# ISTITUTO NAZIONALE DI FISICA NUCLEARE

Sezione di Trieste

M. Budinich

# ON LINEAR SEPARABILITY OF RANDOM SUBSETS OF HYPERCUBE VERTICES

# On Linear Separability of Random Subsets of Hypercube Vertices.

Marco Budinich

Dipartimento di Fisica dell'Università di Trieste and INFN Trieste,
Via Valerio 2, I-34127 Trieste, Italy.

PACS 87.10 - General, theoretical and mathematical biophysics
PACS 05.90 - Other topics in statistical physics and thermodynamics

## Abstract

The classical Cover results on linear separability of points in $R^d$ are a milestone in neural network theory. Nevertheless they are not valid for digital input networks because in this case the points are not in general position being vertices of a d dimensional hypercube. I show here that for large d all Cover findings can be extended to this case. I also show that for $n < O\left((d+1)^{\frac{3}{2}}\right)$ the number of linear separations of n random hypercube vertices tends to that of n points in general position.

Feed-forward neural networks have frequently solicited studies on geometrical properties of their input space.

The values of the d input neurons can be thought as coordinates of d dimensional space $R^d$ and then the set of all possible inputs is a subset of $R^d$ (the pattern space). In the frequent case of digital inputs (0,1 or ±1) the pattern space shrinks to the set of the vertices of the d dimensional hypercube $Q^d \subset R^d$.

The seminal Cover paper[1] [1] showed many interesting properties for sets of n points in general position in $R^d$. The points are in general position if any k-tuple $(k \le d+1)$ of them is linearly independent.

Cover showed that the probability $P(n,d)$ that n random points in general position in $R^d$ are linearly separable is[2]

$$(1) \qquad P(n,d) = \frac{\text{number of linear separations}}{\text{total number of separations}} = \frac{2\sum_{k=0}^{d}\binom{n-1}{k}}{2^n}$$

---

[1] For some more recent works with a similar approach see e. g. [2] and [3].

[2] This is the probability that exists an hyperplane separating a random partition of the n points in two sets. The n points are supposed to be in general position in $R^d$. For more precise definitions see [1].

From this formula Cover derive all of his interesting results directly applicable to one layer feed-forward neural networks (perceptrons). The more important are (all for $d \to \infty$):

- the probability of linear separability of n random points falls to 0 when $n > 2(d+1)$

  $$P(n,d) \to \Phi(-x) \qquad \text{for } d \to \infty \text{ and } n = 2(d+1) + x\sqrt{2(d+1)}$$

  where $\Phi(-x)$ is the cumulative normal distribution;

- the perceptron "capacity" is $2(d+1)$ i.e. two random patterns per weight;
- the probability of "non ambiguous generalization" $\to 0$ if $n < 2(d+1)$

  where n is the number of patterns already "learned" by the network.

If the pattern space is the set of vertices of $Q^d$ (a very common situation in neural networks) (1) and all subsequent results are no more valid. This happens because the points are usually not in general position[1].

In what follows I show that for the identically defined probability $H(n,d)$ that n random vertices of $Q^d$ are linearly separable holds the relation

(2) $\qquad\qquad\qquad H(n,d) \to P(n,d) \qquad \text{when } d \to \infty$

that extends (1) and related results to subsets of vertices of $Q^d$ when $d \to \infty$ (the demonstration is similar to that used by Füredi in [4]).

Let $C_{gp}(n,d)$ and $C(n,d)$ be the number of linear separations of a set $\Pi_n$ of n points in $R^d$ respectively with and without the hypothesis of general position. Füredi [4] obtains the following bounds from the geometrical theorem of Winder [5]

(3) $\qquad\qquad C_{gp}(n,d) - \sum_{k=2}^{d+1} a_k(\Pi_n,d) \leq C(n,d) \leq C_{gp}(n,d)$

where $a_k(\Pi_n,d)$ is the number of linear dependent k-tuples of points of the set $\Pi_n$.

To pass from (3) to the probabilities of (1) and (2) we have to average the quantities $C(n,d)$ and $a_k(\Pi_n,d)$ over all the possible $\Pi_n$ and then to divide by the number of possible partitions i.e. $2^n$. With the hypothesis that the n points are vertices of $Q^d$ we have $\binom{2^d}{n}$ possible choices for the set $\Pi_n$ so (3) gives

---

[1] The d dimensional hypercube is a highly symmetric figure where for example no 2d points in general position exist or where all points with a given number of 1's lay on just one hyperplane.

