

- [Ste78] J.M. Steele. Existence of submatrices with all possible columns. *Journal of Combinatorial Theory(A)*, 24:84–88, 1978.
- [Tal87a] M. Talagrand. Donsker classes and random geometry. *Annals of Probability*, 15:1327–1338, 1987.
- [Tal87b] M. Talagrand. The Glivenko-Cantelli problem. *Annals of Probability*, 15:837–870, 1987.
- [Tal88] M. Talagrand. Donsker classes of sets. *Probability Theory and Related Fields*, 78:169–191, 1988.
- [Tal91] M. Talagrand. Sharper bounds for empirical processes. 1991. manuscript.
- [Vap82] V. N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, New York, 1982.
- [VC71] V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–80, 1971.
- [VW91] L. G. Valiant and M. Warmuth, editors. *Proceedings of the 1991 Workshop on Computational Learning Theory*. Morgan Kaufmann, San Mateo, CA, 1991.
- [Wel88] E. Welzl. Partition trees for triangle counting and other range search problems. In *Proc. 4th Ann. ACM Symp. on Computational Geometry*, pages 23–33, 1988.

- [GZ84] E. Giné and J. Zinn. Some limit theorems for empirical processes. *Annals of Probability*, 12:929–989, 1984.
- [Hau91] David Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 1991. To appear.
- [HKS91] D. Haussler, M. Kearns, and R. Schapire. Bounds on the sample complexity of Bayesian learning using information theory and the VC dimension. In *Proceedings of the Fourth Workshop on Computational Learning Theory*, pages 61–74, 1991.
- [HL91] D. Haussler and P. Long. A generalization of Sauer’s lemma. In *Proc. of Southeastern Conference on Combinatorics*, 1991. To appear.
- [HLW90] David Haussler, Nick Littlestone, and Manfred Warmuth. Predicting $\{0, 1\}$ -functions on randomly drawn points. Technical Report UCSC-CRL-90-54, University of California Santa Cruz, Computer Research Laboratory, December 1990. To appear in *Information and Computation*.
- [HP88] David Haussler and Leonard Pitt, editors. *Proceedings of the 1988 Workshop on Computational Learning Theory*. Morgan Kaufmann, San Mateo, CA, 1988.
- [HW87] David Haussler and Emo Welzl. Epsilon nets and simplex range queries. *Disc. Comp. Geometry*, 2:127–151, 1987.
- [MSW90] J. Matousek, R. Seidel, and E. Welzl. How to net a lot with a little: small epsilon-nets for disks and halfspaces. In *Proc. 6th Ann. ACM Symp. on Computational Geometry*, 1990.
- [Pol84] David Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, 1984.
- [Pol90] David Pollard. *Empirical Processes: Theory and Applications*, volume 2 of *NSF-CBMS Regional Conference Series in Probability and Statistics*. Institute of Math. Stat. and Am. Stat. Assoc., 1990.
- [RHW89] Ron Rivest, David Haussler, and Manfred Warmuth, editors. *Proceedings of the 1989 Workshop on Computational Learning Theory*. Morgan Kaufmann, San Mateo, CA, 1989.
- [Sau72] N. Sauer. On the density of families of sets. *Journal of Combinatorial Theory (Series A)*, 13:145–147, 1972.

- [Ass83] Patrice Assouad. Densité et dimension. *Annales de l'Institut Fourier*, 33(3):233–282, 1983.
- [BEHW89] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the Association for Computing Machinery*, 36(4):929–965, 1989.
- [Bon72] J.A. Bondy. Induced subsets. *J. Comb. Theory (B)*, 12:201–202, 1972.
- [CBL91] N. Cesa-Bianchi and P. Long. Unpublished manuscript, 1991.
- [CF88] B. Chazelle and J. Friedman. A deterministic view of random sampling and its use in geometry. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, pages 539–549. IEEE, 1988.
- [CW89] B. Chazelle and E. Welzl. Quasi-optimal range searching and vc-dimensions. *Discrete and Computational Geometry*, 4:467–490, 1989.
- [Dud78] R. M. Dudley. Central limit theorems for empirical measures. *Ann. Prob.*, 6(6):899–929, 1978.
- [Dud84] R. M. Dudley. A course on empirical processes. *Lecture Notes in Mathematics*, 1097:2–142, 1984.
- [Dud87] R. M. Dudley. Universal Donsker classes and metric entropy. *Ann. Prob.*, 15(4):1306–1326, 1987.
- [EGS88] H. Edelsbrunner, L. Guibas, and M. Sharir. The complexity of many faces in arrangements of lines and of segments. In *Proc. 4th Ann. ACM Symp. on Computational Geometry*, pages 44–55, 1988.
- [FC90] Mark Fulk and John Case, editors. *Proceedings of the 1990 Workshop on Computational Learning Theory*. Morgan Kaufmann, San Mateo, CA, 1990.
- [Fra83] Peter Frankl. On the trace of finite sets. *Journal of Combinatorial Theory(A)*, 34:41–45, 1983.
- [Fra87] Peter Frankl. The shifting technique in extremal set theory. In C. Whitehead, editor, *Surveys in Combinatorics*, pages 81–110. Cambridge University Press, 1987.
- [GKP89] R. Graham, D. E. Knuth, and O. Patashnik. *Concrete Mathematics*. Addison-Wesley, 1989.

