

- [Tou89] David Touretsky. *Advances in Neural Information Processing Systems*, volume 1. Morgan Kaufmann, 1989.
- [Tou90] David Touretsky. *Advances in Neural Information Processing Systems*, volume 2. Morgan Kaufmann, 1990.
- [Val84] Leslie G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, November 1984.
- [Vap82] V. N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, New York, 1982.
- [Vap89] V. N. Vapnik. Inductive principles of the search for empirical dependences (methods based on weak convergence of probability measures). In *Proceedings of the 2nd Workshop on Computational Learning Theory*, San Mateo, CA, 1989. published by Morgan Kaufmann.
- [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*, XVI(2):264–280, 1971.
- [WD81] R. S. Wencur and R. M. Dudley. Some special Vapnik-Chervonenkis classes. *Discrete Mathematics*, 33:313–318, 1981.
- [Wel88] Emo Welzl. Partition trees for triangle counting and other range search problems. In *4th ACM Symp. on Comp. Geometry*, pages 23–33, Urbana, IL, 1988.
- [Whi90a] Halbert White. Connectionist nonparametric regression: Multilayer feedforward networks can learn arbitrary mappings. *Neural Networks*, 3:535–549, 1990.
- [Whi90b] Halbert White. Learning in artificial neural networks: a statistical perspective. *Neural Computation*, 1(4):425–464, 1990.
- [WHR90] A. Weigend, B. Huberman, and D. Rumelhart. Predicting the future: A connectionist approach. *International Journal of Neural Systems*, 1:193–209, 1990.
- [WK91] S. Weiss and C. Kulikowski. *Computer Systems that Learn*. Morgan Kaufmann, San Mateo, CA, 1991.
- [Yam90] Kenji Yamanishi. A learning criterion for stochastic rules. In *Proceedings of the 3rd Workshop on Computational Learning Theory*, pages 67–81. published by Morgan Kaufmann, 1990.

- [Pol84] David Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, 1984.
- [Pol86] David Pollard. Rates of uniform almost-sure convergence for empirical processes indexed by unbounded classes of functions. manuscript, 1986.
- [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.
- [Qui89] Adolfo J. Quiroz. Metric entropy and learnability. Unpublished manuscript, Universidad Simón Bolívar, Caracas, Venezuela, 1989.
- [Ren70] A. Renyi. *Probability Theory*. North Holland, Amsterdam, 1970.
- [Ris86] Jorma Rissanen. Stochastic complexity and modeling. *The Annals of Statistics*, 14(3):1080–1100, 1986.
- [RM86] D. E. Rumelhart and J. L. McClelland. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 1: Foundations*. MIT Press, Cambridge, Mass., 1986.
- [Sau72] N. Sauer. On the density of families of sets. *Journal of Combinatorial Theory (Series A)*, 13:145–147, 1972.
- [Sim63] George F. Simmons. *Introduction to Topology and Modern Analysis*. McGraw-Hill, New York, 1963.
- [Slo88] Robert Sloan. Types of noise in data for concept learning. In *Proc. 1988 Workshop on Comp. Learning Theory*, pages 91–96, San Mateo, CA, 1988. Morgan Kaufmann.
- [STS90] H. Sompolinsky, N. Tishby, and H.S. Seung. Learning from examples in large neural networks. *Phys.Rev.Lett.*, 65:1683–1686, 1990.
- [SV88] George Shackelford and Dennis Volper. Learning k -DNF with noise in the attributes. In *Proceedings of the 1988 Workshop on Computational Learning Theory*, pages 97–103, San Mateo, CA, 1988. published by Morgan Kaufmann.
- [Tal91] M. Talagrand. 1991. manuscript.
- [TLS89] N. Tishby, E. Levin, and S. Solla. Consistent inference of probabilities in layered networks: predictions and generalizations. In *IJCNN International Joint Conference on Neural Networks*, volume II, pages 403–409. IEEE, 1989.

