in the number of programmable parameters?

In this paper we give an affirmative answer, with a polynomial bound in the number of programmable parameters. We believe that the bound can be improved to the one subquadratic in the number of programmable parameters and the number of nodes using a variant of our method. The result is a special case of much more general result about the VC Dimension of the classes defined by certain formulas. In contrast to [KM94], this paper does not use o-minimality and therefore can be applied to more general situations like the Pfaffian functions for which o-minimality is not yet even established(!).

In the case of boolean functions computed by sigmoidal neural networks (cf. [MSS91], [M93b]), our result entails, also for the first time, by a simple counting argument, the fact that not every boolean function can be computed by a single polynomial size sigmoidal or general Pfaffian neural network with an appropriate weight assignment.

We refer to [AB92], [GJ93], and [MS93] for all notions required for the VC Dimension of neural networks, and to [H76] for all notions of differential geometry.

The paper was inspired by the work of Goldberg and Jerrum [GJ93], who could deal with polynomial activation functions. A reference in [GJ93] to Warren's paper [W68] was of particular importance.

The paper is organized as follows. In Section 1, we introduce the necessary formalism for the describing formulas, as well as all preparatory algebraic and topological facts. Section 2 contains the Main Result, and Sections 3 and 4 the applications.

1 The setting

1.1 We shall consider a standard model of a *feedforward network architecture* A with the *activation* function σ (cf., e. g., [M93a], [MS93]) with k *inputs*, m *computational nodes*, and ℓ *weights* (the number of *programmable parameters*). We assume (for simplicity) that the output gate of A has range $\{0,1\}$. We associate with A an exponential formula $\Phi(\bar{v}, \tilde{y}) > 0$ for $\bar{v} \in \mathbb{R}^k$, and $\tilde{y} \in \mathbb{R}^\ell$, Φ being a composition of polynomials, and activation functions over the computation nodes of A . $\Phi(\bar{v}, \tilde{y}) > 0$ represents the function computed by A . Alternatively, and this is crucial in our paper, we describe the computation of A as a Boolean combination of atomic formulas of two forms $\tau(\bar{v}, \tilde{y}) = 0$ or $\tau(\bar{v}, \tilde{y}) > 0$ describing local computations of A at its computational nodes (for appropriate \bar{v} 's, and \tilde{y} 's). The *VC dimension* of the *network* A is the *VC dimension* of the class $\mathcal{C}_\Phi = \{\Phi_{\tilde{\beta}} : \tilde{\beta} \in \mathbb{R}^\ell\}$ for $\Phi_{\tilde{\beta}} = \{\bar{x} \in \mathbb{R}^k : \Phi(\bar{x}, \tilde{\beta}) > 0\}$ the partition of \mathbb{R}^k by A according to the weight assignment $\tilde{\beta}$. (The general reader is referred to [MS93] and [GJ93] for definitions and basic properties of Vapnik-Chervonenkis (VC) dimension. We say a set $S \subseteq \mathbb{R}^k$ is *shattered* by \mathcal{C}_Φ if $\{S \cap C : C \in \mathcal{C}_\Phi\} = P(S)$. The VC dimension of \mathcal{C}_Φ is the maximal size of any set S that can be shattered by \mathcal{C}_Φ , or ∞ if arbitrary large subsets may be shattered.)

We turn our attention now to the analysis of general formulas resulting from the local computation descriptions of A . The method of our analysis is by no means restricted to the network architectures only, and can be applied to a much larger class of formulas, which could be of independent interest.

1.2 We start now with some definitions and notations. Fix integers k, l and C^∞ (infinitely differentiable) functions τ_1, \dots, τ_s from \mathbb{R}^{k+l} to \mathbb{R} . Write τ_i as $\tau_i(v_i, \dots, v_k, y_1, \dots, y_\ell)$ (or $\tau_i(\bar{v}, \tilde{y})$).

Form a first-order language L with primitives $<$ (for order) and function symbols $\bar{\tau}_1, \dots, \bar{\tau}_s$, of arity $k + l$, corresponding to τ_1, \dots, τ_s . (We drop $-$ for readability.)

Let $\Phi(\bar{v}, \tilde{y})$ be a quantifier-free L -formula, so Φ is a Boolean combination of atomic formulas, which can be of two forms:

$$\tau(\bar{v}, \tilde{y}) > 0$$

or

$$\tau(\bar{v}, \tilde{y}) = 0,$$

where τ is an L -term. For this paper, we assume each τ to be one of τ_i 's.

For $\tilde{\beta} \in \mathbb{R}^\ell$, one defines

$$\Phi_{\tilde{\beta}} = \{\bar{v} \in \mathbb{R}^k : \mathbb{R} \models \Phi(\bar{v}, \tilde{\beta})\} \subseteq \mathbb{R}^k,$$

and the family

$$\mathcal{C}_\Phi = \{\Phi_{\tilde{\beta}} : \tilde{\beta} \in \mathbb{R}^\ell\}.$$

In this paper we give *good* explicit bounds on the VC-dimension of \mathcal{C}_Φ , under certain assumptions about the τ_i .