3 Remarks

As argued in the proof of Theorem 1, the result given in Lemma 3 implies that if i_1, \dots, i_m are selected at random without replacement, then the expectation of $\mathbf{Var}(V_{i_m} | V_{i_1}, \dots, V_{i_{m-1}})$ is at most d/m . When $m \gg d$, this means that the value of V_{i_m} is usually highly predictable given the values of $V_{i_1}, \dots, V_{i_{m-1}}$. This is the basis for many of the applications of the Vapnik-Chervonenkis dimension in machine learning and statistics.

In particular, from a Bayesian perspective, we might imagine that a vector $\vec{v} \in V$ is selected at random according to a "prior" distribution P on V and hidden from us, indices i_1, \dots, i_m are selected uniformly at random without replacement, we are given $v_{i_1}, \dots, v_{i_{m-1}}$, and we are asked to predict v_{i_m} . Bayes optimal strategy is to compute the posterior probabilities $P(v_{i_m} = 1 | V_{i_1} = v_{i_1}, \dots, V_{i_{m-1}} = v_{i_{m-1}})$ and $P(v_{i_m} = 0 | V_{i_1} = v_{i_1}, \dots, V_{i_{m-1}} = v_{i_{m-1}})$, and predict according to which of these is larger. The probability that this prediction is wrong when \vec{v} and i_1, \dots, i_m are chosen randomly as above, i.e. the Bayes risk, is the expectation of the minimum of the above two posterior probabilities. Looking into the proof of Lemma 3, it can be seen that this quantity is the same as

$$\frac{1}{m} \sum_{(\vec{u}, \vec{v}) \in E} \min(P(\vec{u}), P(\vec{v})).$$

Hence, by the last series of inequalities in the proof of Lemma 3, the Bayes risk for this prediction problem is at most d/m for any prior P , as was shown in [HKS91]. As pointed out there and in [HLW90], in fact the weighting scheme in [HLW90] gives the stronger result that there exists a (non-Bayesian) prediction strategy such that if i_1, \dots, i_m are chosen randomly without replacement, then given only the values $v_{i_1}, \dots, v_{i_{m-1}}$, the value v_{i_m} can be predicted such that for all $\vec{v} \in V$, the probability of a mistake is at most d/m , i.e. the minimax risk of this prediction problem is at most d/m .

4 Acknowledgements

I would like to thank Phil Long, Michael Kearns, Michel Talagrand, and Andrew Barron for helpful discussions of this material.

References

- [AHW87] N. Alon, D. Haussler, and E. Welzl. Partitioning and geometric embedding of range spaces of finite Vapnik-Chervonenkis dimension. In *Proc. 3rd Symp. on Computational Geometry*, pages 331–340. Waterloo, June 1987.
- [Alo83] Noga Alon. On the density of sets of vectors. *Discrete Math.*, 24:177–184, 1983.

The second bound of the theorem follows easily from this one. \square

We close with the proof of Theorem 2.

Let $W = \{(000\dots 0), (100\dots 0), (110\dots 0), \dots, (111\dots 1)\} \subset \{0, 1\}^s$, and $V = W^d$, the set of all vectors obtained by concatenating d vectors from W . Since $n = sd$, $V \subset \{0, 1\}^n$. It is easy to show that the Vapnik-Chervonenkis dimension of V is d : Say that indices $1 \leq i, j \leq n$ are equivalent if $\lceil i/s \rceil = \lceil j/s \rceil$. Then a sequence of indices is shattered by V if and only if it contains at most one index in each of the d equivalence classes. Thus no set of $d + 1$ indices is shattered. Note also that the size of V is $(s + 1)^d > s^d = (n/d)^d$.