2

$$(4) \qquad P(n,d) - \frac{2\sum\limits_{\Pi_n} \sum\limits_{k=2}^{d+1} a_k(\Pi_n,d)}{2^n \binom{2^d}{n}} \le H(n,d) \le P(n,d) .$$

The quantity

$$\frac{\sum\limits_{\Pi_n} a_k(\Pi_n,d)}{\binom{2^d}{n}\binom{n}{k}}$$

is, by definition, the probability that k points out of the n are not in general position. Since the points are vertices of $Q^d$ this probability is bounded by the probability that a $(d+1)\times(d+1)$ random $\pm 1$ matrix is singular and this probability is known [6] to go as $O\left(\frac{1}{\sqrt{d+1}}\right)$ when $d \to \infty$ so we have

$$(5) \qquad \frac{\sum\limits_{\Pi_n} a_k(\Pi_n,d)}{\binom{2^d}{n}} \le \binom{n}{k} O\left(\frac{1}{\sqrt{d+1}}\right) \qquad \text{when } d \to \infty$$

with this relation, observing that all quantities are positive, (4) gives

$$P(n,d) - O\left(\frac{1}{\sqrt{d+1}}\right)\frac{\sum\limits_{k=2}^{d+1}\binom{n}{k}}{2^{n-1}} \le H(n,d) \le P(n,d)$$

and being the fraction limited between 0 and 1 for every n this proves (2).

A similar argument can be used to study the number of linear separations of vertices of an hypercube[1]. It is intuitive that for n random hypercube vertices two different cases exist. If $n \ll d$ hypercube symmetries are irrelevant and the number of linear separations will equal that of n points in general position while if $n \approx 2^d$ symmetries play a crucial role diminishing the number of linear separations. In what follows I prove a condition that n has to satisfy (in the large d limit) to remain in the case where hypercube symmetries are marginal.

---

[1] In the past a lot of effort has been dedicated to this problem i.e. to count the number of thresholding functions (see e.g. [5]).

Starting from (3) we obtain

$$1 - \frac{\sum\limits_{\Pi_n} \sum\limits_{k=2}^{d+1} a_k(\Pi_n, d)}{C_{gp}(n,d) \binom{2^d}{n}} \leq \frac{<C(n,d)>}{C_{gp}(n,d)} \leq 1$$

where $<C(n,d)>$ is the average value of $C(n,d)$. Using the definition of $C_{gp}(n,d)$ (1) and (5)

$$1 - O\left(\frac{1}{\sqrt{d+1}}\right) \frac{\sum\limits_{k=2}^{d+1} \binom{n}{k}}{\sum\limits_{k=0}^{d} \binom{n-1}{k}} \leq \frac{<C(n,d)>}{C_{gp}(n,d)} \leq 1$$

and from the asymptotic properties of this fraction for $n > 2(d+1)$ and $d \to \infty$ we get

$$1 - O\left(\frac{n}{(d+1)^{\frac{3}{2}}}\right) \leq \frac{<C(n,d)>}{C_{gp}(n,d)} \leq 1$$

that proves that if $n < O\left((d+1)^{\frac{3}{2}}\right)$ the average number of separating hyperplanes of n vertices of $Q^d$ tends to $C_{gp}(n,d)$.

All this shows that as long as $n < O\left((d+1)^{\frac{3}{2}}\right)$ while $d \to \infty$ hypercube symmetries are not important for the average number of separating hyperplanes. From this follows that the probability of linear separability around $n = 2(d+1)$ is not altered by hypercube symmetries. Both these properties derive from the result that the probability of a d×d binary matrix being singular goes as $O\left(\frac{1}{\sqrt{d}}\right)$ when $d \to \infty$.

A final word of caution about the hypothesis of randomness in the choice of the n points that underlies all these results. In real life cases the patterns are highly correlated among themselves and these results do not apply directly.

4

# References

[1]     Cover T.M., *Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition*, IEEE Transactions on Electronic Computers **EC 14** (June 1965), pp. 326-334;

[2]     Baum E.B., *On the Capabilities of Multilayer Perceptrons*, Journal of Complexity, **4** (1988), pp. 193-215;

[3]     Mitchison G.J. and Durbin R.M., *Bounds on the Learning Capacity of Some Multi-Layer Networks*, Biological Cybernetics, **60** (1989), pp. 345-356;

[4]     Füredi Z., *Random Polytopes in the d-Dimensional Cube*, Discrete Computational Geometry **1** (1986), pp. 315-319;

[5]     Winder R.O., *Partitions of N-Space by Hyperplanes*, SIAM Journal on Applied Mathematics, **14** (4) (1966), pp. 811-818;

[6]     Komlós J., *On the Determinant of (0,1)-Matrices*, Studia Scientiarum Mathematicarum Hungarica, **2** (1967), pp. 7-21; these and more recent results are neatly reported in:
Bollobás B., *Random Graphs*, Academic Press 1985, pp. xvi 448, at pages 347-350.