1.3 Assumptions on the τ_i . Let $\bar{\alpha}_1, \dots, \bar{\alpha}_V$ be elements of \mathbb{R}^k . Form the ($\leq sV$ many) C^∞ functions $\tau_i(\bar{\alpha}_j, \tilde{y})$ from \mathbb{R}^ℓ to \mathbb{R} . Choose $\Theta_1, \dots, \Theta_r$ ($r \leq \ell$) from among these, and let

$$F : \mathbb{R}^\ell \rightarrow \mathbb{R}^r$$

be defined by

$$F(\tilde{y}) = \langle \Theta_1(\tilde{y}), \dots, \Theta_r(\tilde{y}) \rangle.$$

By Sard's Theorem [M65], the set of nonregular values $\langle \epsilon_1, \dots, \epsilon_r \rangle$ of F in \mathbb{R}^r has Lebesgue measure 0. Recall that $\langle \epsilon_1, \dots, \epsilon_r \rangle$ is a regular value of F if either

a) $F^{-1}(\langle \epsilon_1, \dots, \epsilon_r \rangle) = \emptyset$,

or

b) $F^{-1}(\langle \epsilon_1, \dots, \epsilon_r \rangle)$ is an $(\ell - r)$ -dimensional C^∞ -submanifold of \mathbb{R}^ℓ .

This motivates the assumption we now impose on the τ_i .

Assumption: There is a bound B , independent of the $\bar{\alpha}_j$, r , and $\epsilon_1, \dots, \epsilon_r$ such that if $F^{-1}(\langle \epsilon_1, \dots, \epsilon_r \rangle)$ is an $(\ell - r)$ -dimensional C^∞ -submanifold of \mathbb{R}^ℓ then $F^{-1}(\langle \epsilon_1, \dots, \epsilon_r \rangle)$ has $\leq B$ connected components.

Fix such a B henceforth.

1.4 Examples. a) The τ_i are polynomials of degree $\leq d$ in \tilde{y} . Then B can be taken as $2 \cdot (2d)^\ell$ by a result of Milnor [M64].

b) Khovanski [K91, p. 91, Corollary 3] proved a basic result about exponential polynomials, namely:

Theorem 1. Suppose $l \geq m$. Let Q_i ($i \leq m$) be elements of $\mathbb{R}[y_1, \dots, y_\ell, e^{\Lambda_1}, \dots, e^{\Lambda_q}]$ where the Λ_i are linear functions of y_1, \dots, y_ℓ . Suppose the Q_i have degree d_i , and let $k = l - m$, and $S = \sum_{i=1}^m d_i + k + 1$. Suppose $(0, \dots, 0)$ is a regular value of (Q_1, \dots, Q_m) , with inverse image a manifold of dimension k . Then that manifold has no more than $2^{q(q-1)/2} \cdot d_1 \cdots d_m \cdot S^k [(k+1)S - k]^q$ connected components.

This gives, for 1.3, if $\tau_i(\bar{v}, \tilde{y})$ is polynomial of degree d in \bar{v}, \tilde{y} and q fixed subterms (independent of i) $\exp(g(\bar{v}, \tilde{y}))$, g linear, $B \leq 2^{q(q-1)/2} \cdot d^l \cdot S^l [lS]^{ql}$ where $S = dl + l + 1$. The q in Theorem 1 becomes ql now, because of the substitutions of $\leq l$ many $\bar{\alpha}_j$ for \tilde{v} .

So

$$\begin{aligned} \log B &\leq (ql)(ql - 1)/2 + l \log d + l \log S + ql \log (lS) \\ &\leq (ql)(ql - 1)/2 + l \log d + l(q + 1) \log S + ql \log l. \end{aligned}$$

c) If the τ_i are definable in an o-minimal expansion of the real field [KPS86], the existence of a B is guaranteed, but *good* bounds are not.

d) Examples for which o-minimality is unknown but where our method applies involved Pfaffian functions (cf. [K91, p. 91, Example 3]). We recall that a sequence of real functions F_1, \dots, F_q is a Pfaffian chain if all partial derivatives of every $F_i, 1 \leq i \leq q$, can be expressed as polynomials in the first i functions in the chain and the coordinate functions. Suppose the $\tau_i(\bar{v}, \tilde{y})$ are polynomials of degree $\leq d$ in the \bar{v}, \tilde{y} and in functions F_1, \dots, F_q which form a Pfaffian chain of length q where the polynomials are of degree $\leq D$. Let $r \leq l$ and $\Theta_1, \dots, \Theta_r$ as in 1.3, defining a manifold of dimension $l - r$. Then if $S = r(d - 1) + lD + 1$, we have

$$B \leq 2^{lq(lq-1)/2} \cdot d^r \cdot S^{l-r} [(l - r + 1)S - (l - r)]^{lq}$$

giving, independent of r a (crude) bound

$$B \leq 2^{lq(lq-1)/2} \cdot d^l \cdot (l(d + D))^l (l^2(d + D))^{lq}.$$

The bound in Theorem 1 corresponds to $D = 1$.

As for the exponential example Khovanski's q becomes in our case lq after the $\bar{\alpha}_j$ get substituted.

2 The Main Result

2.1 We shall prove:

Theorem 2. (Assumption as above).

$$VC\text{-Dimension}(\mathcal{C}_\Phi) \leq 2 \log B + (16 + 2 \log s)l.$$

(Note: In this paper \log is logarithm to base 2.)

2.2 Let $\bar{a}_1, \dots, \bar{a}_V$ be elements of \mathbb{R}^k such that $\{\bar{a}_1, \dots, \bar{a}_V\}$ is shattered by \mathcal{C}_Φ . For each subset E of $\{\bar{a}_1, \dots, \bar{a}_V\}$, pick \tilde{y}_E in \mathbb{R}^l such that $E = \{\bar{a}_j : \mathbb{R} \models \Phi(\bar{a}_j, \tilde{y}_E)\}$. Choose $\epsilon > 0$ such that if any $\tau_i(\bar{a}_j, \tilde{y}_E)$ ($1 \leq i \leq s, 1 \leq j \leq V, E \subseteq \{\bar{a}_1, \dots, \bar{a}_V\}$) is $\neq 0$, then $|\tau_i(\bar{a}_j, \tilde{y}_E)| > \epsilon$.