For each $\vec{v} \in V$ and $1 \leq j \leq n$, let $N(\vec{v}, j)$ be the number of vectors $\vec{u} \in V$ with $\rho(\vec{u}, \vec{v}) = j/n$. Let $C(d, j)$ denote the number of ordered sequences of d non-negative integers that sum to j . We claim that for any $\vec{v} \in V$, $N(\vec{v}, j) \leq C(d, j)2^d$. This follows from the fact that there are at most $C(d, j)$ ways to choose the number of indices on which \vec{u} differs from \vec{v} in each of the d equivalence classes, and given any number of indices on which \vec{u} and \vec{v} must disagree in a given equivalence class, there are at most 2 choices for the values for \vec{u} on the indices in that equivalence class. Hence, using well known identities (see e.g. [GKP89])

$$\begin{aligned} \sum_{j=0}^k N(\vec{v}, j) &\leq 2^d \sum_{j=0}^k C(d, j) \\ &= 2^d \sum_{j=0}^k \binom{j+d-1}{j} \\ &= 2^d \binom{k+d}{k} \\ &< 2^d (\epsilon(k+d)/d)^d \\ &= (2\epsilon(k+d)/d)^d. \end{aligned}$$

Now choose any \vec{v}_1 in V , eliminate all vectors in V within ρ -distance k/n or less of \vec{v}_1 , then choose \vec{v}_2 from the remaining vectors in V and eliminate all vectors within distance k/n of \vec{v}_2 , etc., until V is exhausted. Since we begin with more than $(n/d)^d$ vectors, and each step eliminates at most $(2\epsilon(k+d)/d)^d$ vectors, this process continues for at least

$$\frac{(n/d)^d}{(2\epsilon(k+d)/d)^d} = \left(\frac{n}{2\epsilon(k+d)} \right)^d$$

steps, and the resulting set $\vec{v}_1, \vec{v}_2, \dots$ of vectors is clearly k/n -separated by construction. \square

the obvious way, i.e. $P_{|I}(u_1, \dots, u_m) = P\{\vec{v} \in V : v_{i_j} = u_j, 1 \leq j \leq m\}$. This projection does not change the conditional variances, hence the result follows.

Next we claim that

$$\begin{aligned} \gamma &= m \mathbf{E}[\mathbf{Var}(V_{i_m} | V_{i_1}, \dots, V_{i_{m-1}})] \\ &\geq m \left(\frac{k}{2(n-m+1)} \left(1 - \frac{|V_{\{i_1, \dots, i_{m-1}\}}|}{|V|} \right) \right) \\ &\geq m \left(\frac{k}{2(n-m+1)} \left(1 - \frac{(\epsilon(m-1)/d)^d}{|V|} \right) \right). \end{aligned}$$

The first equality follows by symmetry (and linearity of expectation), the second from Lemma 4, and the third from the Sauer/VC lemma (Lemma 1). Now putting these two claims together, we obtain

$$d \geq m \left(\frac{k}{2(n-m+1)} \left(1 - \frac{(\epsilon(m-1)/d)^d}{|V|} \right) \right)$$

or equivalently,

$$|V| \leq \frac{(\epsilon(m-1)/d)^d}{1 - \frac{2d(n-m+1)}{km}},$$

so long as

$$\frac{2d(n-m+1)}{km} < 1.$$

Now it is clear that

$$m-1 \leq \frac{(2d+2)(n+1)}{k+2d+2},$$

so

$$(\epsilon(m-1)/d)^d \leq \left(\left(\frac{\epsilon}{d} \right) \frac{(2d+2)(n+1)}{k+2d+2} \right)^d = \left((1+1/d) \left(\frac{2\epsilon(n+1)}{k+2d+2} \right) \right)^d \leq e \left(\frac{2\epsilon(n+1)}{k+2d+2} \right)^d.$$

In addition, it is easily verified that

$$\frac{2d(n-m+1)}{km} \leq \frac{2d(n+1 - \frac{(2d+2)(n+1)}{k+2d+2})}{k \left(\frac{(2d+2)(n+1)}{k+2d+2} \right)} = \frac{d}{d+1}.$$