- [Nat88] B. K. Natarajan. Learning over classes of distributions. In *Proceedings of the 1988 Workshop on Computational Learning Theory*, pages 408–409, San Mateo, CA, 1988. published by Morgan Kaufmann.
- [Nat89a] B. K. Natarajan. Probably-approximate learning over classes of distributions. Technical Report HPL-SAL-89-29, Hewlett Packard Labs, Palo Alto, CA, 1989.
- [Nat89b] B. K. Natarajan. Some results on learning. Technical Report CMU-RI-TR-89-6, Carnegie Mellon, 1989.
- [ND90] Andrew Nobel and Amir Dembo. On uniform convergence for dependent processes. Technical Report 74, Stanford University, dept. of Statistics, Stanford, CA, 1990.
- [NH91] S. Nowlan and G. Hinton. Soft weight-sharing. Technical report, Dept. of Comp. Sci., U. Toronto, 1991.
- [Now90] S. Nowlan. Maximum likelihood competitive learning. In D. Touretsky, editor, *Advances in Neural Information Processing Systems*, volume 2, pages 574–582. Morgan Kaufmann, 1990.
- [NP87] Deborah Nolan and David Pollard. U-processes: Rates of convergence. *Annals of Statistics*, 15(2):780–799, 1987.
- [NT88] Balas K. Natarajan and P. Tadepalli. Two new frameworks for learning. In *Proceedings of the 5th International Conference on Machine Learning*, pages 402–415, San Mateo, CA, 1988. published by Morgan Kaufmann.
- [NT89] Kumpati S. Narendra and M. A. L. Thathachar. *Learning Automata – An Introduction*. Prentice Hall, 1989.
- [OH91a] M. Opper and D. Haussler. Calculation of the learning curve of Bayes optimal classification algorithm for learning a perceptron with noise. In *Computational Learning Theory: Proceedings of the Fourth Annual Workshop*, pages 75–87. Morgan Kaufmann, 1991.
- [OH91b] M. Opper and D. Haussler. Generalization performance of Bayes optimal classification algorithm for learning a perceptron. *Physical Review Letters*, 66(20):2677–2680, May 1991.
- [PG89] Tomaso Poggio and Federico Girosi. A theory of networks for approximation and learning. Technical Report A.I. Memo No. 1140, Massachusetts Institute of Technology, Cambridge, MA, 1989.

- [KLPV87] Michael Kearns, Ming Li, Leonard Pitt, and Leslie Valiant. On the learnability of boolean formulae. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, pages 285–295, New York, New York, May 1987.
- [KS90] Michael Kearns and Robert Schapire. Efficient distribution-free learning of probabilistic concepts. In *31th Annual IEEE Symposium on Foundations of Computer Science*, pages 382–391, 1990.
- [KT61] A. N. Kolmogorov and V. M. Tihomirov. ϵ -entropy and ϵ -capacity of sets in functional spaces. *Amer. Math. Soc. Translations (Ser. 2)*, 17:277–364, 1961.
- [Kul59] Solomon Kullback. *Information Theory and Statistics*. Wiley, New York, 1959.
- [Kul89] Sanjeev Kulkarni. On metric entripty, Vapnik-Chervonenkis dimension, and learnability for a class of distributions. Technical Report LIDS-P-1910, Center for Intelligent Control Systems, MIT, 1989.
- [LDS90] Y. LeCun, J. Denker, and S. Solla. Optimal brain damage. In D. Touretsky, editor, *Advances in Neural Information Processing Systems*, volume 2, page 589. Morgan Kaufmann, 1990.
- [Lin90] D. V. Lindley. The present position in Bayesian statistics. *Statistical Science*, 5(1):44–89, 1990.
- [Lit88] Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
- [LMR88] N. Lineal, Y. Mansour, and R. Rivest. Results on learnability and the Vapnik-Chervonenkis dimension. In *Proceedings of the 1988 Workshop on Computational Learning Theory*, pages 56–68, San Mateo, CA, 1988. published by Morgan Kaufmann.
- [Mac92] D. MacKay. *Bayesian Methods for Adaptive Models*. PhD thesis, California Institute of Technology, 1992.
- [Man82] Benoit B. Mandelbrot. *The Fractal Geometry of Nature*. W. H. Freeman, 1982.
- [MD89] J. Moody and C. Darken. Fast learning in networks of locally-tuned processing units. *Neural Computation*, 1(2):281–294, 1989.
- [MN89] P. McCullagh and J. A. Nelder. *Generalized Linear Models*. Chapman and Hall, London, 1989.