Note that for $\tilde{\gamma} \in \mathbb{R}^l$, the set $\{\bar{a}_j : \mathbb{R} \models \Phi(\bar{a}_j, \tilde{\gamma})\}$ depends only on the *signs* (+, -, or 0) taken at $\tilde{\gamma}$ by the functions $\tau_i(\bar{a}_j, \tilde{\gamma})$, ($1 \leq i \leq s, 1 \leq j \leq V$). [The sign of λ is + if $\lambda > 0$, - if $\lambda < 0$, and 0 if $\lambda = 0$].

Because of the \tilde{y}_E one has $\geq 2^V$ such sign series as $\tilde{\gamma}$ varies.

The \tilde{y}_E now show the following:

Lemma 3. If $0 < \epsilon_{ij} < \epsilon$ ($1 \leq i \leq s, 1 \leq j \leq V$) the complement in \mathbb{R}^l of the union of the sets $\{\tilde{y} : \tau_i(\bar{a}_j, \tilde{y}) = \epsilon_{ij}\} \cup \{\tilde{y} : \tau_i(\bar{a}_j, \tilde{y}) = -\epsilon_{ij}\}$ ($1 \leq i \leq s, 1 \leq j \leq V$) has at least 2^V connected components (V, \bar{a}_j, ϵ are fixed as above).

2.3 This can now be combined with Sard [S42], [M65], and a combinatorial idea of Warren [W68], to give Theorem 2.

We use the following cases of Sard's Theorem. We have a C^∞ map $F : \mathbb{R}^m \rightarrow \mathbb{R}^n$. A point p of \mathbb{R}^n is called a *regular value* of F if either $m \leq n$ and $F^{-1}(p)$ is a submanifold of \mathbb{R}^m of dimension $m - n$, or empty, or $m > n$ and $F^{-1}(p)$ is empty. Then the basic result is that the set of q in \mathbb{R}^n which are not regular values of F has Lebesgue measure 0.

(It is easily seen that the normal definition of regular value, in terms of $F^{-1}(p)$ containing no critical points, is equivalent to that given above.)

Now we apply Sard [S42]. Let $P = \{ \langle i, j \rangle : 1 \leq i \leq s, 1 \leq j \leq V \}$. For $\langle i, j \rangle \in P$, let $\tau_{i,j}(\tilde{y}) = \tau_i(\bar{a}_j, \tilde{y})$. For $A \subseteq P$, and $f \in \{1, -1\}^A$, let $F_{A,f}(\tilde{y}) = \langle \cdots, f(\langle i, j \rangle) \cdot \tau_{i,j}(\tilde{y}), \cdots \rangle_{\langle i,j \rangle \in A}$. So $F_{A,f}$ is a C^∞ map from \mathbb{R}^l to \mathbb{R}^A . For $\vec{\epsilon}$ in \mathbb{R}^P , let

$$Z(A, f)(\vec{\epsilon}) = \{ \tilde{y} : \text{for all } \langle i, j \rangle \in A, \tau_{i,j}(\tilde{y}) = f(\langle i, j \rangle) \tau_{i,j}(\tilde{y}) \}.$$

Finally let $I = [-1, 1]$, which has measure 2.

Lemma 4. Let Γ be the set of all $\vec{\epsilon}$ in I^P such that for all $A \subseteq P$ with $\text{card}(A) = j \leq l$, and all $f \in \{1, -1\}^A$, $Z(A, f)(\vec{\epsilon})$ is either empty, or a manifold of dimension $l - j$.

Then Γ has measure $2^{\text{card}(P)}$.

Proof. Look at the $\vec{\epsilon}$ for which the condition fails for some A, f . Let Π_A be the projection of I^P onto I^A . Then $\Pi_A(\vec{\epsilon})$ is not a regular value of $F_{A,f}$, so belongs to a subset of \mathbb{R}_A of measure 0. So the $\vec{\epsilon}$ in I^P for which the condition fails for A, f have measure 0. Since there are only finitely many A, f the result follows. \square

Now a slight refinement.

Lemma 5. Let Γ' be the subset of Γ consisting of all $\vec{\epsilon}$ such that if $\text{card}(A) > l$ and $f \in \{1, -1\}^A$, then $Z(A, f)(\vec{\epsilon})$ is empty. Then Γ' has measure $2^{\text{card}(P)}$.

Proof. Again, consider the $\vec{\epsilon}$ for which condition fails for a fixed A, f . As before, this set has measure 0. Since there are only finitely many A, f , the result follows. \square

We now take up the notations of Lemma 3.

The $\vec{\epsilon}_{ij}$ in I^P with $0 < \epsilon_{ij} < \epsilon$ form a set of measure $\epsilon^{\text{card}(P)}$. (Of course, $\text{card}(P) = sV$). Combining this with Lemma 5, we get that Γ' intersected with the above has measure $\epsilon^{\text{card}(P)}$, and so in particular is nonempty.

Note finally, before we approach Theorem 2 via a theorem of Warren, that for $\vec{\epsilon}$ in Γ' , if $A_1 \subseteq A_2$ and $f_1 \subseteq f_2$ then $Z(A_2, f_2)(\vec{\epsilon})$ is a *submanifold* of $Z(A_1, f_1)(\vec{\epsilon})$.