Hence

$$\frac{1}{1 - \frac{2d(n-m+1)}{km}} \leq d+1.$$

Putting these together, this gives the final bound

$$|V| \leq e(d+1) \left(\frac{2\epsilon(n+1)}{k+2d+2} \right)^d.$$

probability that $\vec{u} \neq \vec{v}$, or $\epsilon n(1 - 1/N_j)/(n - m + 1)$. The variance $p(1 - p)$ of a Bernoulli random variable is just half the probability that the value of this random variable differs on two independent trials. Hence

$$\mathbf{E}[\mathbf{Var}(V_{i_m} | \vec{v} \in C_j)] \geq \frac{\epsilon n}{2(n - m + 1)} \left(1 - \frac{1}{N_j}\right),$$

where the expectation is over the random choice of i_m . For real x , let $x^+ = x$ if $x \geq 0$, else $x^+ = 0$. From the above we have

$$\begin{aligned} \mathbf{E}[\mathbf{Var}(V_{i_m} | V_{i_1}, \dots, V_{i_{m-1}})] &= \sum_{j=1}^M P(C_j) \mathbf{E}[\mathbf{Var}(V_{i_m} | \vec{v} \in C_j)] \\ &\geq \sum_{j=1}^M \left(\frac{N_j}{N}\right) \frac{\epsilon n}{2(n - m + 1)} \left(1 - \frac{1}{N_j}\right) \\ &= \frac{\epsilon n}{2(n - m + 1)N} \sum_{j=1}^M (N_j - 1)^+ \\ &\geq \frac{\epsilon n}{2(n - m + 1)N} \sum_{j=1}^M (N_j - 1) \\ &= \frac{\epsilon n}{2(n - m + 1)} \left(1 - \frac{M}{N}\right). \end{aligned}$$

□

We can now complete the proof of Theorem 1. Without loss of generality, let us assume that V itself is ϵ -separated, and obtain an upper bound on $|V|$. Let P be the uniform distribution on V . Recall that $k = \epsilon n$. We can assume that $k \geq 3$, since it can be verified that the upper bound given in the statement of the theorem is greater than the trivial upper bound from Lemma 1 when $k = 1$ or $k = 2$. Let us choose

$$m = \left\lceil \frac{(2d + 2)(n + 1)}{k + 2d + 2} \right\rceil$$

indices i_1, \dots, i_m uniformly at random without replacement³ from $\{1, \dots, n\}$ and look at

$$\gamma = \mathbf{E} \left[\sum_{j=1}^m \mathbf{Var}(V_{i_j} | V_{i_1}, \dots, V_{i_{j-1}}, V_{i_{j+1}}, \dots, V_{i_m}) \right].$$

We first claim that Lemma 3 implies that $\gamma \leq d$. This can be verified by projecting V onto $I = (i_1, \dots, i_m)$ and then defining the induced probability distribution $P|_I$ on $V|_I$ in

³Since $k \geq 3$ and $n \geq d, k$, it is easy to see that $m \leq n$.

To improve this upper bound to d , instead of directing the edges of E , we appeal to Lemma 2.7 of [HLW90], where it is shown that the vectors \vec{u}, \vec{v} of each edge $(\vec{u}, \vec{v}) \in E$ can be weighted with non-negative weights $w_{(\vec{u}, \vec{v})}(\vec{u})$ and $w_{(\vec{u}, \vec{v})}(\vec{v})$, resp., such that $w_{(\vec{u}, \vec{v})}(\vec{u}) + w_{(\vec{u}, \vec{v})}(\vec{v}) = 1$ for all $(\vec{u}, \vec{v}) \in E$, and for any vector $\vec{v} \in V$

$$\sum_{\vec{u} \in V: (\vec{u}, \vec{v}) \in E} w_{(\vec{u}, \vec{v})}(\vec{v}) \leq d.$$

For the sake of brevity, the argument needed to establish this is not repeated here. We then argue that

$$\begin{aligned} \sum_{i=1}^n \mathbf{Var}(V_i | V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n) &\leq \sum_{(\vec{u}, \vec{v}) \in E} \min(P(\vec{u}), P(\vec{v})) \\ &\leq \sum_{(\vec{u}, \vec{v}) \in E} w_{(\vec{u}, \vec{v})}(\vec{u})P(\vec{u}) + w_{(\vec{u}, \vec{v})}(\vec{v})P(\vec{v}) \\ &= \sum_{\vec{v} \in V} P(\vec{v}) \left(\sum_{\vec{u} \in V: (\vec{u}, \vec{v}) \in E} w_{(\vec{u}, \vec{v})}(\vec{v}) \right) \\ &\leq d \sum_{\vec{v} \in V} P(\vec{v}) \\ &= d \end{aligned}$$

□

The final lemma we will need in order to prove Theorem 1 is the following.

Lemma 4 *Suppose that V is an ϵ -separated subset of $\{0, 1\}^n$. Let P be the uniform distribution on V . For any integer m , $1 \leq m \leq n$, fix a sequence $I = (i_1, \dots, i_{m-1})$ of $m-1$ distinct indices between 1 and n and draw index i_m uniformly at random from the remaining $n-m+1$ indices. Then*

$$\mathbf{E}[\mathbf{Var}(V_{i_m} | V_{i_1}, \dots, V_{i_{m-1}})] \geq \frac{\epsilon n}{2(n-m+1)} \left(1 - \frac{|V_{I'}|}{|V|} \right),$$

where \mathbf{E} denotes expectation over the random choice of i_m .