- [Far82] J. Doyne Farmer. Information dimension and the probabilistic structure of chaos. *Z. Naturforsch. A*, 37:1304–1325, 1982.
- [Fer67] Thomas Ferguson. *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press, 1967.
- [FOY83] J. Doyne Farmer, Edward Ott, and James A. Yorke. The dimension of chaotic attractors. *Physica D*, 7:153–180, 1983.
- [GT90] G. Gyorgyi and N. Tishby. Statistical theory of learning a rule. In K. Thue-mann and R. Koeberle, editors, *Neural Networks and Spin Glasses*. World Scientific, 1990.
- [Gul90] V. Gullapalli. A stochastic reinforcement algorithm for learning real-valued functions. *Neural Networks*, 3(6):671–692, 1990.
- [Hau88] David Haussler. Quantifying inductive bias: AI learning algorithms and Valiant’s learning framework. *Artificial Intelligence*, 36:177–221, 1988.
- [Hau89] David Haussler. Learning conjunctive concepts in structural domains. *Machine Learning*, 4:7–40, 1989.
- [Hau90] David Haussler. Decision theoretic generalizations of the pac learning model. In *Proc. First Workshop on Algorithmic Learning Theory*, pages 21–41, Yokyo, Japan, 1990.
- [HKLW91] David Haussler, Michael Kearns, Nick Littlestone, and Manfred K. Warmuth. Equivalence of models for polynomial learnability. *Information and Computation*, 95:129–161, 1991.
- [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.
- [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*.
- [HW87] David Haussler and Emo Welzl. Epsilon nets and simplex range queries. *Disc. Comp. Geometry*, 2:127–151, 1987.
- [Kie87] Jack Kiefer. *Introduction to Statistical Inference*. Springer-Verlag, 1987.

- [BW91] Wray Buntine and Andreas Weigend. Bayesian back propagation. Unpublished manuscript, 1991.
- [CB] Bertrand Clarke and Andrew Barron. Entropy, risk and the Bayesian central limit theorem. manuscript.
- [CB90] Bertrand Clarke and Andrew Barron. Information-theoretic asymptotics of Bayes methods. *IEEE Transactions on Information Theory*, 36(3):453–471, 1990.
- [Cov68] Thomas M. Cover. Capacity problems for linear machines. In L. Kanal, editor, *Pattern Recognition*, pages 283–289. Thompson Books, 1968.
- [Dev88] Luc Devroye. Automatic pattern recognition: A study of the probability of error. *IEEE Trans. on Pattern Anal. and Mach. Intelligence*, 10(4):530–543, 1988.
- [DH73] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. Wiley, 1973.
- [DR89] R. Durbin and D. E. Rumelhart. Product units: A computationally powerful and biologically plausible extension to backpropagation networks. *Neural Computation*, 1(1):133–142, 1989.
- [DSW⁺87] J. Denker, D. Schwartz, B. Wittner, S. Solla, R. Howard, L. Jackel, and J. Hopfield. Automatic learning, rule extraction and generalization. *Complex Syst.*, 1:877–922, 1987.
- [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.
- [Ede87] Herbert Edelsbrunner. *Algorithms in Combinatorial Geometry*. Springer-Verlag, 1987.
- [EHKV89] Andrzej Ehrenfeucht, David Haussler, Michael Kearns, and Leslie Valiant. A general lower bound on the number of examples needed for learning. *Information and Computation*, 82:247–261, 1989.

- [Ass83] Patrice Assouad. Densité et dimension. *Annales de l'Institut Fourier*, 33(3):233–282, 1983.
- [AST90] Martin Anthony and John Shawe-Taylor. A result of Vapnik with applications. Technical Report CSD-TR-628, University of London, Surrey, England, 1990.
- [AV79] Dana Angluin and Leslie G. Valiant. Fast probabilistic algorithms for hamiltonian circuits and matchings. *Journal of Computer and System Sciences*, 18(2):155–193, April 1979.
- [AW90] Naoki Abe and Manfred Warmuth. On the computational complexity of approximating distributions by probabilistic automata. In *Proceedings of the 3rd Workshop on Computational Learning Theory*, pages 52–66. published by Morgan Kaufmann, 1990.
- [BA85] A. G. Barto and P. Anandan. Pattern recognizing stochastic learning automata. *IEEE Trans. on Systems, Man and Cybernetics*, 15:360–374, 1985.
- [Bar89] Andrew Barron. Statistical properties of artificial neural networks. In *28th Conference on Decision and Control*, pages 280–285, 1989.
- [BC90] Andrew Barron and Thomas Cover. Minimum complexity density estimation. *IEEE Transactions on Information Theory*, 1990. To appear.
- [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.
- [Ber85] J. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag, New York, 1985.
- [BFOS84] Leo Breiman, Jerome H. Friedman, Richard A. Olshen, and Charles J. Stone. *Classification and Regression Trees*. Wadsworth International Group, 1984.
- [BH89] Eric Baum and D. Haussler. What size net gives valid generalization? *Neural Computation*, 1(1):151–160, 1989.
- [BI88] Gyora M. Benedek and Alon Itai. Learnability by fixed distributions. In *Proc. 1988 Workshop on Comp. Learning Theory*, pages 80–90, San Mateo, CA, 1988. Morgan Kaufmann.
- [Bil86] Patrick Billingsley. *Probability and Measure*. Wiley, New York, 1986.
- [Bun90] W.L. Buntine. *A Theory of Learning Classification Rules*. PhD thesis, University of Technology, Sydney, 1990.