Warren [W68] proved:

Theorem 6. Let \mathcal{M} be a connected topological n -manifold, and let $\mathcal{M}_1, \dots, \mathcal{M}_n$ be connected $(n - 1)$ -manifolds which are submanifolds of \mathcal{M} so that

- (1) The \mathcal{M}_i are closed in \mathcal{M} ;
- (2) The intersection of any given j of the \mathcal{M}_i , $1 \leq j \leq n$ is either empty, or is an $(n - j)$ -submanifold of the intersection of any $(j - 1)$ of the \mathcal{M}_i ;
- (3) Any intersection of more than n of the \mathcal{M}_i is nonempty.

Let b_j ($0 \leq j \leq n$) be the number of connected components among all intersections of any j_n of the \mathcal{M}_i .

Then $\mathcal{M} - \bigcup_{i=1}^n \mathcal{M}_i$ has no more than $\sum_{j=0}^n b_j$ connected components.

Proof. See [W68, Theorem 1].

We want to apply this by fixing \vec{e} in Γ' , taking $\mathcal{M} = \mathbb{R}^l$, and the \mathcal{M}_i as the zerosets of the $\tau_{i,j}(\tilde{y}) \pm \epsilon_{ij}$. All that is missing is that we did not guarantee that these zerosets are connected. But if we rather take the \mathcal{M}_i as the connected components of the zerosets, the hypotheses of Theorem 6 are satisfied. Indeed, Warren's result clearly remains true if the condition on connectedness of the \mathcal{M}_i is dropped.

Going back to Lemma 3, we seek to bound 2^V by the number of connected components of the complement in \mathbb{R}^l of the union of the sets $\{\tilde{y} : \tau_{ij}(\tilde{y}) = \epsilon_{ij}\} \cup \{\tilde{y} : \tau_{ij}(\tilde{y}) = -\epsilon_{ij}\}$, where each ϵ_{ij} is between 0 and ϵ , and \vec{e} is in Γ' . To apply Warren, we have to bound the b_j , for $0 \leq j \leq l$. Of course $b_0 = 1$.

Now $n = 2sV$. Let $1 \leq j \leq l$. There are n many zero sets, but of course any intersection $\{\tilde{y} : \tau_{ij}(\tilde{y}) = \epsilon_{ij}\} \cap \{\tilde{y} : \tau_{ij}(\tilde{y}) = -\epsilon_{ij}\} = \emptyset$. Any intersection of no more than j of the zerosets has $\leq B$ connected components (original assumption). So by these two remarks

$$b_j \leq 2^j \cdot \binom{sV}{j} \cdot B,$$

giving

$$\begin{aligned} \sum_{j=0}^l b_j &\leq B \cdot \sum_{j=0}^l 2^j \cdot \binom{sV}{j} \\ &\leq B \cdot \left(\frac{2sVe}{l}\right)^l \end{aligned}$$

by [W68].

So now we have

$$2^V \leq B \cdot \left(\frac{2sVe}{l}\right)^l$$

Conclusion of Proof of Theorem 2.

Case 1. $V \leq 4s\epsilon l$

Then

$$2^V \leq B(8s^2\epsilon^2)^l \leq B(4s\epsilon)^{2l},$$

so

$$V \leq \log B + 2l \log(4s\epsilon) \leq \log B + 10l + 2l \log s$$

Case 2. $V > 4sel$

Then

$$2^V \leq B \left(\frac{V}{l} \right)^{2l},$$

so

$$2^{\frac{V}{l}} \leq B^{\frac{1}{l}} \left(\frac{V}{l} \right)^2.$$

Now $2^{\frac{V}{2l}} > \left(\frac{V}{l} \right)^2$ if $V > 16l$,

so either $2^{\frac{V}{2l}} < B^{\frac{1}{l}}$, or $V \leq 16l$.

So either $2^V < B^2$, or $V \leq 16l$.

So either $V < 2 \log B$, or $V \leq 16l$. \square

3 Applications

3.1 If we now work with polynomials, and Milnor's bound for B , we get the results from [GJ93].

3.2 An example involving exponentiation. Fix q and linear functions $\Lambda_1, \dots, \Lambda_q$ of \bar{v}, \tilde{y} . Let $\tau_i(\bar{v}, \tilde{y})$, $1 \leq i \leq s$, be polynomials, of total degree d_i , in \bar{v}, \tilde{y} and the e^{Λ_i} 's.

We showed after Theorem 1 (1.4) that

$$\log B \leq (ql)(ql - 1)/2 + l \log d + l(q + 1) \log S + ql \log l,$$

where

$$S = dl + l + 1 \leq (d + 1)(l + 1)$$

So

$$\begin{aligned} VC - \text{Dim}(\mathcal{C}_\Phi) & \\ & \leq (ql)(ql - 1) + 2l \log d + 2l(q + 1) \log S \\ & \quad + 2ql \log l + (16 + 2 \log s)l \\ & \leq (ql)(ql - 1) + 2l \log d + 2l(q + 1) \log(d + 1) \\ & \quad + 2l(q + 1) \log(l + 1) + 2ql \log l + (16 + 2 \log s)l. \end{aligned}$$

So

$$\begin{aligned} VC - \text{Dim}(\mathcal{C}_\Phi) & \\ & \leq (ql)(ql - 1) + 4l(q + 1) \log(l + 1) + 2l(q + 2) \log(d + 1) + (16 + 2 \log s)l. \end{aligned}$$

3.3 Application to sparse formulas. Since Khovanski's [K91] one has known how to use Finiteness Theorems about exponentiation to give uniform estimates in problems involving families of polynomials where there is an absolute bound to the number of nonzero coefficients occurring, but none on the degrees involved. So this is all we now assume about the $\tau_i(\bar{v}, \tilde{y})$.

The strategy is to break the \tilde{y} -space \mathbb{R}^ℓ into 3^ℓ pieces according to $y_j < 0$, $y_j = 0$, $y_j > 0$.