Proof. Let us consider two vectors in V to be equivalent if they have the same value on all of the indices i_1, \dots, i_{m-1} in I . Suppose that this partitions V into $|V_{I'}| = M$ equivalence classes C_1, \dots, C_M . Let $N_j = |C_j|$ and $N = |V|$. Now let us focus on a single equivalence class C_j . Suppose that an additional index i_m is selected at random from the remaining $n-m+1$ indices, and two vectors \vec{u}, \vec{v} are selected uniformly at random with replacement from C_j . Since C_j is ϵ -separated, if $\vec{u} \neq \vec{v}$ then they differ on at least ϵn of the remaining $n-m+1$ indices. Hence the probability that $u_{i_m} \neq v_{i_m}$ is at least $\epsilon n / (n-m+1)$ times the

incident on \vec{v} . Since the density $|E|/|V|$ of the graph (V, E) is at most d by Lemma 2, if vectors are drawn uniformly from V , the average degree of a vector $\vec{v} \in V$ with respect to this graph is at most $2d$, since each edge gets counted twice when you sum the degrees of the nodes. As in [HLW90], we can thus direct the edges of this graph so that the outdegree of \vec{v} (number of edges in E directed away from \vec{v}) is at most $2d$ as follows. First find a vector $\vec{v} \in V$ whose degree is at most $2d$ and direct all the edges incident on \vec{v} away from \vec{v} . Then remove \vec{v} from V and iterate this construction on the subgraph induced by the remaining vectors in V . At each step we are guaranteed of finding a \vec{v} with degree at most $2d$ in the remaining graph because the density bound holds not only for the subgraph of the n -cube induced by V , but also for the subgraph induced by any subset of V . This is because any subset of V also has Vapnik-Chervonenkis dimension at most d , and hence Lemma 2 applies to it as well. When the construction is finished, it is clear that every edge has been directed, and no $\vec{v} \in V$ has outdegree more than $2d$ in the original graph. For each vector $\vec{v} \in V$ let $outdeg(\vec{v})$ denote the outdegree of \vec{v} and for each edge $(\vec{u}, \vec{v}) \in E$, let $tail(\vec{u}, \vec{v})$ denote the vector in the pair \vec{u}, \vec{v} that the edge is directed away from. We will use the directions on the edges in E shortly.

Now let us consider $\mathbf{Var}(V_i|V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n)$. Partition E into E_1, \dots, E_n by letting E_i be the edges that cross the i th dimension of the n -cube, i.e. $E_i = \{(\vec{u}, \vec{v}) \in E : u_j = v_j, j \neq i\}$. It is readily verified that

$$\begin{aligned} \mathbf{Var}(V_i|V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n) &= \sum_{(\vec{u}, \vec{v}) \in E_i} (P(\vec{u}) + P(\vec{v})) \frac{P(\vec{u})}{(P(\vec{u}) + P(\vec{v}))} \frac{P(\vec{v})}{(P(\vec{u}) + P(\vec{v}))} \\ &= \sum_{(\vec{u}, \vec{v}) \in E_i} \frac{P(\vec{u})P(\vec{v})}{(P(\vec{u}) + P(\vec{v}))}. \end{aligned}$$

Hence

$$\sum_{i=1}^n \mathbf{Var}(V_i|V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n) = \sum_{(\vec{u}, \vec{v}) \in E} \frac{P(\vec{u})P(\vec{v})}{(P(\vec{u}) + P(\vec{v}))}.$$

Now note that for any $x, y > 0$, $xy \leq (x + y)\min(x, y)$. Hence

$$\begin{aligned} \sum_{i=1}^n \mathbf{Var}(V_i|V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n) &\leq \sum_{(\vec{u}, \vec{v}) \in E} \min(P(\vec{u}), P(\vec{v})) \\ &\leq \sum_{(\vec{u}, \vec{v}) \in E} P(tail(\vec{u}, \vec{v})) \\ &= \sum_{\vec{v} \in V} P(\vec{v}) outdeg(\vec{v}) \\ &\leq 2d \sum_{\vec{v} \in V} P(\vec{v}) \\ &= 2d \end{aligned}$$

of edges of the subgraph of the n -cube induced by W . By the above results, $|W| = |V|$, $|F| \geq |E|$, and the Vapnik-Chervonenkis dimension of W is at most d .