we have

$$4(2eNm/W)^W e^{-\alpha^2 \nu m/2} \leq \delta.$$

To see this, first note that by rearranging, we get

$$\alpha^2 \nu m/2 \geq W \ln(2eNm/W) + \ln(4/\delta).$$

thus it suffices to show

$$\alpha^2 \nu m/4 \geq W \ln(2eNm/W) \text{ and } \alpha^2 \nu m/4 \geq \ln(4/\delta).$$

The latter inequality is assured by the last term in the formula for m . For the former inequality, let us take m equal to the first term only, i.e.

$$m = \frac{8W}{\alpha^2 \nu} \ln \frac{16N}{\alpha^2 \nu}.$$

Substituting this into the former inequality and simplifying, we get

$$2 \ln \frac{16N}{\alpha^2 \nu} \geq \ln \left(\frac{16eN}{\alpha^2 \nu} \ln \frac{16N}{\alpha^2 \nu} \right).$$

This further simplifies to

$$\frac{16N}{\alpha^2 \nu} \geq e \ln \frac{16N}{\alpha^2 \nu},$$

which holds, since $x \geq e \ln x$ for any x . Finally, since this inequality holds for the given m , it is easy to see that it will also hold for larger m . \square

References

- [AL88] Dana Angluin and Philip Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988.
- [Ale85] K. Alexander. Rates of growth for weighted empirical processes. In *Proc. of Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, volume 2, pages 475–493, 1985.
- [Ale87] K. Alexander. Rates of growth and sample moduli for weighted empirical processes indexed by sets. *Probability Theory and Related Fields*, 75:379–423, 1987.
- [Ang88] Dana Angluin. Queries and concept learning. *Machine Learning*, 2:319–342, 1988.

non-increasing. Different σ 's can be used for different units. Let \mathcal{H} be all functions from X into A representable on \mathcal{A} by varying the weights and biases. Let $\nu > 0$, $0 < \alpha < 1$ and $m \geq 1$. Suppose that \vec{z} is generated by m independent random draws from a probability measure P on $Z = X \times Y$. Then

$$\Pr(\exists h \in \mathcal{H} : d_\nu(\hat{\mathbf{r}}_{h,l}(\vec{z}), \mathbf{r}_{h,l}(P)) > \alpha) \leq 4(2eNm/W)^W e^{-\alpha^2 \nu m/2}.$$

This probability is at most δ for a sample size m that is

$$O\left(\frac{1}{\alpha^2 \nu} \left(W \log \frac{N}{\alpha \nu} + \log \frac{1}{\delta}\right)\right).$$

Proof. Each computation unit in the network with k weights is associated with a class of $\{0, 1\}$ -valued functions of the form

$$f(\vec{x}) = \text{sign}(\sigma(\phi_1(\vec{x}), \dots, \phi_k(\vec{x})) + \theta + \sum_{j=1}^k w_j \phi_j(\vec{x}))$$

where θ, w_1, \dots, w_k are adjustable real-valued parameters and $\phi_1, \dots, \phi_k, \mu$ and σ are fixed functions, the latter monotone. By Theorems 4 and 5 this class of functions has pseudo dimension at most $k + 1$. Since the pseudo dimension is the same as the Vapnik Chervonenkis dimension for classes of indicator functions, this implies that the class has Vapnik Chervonenkis dimension at most $k + 1$. Now let d be the sum of all the Vapnik Chervonenkis dimensions of all the classes of functions associated with the computation units of the architecture \mathcal{A} . It follows that $d \leq W$, the total number of weights and biases in the network.