Having chosen for each j one such sign, one changes to variables y'_j with $y'_j = \log(-y_j)$ if $y_j < 0$, $y'_j = y_j$ if $y_j = 0$, and $y'_j = \log(y_j)$ if $y_j > 0$. Then $\tau_i(\bar{v}, \tilde{y})$ transforms to a function *linear* in no more than q_i exponentials of linear functions of the \tilde{y}' , where q_i is the number of nonzero coefficients of τ_i . In particular any $\tau_i(\bar{a}_j, \tilde{y})$ will satisfy the hypotheses of Khovanski's Theorem 1, with $d_i = 1$.

So we can apply 3.2 3^l times. After taking $\log s$ we get for $VC - Dim(\mathcal{C}_\Phi)$ the bound

$$(ql)(ql - 1) + 2l(q + 1) + 2l(q + 1) \log(l + 1) \\ + 2ql \log l + (16 + 2 \log s)l + l \log 3.$$

4 Application to Sigmoidal Neural Networks

4.1 Let us recall again [MS92] the definition of a sigmoidal network architecture A . The data involves:

- a) A directed acyclic graph G , labelled by variables and polynomials as explained below;
- b) an integer ℓ , the dimension of the space of *weights* (the number of *programmable* parameters), and the weight variables y_1, \dots, y_ℓ ;
- c) if there are k input nodes (i.e. nodes of in-degree 0) these are labelled by variables v_1, \dots, v_k ;
- d) there is exactly one output node (i.e. a node of out-degree zero);
- e) those nodes which are not input nodes are called computation nodes, and the m^{th} such N_m is labelled by a variable z_m , and a polynomial

$$P_{N_m}(v_{t_1}, \dots, v_{t_\rho}, z_{u_1}, \dots, z_{u_\gamma}, y_{\lambda_1}, \dots, y_{\lambda_\delta})$$

where the y 's are a subset of the weight variables, the v 's correspond to the input nodes immediately below N_m (i.e. connected to N_m) and the z 's correspond to the computation nodes immediately below N_m .

One now fixes an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, in our case the function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \ .$$

Then A computes a function $\beta_A : \mathbb{R}^{k+\ell} \rightarrow \mathbb{R}$, as described recursively below:

- a) If N is an input node, with associated variable v_i , $f_N(\bar{v}, \tilde{y}) = v_i$,

b) If N is a computation node with variable z_m

$$f_N(\bar{v}, \tilde{y}) = P_N(v_{t_1}, \dots, v_{t_\rho}, \sigma(f_{N_1}(\bar{v}, \tilde{y})), \dots, \sigma(f_{N_\gamma}(\bar{v}, \tilde{y})), y_{\lambda_1}, \dots, y_{\lambda_s}) \quad (\#)$$

where N_i corresponds to z_{u_i} , $1 \leq i \leq \gamma$.

Then β_A is f_{N_w} , where N_w is the output node.

Now, if we work in a language with $+$, $-$, \cdot , 0 , 1 , $<$ and a symbol σ for the activation function, then $f_{N_w}(\bar{v}, \tilde{y})$ is given by a term $\tau(\bar{v}, \tilde{y})$, by transcribing naively the above recursion. Let $\Phi(\bar{v}, \tilde{y})$ be

$$\tau(\bar{v}, \tilde{y}) > 0 \quad .$$

Then (by definition) the VC -dimension of A is the VC -dimension of \mathcal{C}_Φ . By [L92] (which appeals to [W94]) this dimension is finite, since σ is definable in $+$, $-$, \cdot , 0 , 1 , $<$, e^x .

We now apply our method to get a good polynomial bound for $VC - \dim(A)$. So we need to know a bound on the number of connected components of a manifold of dimension $\ell - j$ defined by the conditions

$$\tau(\bar{\alpha}_i, \tilde{y}) = \varepsilon_i \quad 1 \leq i \leq j \quad (\leq \ell) \quad .$$

We are aware of several approaches to this computation, and may in future look more closely at the relative merits of various methods. For now we appeal directly to the Khovanski estimates previously used, but now applied in a high-dimensional space.

For each i with $1 \leq i \leq j$, and each computation node N we add variables $Z_{N,i}$ and $\hat{Z}_{N,i}$. Among these are the output variables $Z_{w,i}$ each i . Finally, we add input variables

$$v_{c,i} \quad \text{for} \quad c \leq k, \quad i \leq j.$$

Now consider the system of equations

$$Z_{N,i} = P_N(v_{t_1,i}, \dots, v_{t_\rho,i}, \hat{Z}_{N_1,i}, \dots, \hat{Z}_{N_\gamma,i}, y_{\lambda_1}, \dots, y_{\lambda_s})$$

$$1 = \hat{Z}_{N,i}(1 + e^{-Z_{N,i}})$$

as N ranges over computation nodes, and $1 \leq i \leq j$. To see the meaning, refer to (#).

Write the system as

$$S(\bar{v}_1, \dots, \bar{v}_j, z_{w,1}, \dots, z_{w,j}, \tilde{y}, \tilde{w})$$

where $\bar{v}_c = (v_{c,1}, \dots, v_{c,j})$ and \tilde{w} denotes all the remaining variables.

The essential points are:

- (1) $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w}) \Rightarrow \tau(\bar{\alpha}_i, \tilde{y}) = \varepsilon_i \quad 1 \leq i \leq j$;
(2) If $\tau(\bar{\alpha}_i, \tilde{y}) = \varepsilon_i$ for all $1 \leq i \leq j$, then there are unique \tilde{w} such that $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w})$;
(3) The set in \mathbb{R}^ℓ defined by the conditions $\tau(\bar{\alpha}_i, \tilde{y}) = \varepsilon_i, \quad 1 \leq i \leq j$ is homeomorphic to that in (\tilde{y}, \tilde{w}) space defined by $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w})$, so either both or neither are manifolds.