Let us say that $\vec{u} \leq \vec{v}$ if $u_i \leq v_i$ for all i , $1 \leq i \leq n$. We claim that W is closed downward under the ordering \leq , in the sense that if $\vec{v} \in W$, then $\vec{u} \in W$ for all $\vec{u} \leq \vec{v}$. It is clear that if $\vec{u} \leq \vec{v} \in W$ and \vec{u} differs from \vec{v} on only one index i , then $\vec{u} \in W$: otherwise one more non-trivial shift of W would be possible. The claim follows by induction. It follows from this that if $\vec{v} \in W$, then the set of indices i for which $v_i = 1$ is shattered by W . Since the Vapnik-Chervonenkis dimension of W is at most d , this implies that no vector in W contains more than d ones. Therefore

$$|V| = |W| \leq \sum_{i=1}^d \binom{n}{i}$$

(which is the Sauer/VC lemma (Lemma 1)) and

$$|E|/|V| \leq |F|/|W| \leq d.$$

The last inequality can be verified by noting that a vector in $\{0, 1\}^n$ with at most d ones can have n -cube edges to at most d vectors with fewer ones. \square

Lemma 2 is the key in proving the next result, which is main tool we use in the proof of Theorem 1. It is closely related to the results obtained in [HKS91]. Let P be a probability distribution on V . Hence V can now be viewed as a vector-valued random variable. For each i , $1 \leq i \leq n$, let V_i be the i th component of the random variable V . Thus V_1, \dots, V_n are correlated Bernoulli random variables, and the value of V_i is determined by choosing $\vec{v} \in V$ at random by the distribution P , and letting $V_i = v_i$. Recall that for any Bernoulli random variable B , the variance of B is $p(1-p)$, where p is $P(B=1)$, and for Bernoulli random variables B_1, \dots, B_m , the *conditional variance* of B_m given B_1, \dots, B_{m-1} is defined by

$$\mathbf{Var}(B_m | B_1, \dots, B_{m-1}) = \sum_{\vec{v} \in \{0,1\}^{m-1}} P(\vec{v}) P(B_m = 1 | \vec{v}) (1 - P(B_m = 1 | \vec{v})),$$

where for $\vec{v} = (v_1, \dots, v_{m-1})$, $P(\vec{v}) = P(B_1 = v_1, \dots, B_{m-1} = v_{m-1})$ and $P(B_m = 1 | \vec{v}) = P(B_m = 1 | B_1 = v_1, \dots, B_{m-1} = v_{m-1})$.

Lemma 3 *For any probability distribution P on V ,*

$$\sum_{i=1}^n \mathbf{Var}(V_i | V_1, \dots, V_{i-1}, V_{i+1}, \dots, V_n) \leq d.$$

Proof. We will show here how an upper bound of $2d$ can be obtained on the above quantity, and then indicate how further results from [HLW90] can be used to improve this bound to d . Let E be the set of edges of the subgraph of the n -cube induced by V , as above. Let the *degree* of the vector \vec{v} in the graph (V, E) be the number of edges in E

Lemma 2 ([HLW90])

$$|E|/|V| \leq d.$$

Although this result is already proved in [HLW90], for completeness we provide an alternate proof here. This proof was suggested to us by Nati Lineal, and uses the simple technique of *shifting* [Fra87,Ste78,Alo83,Tal88] in place of the recursion in [HLW90].

Proof. For each index i , $1 \leq i \leq n$, and each $\vec{v} \in V$, if $v_i = 1$ and the vector $\vec{v}' = (v_1, \dots, v_{i-1}, 0, v_{i+1}, \dots, v_n)$ is not in V , then let $S_{i,V}(\vec{v}) = \vec{v}'$ (here we say that \vec{v} is *shifted* to \vec{v}'), otherwise let $S_{i,V}(\vec{v}) = \vec{v}$. We define the *shift* of V on index i , denoted $S_i(V)$, by

$$S_i(V) = \{S_{i,V}(\vec{v}) : \vec{v} \in V\}.$$

Let $S_i(E)$ denote the set of edges in the subgraph of the n -cube induced by $S_i(V)$. We claim that

1. $|S_i(V)| = |V|$,
2. $|S_i(E)| \geq |E|$, and
3. for any index set I , if I is shattered by $S_i(V)$ then I is shattered by V . Hence the Vapnik-Chervonenkis dimension of $S_i(V)$ is no more than that of V [Alo83].