For each $h \in \mathcal{H}$ let l_h be the loss function associated with h for the discrete loss l , i.e. $l_h(x, y) = 1$ if $y \neq h(x)$, $l_h(x, y) = 0$ if $y = h(x)$. Let $\mathbf{F} = l_{\mathcal{H}} = \{l_h : h \in \mathcal{H}\}$. Let $\vec{z} = ((x_1, y_1), \dots, (x_m, y_m))$ be any fixed sample and $\vec{x} = (x_1, \dots, x_m)$. It is easily verified that $|\mathbf{F}_{|\vec{z}}| = |\mathcal{H}_{|\vec{x}}|$. It is shown in [BH89], Theorem 1, (and is also implied directly from results in [Cov68]) that for any class \mathcal{H} as above, $|\mathcal{H}_{|\vec{x}}| \leq (Nem/d)^d$ for all $\vec{x} = (x_1, \dots, x_m)$, where d and N are as above. This implies that

$$|\mathbf{F}_{|\vec{z}}| \leq (Nem/d)^d \leq (Nem/W)^W$$

for all samples \vec{z} of length m . Since each $l_h \in \mathbf{F}$ is a random variable that is bounded between 0 and 1, Theorem 3 (second part) shows that

$$\Pr(\exists h \in \mathbf{F} : d_\nu(\hat{\mathbf{E}}_{\vec{z}}(h), \mathbf{E}(h)) > \alpha) \leq 4(2eNm/W)^W e^{-\alpha^2 \nu m/2}.$$

This gives the first bound.

For the second bound, it can be shown that for sample size

$$m = \frac{4}{\alpha^2 \nu} \left(2W \ln \frac{16N}{\alpha^2 \nu} + \ln \frac{4}{\delta}\right)$$

we must have

$$d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f^*), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f^*)) > \alpha/4.$$

Thus if N is an $\alpha/4$ -cover for $\mathbf{F}|_{\mathcal{Z}}$ with respect to the metric \vec{d}_ν , then whenever there exists $f \in \mathbf{F}$ with

$$d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2,$$

there exists $f^* \in N$ with

$$d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f^*), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f^*)) > \alpha/4.$$

For fixed f^* and random σ , the probability of the latter event is at most $2e^{-\alpha^2\nu m/8M}$ by Lemma 11. Hence for any fixed $\vec{z} \in Z^{2m}$, if we select a permutation $\sigma \in \mathcal{S}_{2m}$ uniformly at random,

$$\Pr(\exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2) \leq \min(2\mathcal{N}(\alpha/4, \mathbf{F}|_{\mathcal{Z}}, \vec{d}_\nu)e^{-\alpha^2\nu m/8M}, 1).$$

The remainder of the proof is as above. \square

We need only one more lemma to complete our proof; one that can be used to relate the \vec{d}_ν covering numbers to the d_{L^1} covering numbers.

Lemma 14 *For any $m \geq 1$, $\vec{x}, \vec{y} \in (\mathbb{R}^+)^{2m}$ and $\nu > 0$, $\vec{d}_\nu(\vec{x}, \vec{y}) \leq \frac{2}{\nu}d_{L^1}(\vec{x}, \vec{y})$.*

Proof. For any $\sigma \in \mathcal{S}_{2m}$

$$\begin{aligned} & d_\nu(\mu_1(\vec{x}, \sigma), \mu_1(\vec{y}, \sigma)) + d_\nu(\mu_2(\vec{x}, \sigma), \mu_2(\vec{y}, \sigma)) \\ &= \frac{|\sum_{i=1}^m (x_{\sigma(i)} - y_{\sigma(i)})|}{\nu m + \sum_{i=1}^m (x_{\sigma(i)} + y_{\sigma(i)})} + \frac{|\sum_{i=m+1}^{2m} (x_{\sigma(i)} - y_{\sigma(i)})|}{\nu m + \sum_{i=m+1}^{2m} (x_{\sigma(i)} + y_{\sigma(i)})} \\ &\leq \frac{\sum_{i=1}^{2m} |x_{\sigma(i)} - y_{\sigma(i)}|}{\nu m} = \frac{2}{\nu}d_{L^1}(\vec{x}, \vec{y}). \end{aligned}$$

The result follows. \square

The theorem follows easily from the last two lemmas. \square

10.5 Bounds on sample size for learning in feedforward nets with sharp thresholds

In this section we give the bounds for uniform convergence of empirical estimates in neural networks with sharp threshold functions claimed in section 7.