So now we can use the Khovanski estimates on S , assuming $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w})$ defines a manifold of dimension $l - j$. Note that there are $l + (2m + 1) \cdot j$ variables among (\tilde{y}, \tilde{w}) , if m is the number of nonoutput computation nodes of A . $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w})$ is defined by $2mj$ equations, and of course $l + (2m + 1)j - 2mj = l - j$.

Let d be a bound for the degree of all P_N .

Then, by Khovanski, $S(\bar{\alpha}_1, \dots, \bar{\alpha}_j, \varepsilon_1, \dots, \varepsilon_j, \tilde{y}, \tilde{w})$ defines a set with no more than $2^{((m+1)j)((m+1)j-1)/2} d^{2mj} (2mjd + (l-j) + 1)^{l-j} ((l-j+1)2mjd - (l-j))^{(m+1)j}$ connected components.

So this gives us a bound B for the τ -problem, namely:

$$B \leq 2^{nl(nl-1)/2} \cdot d^{2nl} \cdot (l \cdot (2nd + 1))^l \cdot (2nl^2 d)^{nl}$$

where $n = m + 1 =$ number of computation nodes of A . So

$$\log B \leq nl(nl-1)/2 + 2nl \log d + l \log l + l \log(2nd+1) + 2nl \log l + nl \log(2nd) = \beta(A),$$

say.

Now, applying Theorem 2, we get:

Theorem 7. The VC-Dimension of A is bounded above by

$$2\beta(A) + 16l.$$

The term $(nl)(nl-1)/2$ is obviously the dominant term, if d is small. Since in the general case l could majorize n , one can argue that our bound is of degree 4 as a function of l only.

4.2 Generalizations: The estimation above with a dominant term $(m\ell)^2$ does not depend essentially on the type of the activation function used. An alternative approach to the above result works directly with the function $f_{N_w}(\bar{v}, \tilde{y})$, and uses the fact that f_{N_w} is a *Pfaffian* function. For the fundamental work on Pfaffian functions one should consult [K91].

$\sigma(x)$ is Pfaffian, since $\sigma'(x) = \sigma(x) - (\sigma(x))^2$.

Clearly $f_N(\bar{v}, \tilde{y})$ is Pfaffian, for N an input node, for $f_N(\bar{v}, \tilde{y}) = v_i$, where v_i is the input variable corresponding to N . Using (#) we have, for a computation node N ,

$$\begin{aligned}
\frac{\partial}{\partial y_j} f_N(\bar{v}, \tilde{y}) &= \frac{\partial P_N}{\partial Z_{u_1}} \cdot \frac{\partial f_{N_1}}{\partial y_j} \cdot (\sigma(f_{N_1}) - \sigma(f_{N_1})^2) \\
&+ \frac{\partial P_N}{\partial Z_{u_2}} \cdot \frac{\partial f_{N_2}}{\partial y_j} \cdot (\sigma(f_{N_2}) - \sigma(f_{N_2})^2) \\
&+ \\
&\vdots \\
&+ \\
&+ \frac{\partial P_N}{\partial Z_{u_r}} \cdot \frac{\partial f_{N_r}}{\partial y_j} \cdot (\sigma(f_{N_r}) - \sigma(f_{N_r})^2) \\
&+ \sum \frac{\partial P_N}{\partial y_{\lambda_r}} \cdot \frac{\partial y_{\lambda_r}}{\partial y_j}
\end{aligned}$$

and

$$\frac{\partial}{\partial y_j} \sigma(f_N(\bar{v}, \tilde{y})) = \frac{\partial}{\partial y_j} f_N(\bar{v}, \tilde{y}) \cdot (\sigma(f_N) - \sigma(f_N)^2).$$

From this one sees that if $\alpha_1, \dots, \alpha_r$ ($r \leq l$) are arbitrary values of \bar{v} then the collection of all $f_N(\bar{\alpha}_i, \tilde{y})$ and $\sigma(f_N(\bar{\alpha}_i, \tilde{y}))$, for $r \leq l$ and N an input or computation node, from a Pfaffian chain of length $2ml$, in which all polynomials have degree $\leq d + 2$. Finally, let $\Theta_i(\tilde{y})$ be $f_{N_w}(\bar{\alpha}_i, \tilde{y})$, a polynomial of degree $\leq d$ in the variables and the elements of the chain. Our task was to bound the number of connected components of

$$\{\tilde{y} : \Theta(\tilde{y}) = \varepsilon_1, \dots, \Phi_r(\tilde{y}) = \varepsilon_r\}$$

under the assumption this is an $(l - r)$ -submanifold of \mathbb{R}^ℓ . We can apply [K91][p.91, Example 3], described in 1.3.

So we get, in the present case,

$$B \leq 2^{2ml(2ml-1)/2} \cdot (d + 2)^l \cdot S^{l-1} (lS)^{2ml},$$

where

$$S \leq (d + 3)(l + 1).$$

This is slightly inferior to the bound given in Theorem 7. However, the method used here clearly generalizes to give a huge variety of examples in which, as in 3.2 or Theorem 7, we get a dominant term quadratic in ql . (In the above $q = 2ml$).

In particular, the analogue of Theorem 3.2 is:

Let $\tau_i(\bar{v}, \tilde{y})$, $1 \leq i \leq s$, be polynomials of degree $\leq d$ in the \bar{v} , \tilde{y} and functions f_1, \dots, f_q in a Pfaffian chain of length q and degree $\leq D$.