The first claim is obvious. To verify the second claim, we map the edges of E in a 1-1 manner into the edges of $S_i(E)$. Assume $(\vec{u}, \vec{v}) \in E$. If neither \vec{u} nor \vec{v} are shifted then this edge is unaffected by the shift, so map it to itself. If both \vec{u} and \vec{v} are shifted then this edge is simply mapped to the edge $(S_{i,V}(\vec{u}), S_{i,V}(\vec{v}))$. Finally, let us assume that \vec{v} is shifted, but \vec{u} is not. In this case \vec{u} and \vec{v} must differ on some index $j \neq i$, and we must have $u_i = v_i = 1$. Since \vec{u} is not shifted, $\vec{u}' = (u_1, \dots, u_{i-1}, 0, u_{i+1}, \dots, u_n) \in V$. It follows that $(\vec{u}', S_{i,V}(\vec{v})) \in S_i(E)$. Hence we can map (\vec{u}, \vec{v}) to $(\vec{u}', S_{i,V}(\vec{v}))$. It is easily verified that the resulting map is 1-1.

To verify the third claim, suppose that a sequence I of k indices is shattered by $S_i(V)$. If i is not in I , then clearly I is also shattered by V , since $V|_I = S_i(V)|_I$ in this case. So let us assume that i is in I . Without loss of generality, we may assume that $i = 1$ and $I = (1, \dots, k)$. Since I is shattered by $S_i(V)$, for every $\vec{u} \in \{0, 1\}^k$ there is a $\vec{v} \in S_i(V)$ with $v_j = u_j$, $1 \leq j \leq k$. However, if $u_1 = 1$ then we must have \vec{v} and $\vec{v}' = (0, v_2, \dots, v_n)$ both in V , otherwise \vec{v} would have been shifted, and hence not be in $S_i(V)$. This implies that I is shattered by V , establishing the last claim.

Now beginning with V , simply shift V repeatedly on any sequence of (not necessarily distinct) indices until no more non-trivial shifts are possible, i.e. until you obtain a set W such that $S_i(W) = W$ for all $1 \leq i \leq n$. This must happen eventually, since each non-trivial shift reduces the total number of ones in the vectors of V . Let F be the set

This shows that the L_1 sphere packing numbers for arbitrary sets with Vapnik-Chervonenkis dimension d behave something like L_1 sphere packing numbers for bounded regions of d -dimensional Euclidean space.

Note that for $k = 1$ (i.e. $\epsilon = 1/n$), any two distinct vectors in V are ϵ -separated, thus the first bound gives a result similar to the Sauer/VC bound (Lemma 1) in this case, although not as tight. However, for larger values of ϵ , Theorem 1 improves on the best previous result, which is

$$\mathcal{M}(\epsilon, V) \leq \left(\frac{c_0}{\epsilon} \log \frac{1}{\epsilon} \right)^d,$$

where c_0 is some constant, obtained using the method of Dudley [Dud78] (see [Hau91] for a bound on the constants in Dudley's result). By getting rid of the extra log factor in Dudley's result, certain key bounds in the theory of empirical processes can also be improved by a logarithmic factor [Tal91].

It is likely that the constant 2ϵ in our result can be further improved. (It certainly can be improved for small d by using more precise upper estimates of $\sum_{i=0}^d \binom{n}{i}$ than that given in Lemma 1.) However, to within some multiplicative constant, the general form of the first bound of Theorem 1 is as tight as possible. This follows from the following lower bound.

Theorem 2 *For every natural numbers $d, s \geq 1$ there exists a subset $V \subset \{0, 1\}^n$, where $n = sd$, with Vapnik-Chervonenkis dimension d such that for each k , $1 \leq k \leq n$,*

$$\mathcal{M}(k/n, V) \geq \left(\frac{n}{2\epsilon(k+d)} \right)^d.$$

This leaves a gap from $1/2\epsilon$ to 2ϵ for the best universal value of the key constant in the bound of Theorem 1. Again, it is likely that the lower bound of Theorem 2 can be improved as well. However, at this time we do not have a good guess as to what the best possible constant is. This remains an intriguing open problem. It is also open² whether similar results hold for any of the various generalizations of the Vapnik-Chervonenkis dimension and the Sauer/VC lemma that have been studied (e.g. [Fra83,Dud87,Pol90,Hau91,HL91]).

2 Proofs of the Results

Throughout this section we assume that $V \subseteq \{0, 1\}^n$ and the Vapnik-Chervonenkis dimension of V is d . We begin with the following simple lemma from [HLW90].