Theorem 13 *Let \mathcal{A} be a feedforward architecture as defined in section 7 with $n \geq 1$ inputs, one output, $N \geq 2$ computation units and a total of W weights and biases. Let $X = \mathbb{R}^n$, $Y = A = \{0, 1\}$, and l be the discrete loss function. Assume that the squashing function for each computation unit has the form $\text{sign} \circ \sigma$, where σ is non-decreasing or*

Proof. First note that both bounds are trivial if $m < 2M/(\alpha^2\nu)$, so we may assume $m \geq 2M/(\alpha^2\nu)$. Hence by Lemma 12,

$$p(\alpha, \nu, m) \leq 2\Pr \left\{ \vec{z} \in Z^{2m} : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_z(f), \widehat{\mathbf{E}}''_z(f)) > \alpha/2 \right\}.$$

Thus it suffices to obtain bounds for the latter quantity.

We begin with the second bound, for the case when $\mathbf{F}|_z$ is always finite. For any sample $\vec{z} = (z_1, \dots, z_{2m}) \in Z^{2m}$ and $\sigma \in \mathcal{S}_{2m}$, let $\sigma(\vec{z}) = (z_{\sigma(1)}, \dots, z_{\sigma(2m)})$. For any fixed function $f \in \mathbf{F}$ and fixed $\vec{z} \in Z^{2m}$, if we select a permutation $\sigma \in \mathcal{S}_{2m}$ uniformly at random,

$$\Pr \left(d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2 \right) \leq 2e^{-\alpha^2\nu m/2M} \quad (7)$$

by Lemma 11. Hence for any fixed $\vec{z} \in Z^{2m}$, if we select a permutation $\sigma \in \mathcal{S}_{2m}$ uniformly at random,

$$\Pr \left(\exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2 \right) \leq \min(2|\mathbf{F}|_z|e^{-\alpha^2\nu m/2M}, 1). \quad (8)$$

Thus if we draw \vec{z} at random from Z^{2m} and independently select a permutation $\sigma \in \mathcal{S}_{2m}$ uniformly at random,

$$\Pr \left(\exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2 \right) \leq \mathbf{E}(\min(2|\mathbf{F}|_z|e^{-\alpha^2\nu m/2M}, 1)).$$

However, since each of the $2m$ observations in \vec{z} are independent, each of the samples $\sigma(\vec{z})$ for $\sigma \in \mathcal{S}_{2m}$ are equally likely. Hence

$$\Pr \left(\exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2 \right),$$

where both \vec{z} and σ are chosen at random, is the same as

$$\Pr \left(\exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_z(f), \widehat{\mathbf{E}}''_z(f)) > \alpha/2 \right),$$

where only the sample \vec{z} is chosen at random. The second bound follows.

The proof of the first bound is similar, except for steps 7 and 8. From Lemma 10, using the extension \vec{d}_ν of the metric d_ν to even-length sequences of reals given in Definition 12 above, if $f, f^* \in \mathbf{F}$ are such that

$$\vec{d}_\nu((f(z_1), \dots, f(z_{2m})), (f^*(z_1), \dots, f^*(z_{2m}))) \leq \alpha/4,$$

then for any $\sigma \in \mathcal{S}_{2m}$ such that

$$d_\nu(\widehat{\mathbf{E}}'_{\sigma(\vec{z})}(f), \widehat{\mathbf{E}}''_{\sigma(\vec{z})}(f)) > \alpha/2,$$

Lemma 12 *Let \mathbf{F} be a permissible set of functions on Z with $0 \leq f(z) \leq M$ for all $f \in \mathbf{F}$ and $z \in Z$. Assume $\nu > 0$, $0 < \alpha < 1$ and $m \geq 2M/(\alpha^2\nu)$. Then*

$$\begin{aligned} & \Pr \left\{ \vec{z} \in Z^m : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha \right\} \\ & \leq 2\Pr \left\{ \vec{z} \in Z^{2m} : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\vec{z}}(f), \widehat{\mathbf{E}}''_{\vec{z}}(f)) > \alpha/2 \right\}. \end{aligned}$$