Then $VC - Dim(\mathcal{C}_\Phi)$

$$\begin{aligned} &\leq 2(ql)(ql - 1) + 2l \log d + 2l \log(ld + lD + 1) \\ &\quad 2ql \log l + 2ql \log(ld + lD + 1) + l(16 + 2 \log S). \end{aligned}$$

As for Theorem 7, it generalizes to architectures with Pfaffian activation functions. The only difference is that a q and D appear. suppose that the activation functions of A are all members of a Pfaffian chain of length q and degree D . Then the argument outlined earlier for the sigmoid case gives

$$\begin{aligned} B &\leq 2^{lmq(lmq-1)/2} \cdot d^l \cdot l^l (d + D)^l \cdot (l^2(d + D))^l m^q \\ &\leq 2^{lmq(lmq-1)/2} \cdot d^l \cdot l^{l+2lmq} \cdot (d + D)^{l+2lmq} \end{aligned}$$

so

$$\log B \leq lmq(lmq - 1)/2 + l \log d + (l + 2lmq) \log l + (l + 2lmq) \log(d + D),$$

given the VC bound

$$\begin{aligned} &\leq lmq(lmq - 1) + (2 \log d + 16 + 2 \log s)l \\ &\quad + 2(l + 2lmq) \log l + 2(l + 2lmq) \log(d + D). \end{aligned}$$

Thus there is a quadratic effect from q , but only a logarithmic one from D .

4.3 Arctangent. A special case is worth recording. Take arctangent as the activation function of a network architecture. The Pfaffian chain is $\frac{1}{1+x^2}$, $\arctan x$, so $q = 2$, and one readily verifies $D = 2$. so one has for arctangent activation the dominant term $4lm$, rather than lm for the sigmoid.

4.4 Sparse Networks. We maintain the notations of 4.1., but now we consider families of A 's, based on same graph and σ , but where the P_N can vary, subject to the restriction that none of them have more than Δ many nonzero coefficients. Combining the ideas of 3.3. and 4.1. we easily get for $\log B$ a bound with dominant term quadratic in $ln\Delta$, and this is of course dominant in the VC-dimension bound for the A 's in the family.

4.5 Haussler's Pseudodimension. We refer to [MS93] for the definition of the pseudo-dimension of an architecture. Since the pseudo-dimension of an architecture A is bounded by the VC-Dimension of a new architecture A' (see [MS93]) got directly from A , we get polynomial bounds for the pseudo-dimension. This answers affirmatively the second part of Problem 10 in [M93a].

4.6 Boolean Functions. We are interested now in computation of boolean functions $f : \{0, 1\}^k \rightarrow \{0, 1\}$ by neural networks (cf. [MSS91], [M93b]). It is known that applying some single non-boolean activation functions enhances, sometimes dramatically, the computational power of a neural network (cf. [MSS91]) even if restricted

to the boolean functions. However it has been open for sometime now how much this increase in computational power of a neural network could be. The fundamental inability to answer to this problem was caused by the lack of a method bounding the amount of information that can be encoded in the weights of a neural network. Particularly, no known methods were sufficient even to show that there always exists a boolean function $f : \{0, 1\}^k \rightarrow \{0, 1\}$ which *cannot* be computed by single constant depth, polynomial size (number of nodes and programmable parameters) neural network with sigmoidal activation function with an appropriate weight assignment. Main results of this paper entail a solution to this problem. In fact the polynomial bounds on the VC Dimension entail that no subexponential size $2^{o(k)}$ sigmoidal or general Pfaffian neural network can compute all boolean function $f : \{0, 1\}^k \rightarrow \{0, 1\}$ under appropriate weight assignments. Let A be a sigmoidal or general Pfaffian neural network with m nodes and ℓ programmable parameters. Denote by \mathbb{B}_A the set of *all boolean* functions computed by A under an appropriate weight assignment, and by d the VC Dimension of A .

Observe that also the VC Dimension of A restricted to the boolean functions is bounded by d . We have $\ln(|\mathbb{B}_A|) \leq O(kd)$ (cf., e.g., [AB92]). Our $O((\ell m)^2)$ upper bounds on the VC Dimension d of A entail now the following formula for the number $|\mathbb{B}_A|$ of different boolean functions computed by A : $|\mathbb{B}_A| \leq 2^{O(k\ell^2 m^2)}$.

4.7 Multivariate activation. There is also more remarkable further generalization. There is an obvious way to consider network architectures with multivariate activation functions. If these are Pfaffian, we still get a quadratic dominant term. We will elaborate this in a future publication.

5 Optimality of Khovanski's $2^{q(q-1)/2}$ bound ?

We strongly suspect that this bound can be lowered to the order q^q ($\sim 2^{q \log q}$). Obviously this would improve our upper bounds on the VC Dimension. The best lower bound on the VC Dimension of neural networks is $\Omega(l \log l)$ (cf. [M93a], [M94]) for the threshold, and $\Omega(l^2)$ (cf. [KS95]) for piecewise polynomial and sigmoidal activation functions. There is still a large gap between $\Omega(l^2)$ lower bound and our $O(l^4)$ upper bound for sigmoidal and Pfaffian activation functions. The current bound on B in our paper comes because of Khovanski's technique of removing one variable at a time (cf. [K91, p.13]). We are looking closely at a method for getting to a kind of *Bezout's estimate* in one step, removing all variables simultaneously. \square

Acknowledgement.

We thank Gregory Cherlin, Mark Jerrum and Eduardo Sontag for stimulating remarks and discussions. In particular, the current organization of the paper was suggested by Eduardo Sontag.