Let E be the set of all pairs (\vec{u}, \vec{v}) with $\vec{u}, \vec{v} \in V$ such that $\rho(\vec{u}, \vec{v}) = 1/n$. Thus E is the set of edges in the subgraph of the Boolean n -cube induced by V (see also [Bon72,AHW87]).

²This question has been resolved recently for the pseudodimension [Pol90] by Phil Long and Nicolò Cesa-Bianchi [CBL91]. They show how Theorem 1 can be extended to get similar bounds on the L_1 sphere packing numbers for sets of vectors in n dimensional Euclidean space with pseudodimension d .

1 Statement of Results

Let n be natural number greater than zero. Let $V \subseteq \{0, 1\}^n$. For a sequence of indices $I = (i_1, \dots, i_k)$, with $1 \leq i_j \leq n$, let $V|_I$ denote the projection of V onto I , i.e.

$$V|_I = \{(v_{i_1}, \dots, v_{i_k}) : (v_1, \dots, v_n) \in V\}.$$

If $V|_I = \{0, 1\}^k$ then we say that V *shatters* the index sequence I . The *Vapnik-Chervonenkis dimension* of V is the size of the longest index sequence I that is shattered by V [VC71] (this terminology comes from [HW87]). We will denote this number by d . Hence

$$d = \max\{k : \exists I = (i_1, \dots, i_k), 1 \leq i_j \leq n, \text{ with } V|_I = \{0, 1\}^k\}.$$

This quantity plays a important role in certain areas of statistics, in particular in the theory of empirical processes [Dud78, Vap82, GZ84, Dud84, Pol84, Tal87a, Tal87b, Tal88, Pol90]. It has also been used recently in the fields of computational geometry [HW87] [Wel88] [MSW90] [EGS88] [CF88] [CW89] and machine learning [BEHW89, HP88, RHW89, FC90, VW91].

Let $|V|$ denote the cardinality of V . The following result is well known, and was independently discovered by several people, including Sauer [Sau72] and Vapnik and Chervonenkis (see [Ass83] for a review, and also [Dud84]).

Lemma 1 (*Sauer/VC*) *If the Vapnik-Chervonenkis dimension of V is d , then*

$$|V| \leq \sum_{i=0}^d \binom{n}{i} \leq (en/d)^d,$$

where e is the base of the natural logarithm.

For vectors $\vec{u}, \vec{v} \in \{0, 1\}^n$, let

$$\rho(\vec{u}, \vec{v}) = \frac{1}{n} \sum_{i=1}^n |u_i - v_i|.$$

For any $\epsilon > 0$, a set of vectors $W \subset \{0, 1\}^n$ is ϵ -*separated* if for all distinct $\vec{u}, \vec{v} \in W$, $\rho(\vec{u}, \vec{v}) \geq \epsilon$. The ϵ *packing number* for a set $V \subseteq \{0, 1\}^n$, denoted $\mathcal{M}(\epsilon, V)$, is the cardinality of the largest ϵ -separated subset W of V . Thus for integer r , $\mathcal{M}((2r+1)/n, V)$ is the largest set of disjoint L_1 balls of radius r/n with centers in V , or equivalently, the size of the largest r -bit error correcting code contained in V . In this paper we demonstrate the following result.

Theorem 1 *If the Vapnik-Chervonenkis dimension of V is d and $\epsilon = k/n$ for integer k , $1 \leq k \leq n$, then*

$$\mathcal{M}(\epsilon, V) \leq \epsilon(d+1) \left(\frac{2e(n+1)}{k+2d+2} \right)^d \leq \epsilon(d+1) \left(\frac{2e}{\epsilon} \right)^d.$$

**Sphere Packing Numbers
for Subsets of the Boolean n -Cube
with Bounded Vapnik-Chervonenkis Dimension**

David Haussler¹
haussler@cse.ucsc.edu

UCSC-CRL-91-41
October, 1991, revised March, 1992

Department of Computer and Information Sciences
University of California, Santa Cruz, CA 95064
and
Mathematical Sciences Research Institute
Berkeley, CA

Abstract: Let $V \subseteq \{0, 1\}^n$ have Vapnik-Chervonenkis dimension d . Let $\mathcal{M}(k/n, V)$ denote the cardinality of the largest $W \subseteq V$ such that any two distinct vectors in W differ on at least k indices. We show that $\mathcal{M}(k/n, V) \leq (cn/(k+d))^d$ for some constant c . This improves on the previous best result of $((cn/k) \log(n/k))^d$. This new result has applications in the theory of empirical processes.

¹The author gratefully acknowledges the support of the Mathematical Sciences Research Institute at UC Berkeley and ONR grant N00014-91-J-1162.