References

- [AB92] M. Anthony, N. Biggs, *Computational Learning Theory: An Introduction*, Cambridge University Press, 1992.
- [AS93] M. Anthony, J. Shawe-Taylor, A Result of Vapnik with Applications, *Discrete Applied Math.* **47** (1993), pp. 207–217.
- [BT90] A. Borodin, P. Tiwari, On the Decidability of Sparse Univariate Polynomial Interpolation, *Proc. 22nd ACM STOC* (1990), pp. 535–545.
- [D92] L. van den Dries, *Tame Topology and 0-minimal Structures*, preprint, University of Illinois, Urbana, 1992; to appear as a book.
- [DMM94] L. van den Dries, A. Macintyre and D. Marker, The Elementary Theory of Restricted Analytic Fields with Exponentiation, *Annals of Mathematics* **140** (1994), pp 183-205.
- [GJ93] P. Goldberg and M. Jerrum, Bounding the Vapnik Chervonenkis Dimension of Concept Classes Parametrized by Real Numbers. *Machine Learning*, 1994 (to appear). A preliminary version appeared in *Proc. 6th ACM Workshop on Computational Learning Theory*, pp. 361–369, 1993.
- [H12] G. H. Hardy, Properties of Logarithmic-Exponential Functions, *Proc. London Math. Soc.* **10** (1912), pp. 54–90.
- [H92] D. Haussler, Decision Theoretic Generalizations of the PAC Model for Neural Net and other Learning Applications, *Information and Computation* **100**, (1992), pp. 78–150.
- [HKP91] J. Hertz, A. Krogh and R. G. Palmer, *Introduction to the Theory of Neural Computation*, Addison-Wesley, 1991.
- [H76] M. W. Hirsch, *Differential Topology*, Springer-Verlag, 1976.
- [KM94] M. Karpinski and A. Macintyre, Polynomial Bounds for VC Dimension of Sigmoidal Neural Networks, *Proc. 27th ACM STOC* (1995), pp.200-208.
- [KM95] M. Karpinski and A. Macintyre, Bounding VC Dimension for Neural Networks: Progress and Prospects (Invited Lecture), *Proc. EuroCOLT'95, Lecture Notes in Artificial Intelligence Vol.904*, Springer-Verlag, 1995, pp. 337-341.
- [KW93] M. Karpinski and T. Werther, VC Dimension and Uniform Learnability of Sparse Polynomials and Rational Functions, *SIAM J. Computing* **22** (1993), pp 1276–1285.

- [K91] A. G. Khovanski, *Fewnomials*, American Mathematical Society, Providence, R.I., 1991.
- [KPS86] J. Knight, A. Pillay and C. Steinhorn, *Definable Sets and Ordered Structures II*, *Trans. American Mathematical Society* **295** (1986), pp.593-605.
- [KS95] P. Koiran and E.D. Sontag, *Neural Networks with Quadratic VC Dimension to appear in Advances in Neural Information Processing Systems (NIPS '95)*, 1995.
- [L92] M. C. Laskowski, *Vapnik-Chervonenkis Classes of Definable Sets*, *J.London Math. Society* **45** (1992), pp 377-384.
- [M93a] W. Maass, *Perspectives of Current Research about the Complexity of Learning on Neural Nets*, in: *Theoretical Advances in Neural Computation and Learning*, V. P. Roychowdhury, K. Y. Siu, A. Orlitsky (Editors), Kluwer Academic Publishers, 1994, pp. 295-336.
- [M93b] W. Maass, *Bounds for the Computational Power and Learning Complexity of Analog Neural Nets*, *Proc. 25th ACM STOC* (1993), pp. 335-344.
- [M94] W. Maass, *Neural Nets with Superlinear VC-Dimension*, *Proc. of the International Conference on Artificial Neural Networks 1994 (ICANN '94)*, Springer (Berlin 1994), pp. 581-584; journal version appeared in *Neural Computation* **6** (1994), pp. 875-882.
- [MSS91] W. Maass, G. Schnitger and E. D. Sontag, *On the Computational Power of Sigmoidal versus Boolean Threshold Circuits*, *Proc. 32nd IEEE FOCS* (1991), pp. 767-776.
- [MS93] A. J. Macintyre and E. D. Sontag, *Finiteness results for Sigmoidal Neural Networks*, *Proc. 25th ACM STOC* (1993), pp.325-334.
- [M64] J. Milnor, *On the Betti Numbers of Real Varieties*, *Proc. of the American Mathematical Society* **15** (1964), pp 275-280.
- [M65] J. Milnor, *Topology from the Differentiable Viewpoint*, Univ.Press, Virginia, 1965.
- [S42] A. Sard, *The Measure of the Critical Points of Differentiable Maps*, *Bull. Amer. Math. Soc.* **48** (1942), pp. 883-890.
- [S-T94] J. Shawe-Taylor, *Sample Sizes for Sigmoidal Neural Networks*, Preprint, University of London, 1994, to appear in *Proc. ACM COLT*, 1995.
- [S92] E. D. Sontag, *Feedforward Nets for Interpolation and Classification*, *J. Comp. Syst. Sci.* **45** (1992), pp. 20-48.

- [TV94] G. Turan and F. Vatan, On the Computation of Boolean Functions by Analog Circuits of Bounded Fan-in, Proc. 35th IEEE FOCS (1994), pp. 553–564.
- [W68] H. E. Warren, Lower Bounds for Approximation by Non-linear Manifolds, Trans. of the AMS **133** (1968), pp. 167–178.
- [W94] A. J. Wilkie, Model Completeness Results of Restricted Pfaffian Functions and the Exponential Function; to appear in Journal of the AMS, 1994.