Proof. If Z and \mathbf{F} are uncountable, the assumption of permissibility guarantees that these probabilities are well-defined (see section 10.2 and [Pol84]). From Chebyshev's inequality (see Lemma 9, part (1) in section 10.3), for each individual $f \in \mathbf{F}$,

$$\Pr \left\{ \vec{z} \in Z^m : d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha/2 \right\} \leq M/(\alpha^2\nu m).$$

Since $m \geq 2M/(\alpha^2\nu)$, this probability is at most $1/2$. Now consider any $f \in \mathbf{F}$ and sample $\vec{z}' \in Z^m$ such that $d_\nu(\widehat{\mathbf{E}}_{\vec{z}'}(f), \mathbf{E}(f)) > \alpha$. If we draw an independent random sample $\vec{z}'' \in Z^m$, then with probability at least $1/2$, $d_\nu(\widehat{\mathbf{E}}_{\vec{z}''}(f), \mathbf{E}(f)) \leq \alpha/2$. Whenever this happens we have $d_\nu(\widehat{\mathbf{E}}_{\vec{z}'}(f), \widehat{\mathbf{E}}_{\vec{z}''}(f)) > \alpha/2$ by the triangle inequality for d_ν . Thus

$$\begin{aligned} & \Pr \left\{ \vec{z} \in Z^{2m} : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}'_{\vec{z}}(f), \widehat{\mathbf{E}}''_{\vec{z}}(f)) > \alpha/2 \right\} \\ & \geq \Pr \left\{ \vec{z}'\vec{z}'' \in Z^{2m} : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}_{\vec{z}'}(f), \mathbf{E}(f)) > \alpha \text{ and } d_\nu(\widehat{\mathbf{E}}_{\vec{z}''}(f), \mathbf{E}(f)) \leq \alpha/2 \right\} \\ & \geq \frac{1}{2} \Pr \left\{ \vec{z}' \in Z^m : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}_{\vec{z}'}(f), \mathbf{E}(f)) > \alpha \right\}. \end{aligned}$$

Again, when Z and \mathbf{F} are uncountable, permissibility guarantees that the implied use of Fubini's theorem in obtaining the above inequalities is justified. \square

We are now in a position to prove the following version of our theorem, using the extended metric \vec{d}_ν in place of the L^1 metric to measure covering numbers.

Lemma 13 *Let \mathbf{F} be a permissible set of functions on Z with $0 \leq f(z) \leq M$ for all $f \in \mathbf{F}$ and $z \in Z$. Assume $\nu > 0$, $0 < \alpha < 1$ and $m \geq 1$. Let*

$$p(\alpha, \nu, m) = \Pr \left\{ \vec{z} \in Z^m : \exists f \in \mathbf{F} \text{ with } d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha \right\}.$$

Then

$$p(\alpha, \nu, m) \leq 2\mathbf{E}(\min(2\mathcal{N}(\alpha/4, \mathbf{F}|_{\vec{z}}, \vec{d}_\nu)e^{-\alpha^2\nu m/8M}, 1)),$$

and if in addition $\mathbf{F}|_{\vec{z}}$ is finite for all $\vec{z} \in Z^{2m}$ then

$$p(\alpha, \nu, m) \leq 2\mathbf{E}(\min(2|\mathbf{F}|_{\vec{z}}|e^{-\alpha^2\nu m/2M}, 1)),$$

where the expectations are over \vec{z} drawn randomly from Z^{2m} .

Proof. For each i , $1 \leq i \leq m$, let Y_i be an independent random variable such that $Y_i = x_i - x_{m+i}$ with probability $1/2$ and $Y_i = x_{m+i} - x_i$ with probability $1/2$. Note that for any $\sigma \in \binom{[2m]}{m}$,

$$d_\nu(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) = \frac{|\frac{1}{m} \sum_{i=1}^m x_{\sigma(i)} - \frac{1}{m} \sum_{i=m+1}^{2m} x_{\sigma(i)}|}{\nu + \frac{1}{m} \sum_{i=1}^{2m} x_{\sigma(i)}} = \frac{|\sum_{i=1}^m (x_{\sigma(i)} - x_{\sigma(m+i)})|}{\nu m + \sum_{i=1}^{2m} x_i}.$$

Hence

$$\begin{aligned} & \Pr(d_\nu(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) > \alpha) \\ &= \Pr\left(\left|\sum_{i=1}^m (x_{\sigma(i)} - x_{\sigma(m+i)})\right| > \alpha(\nu m + \sum_{i=1}^{2m} x_i)\right) \\ &= \Pr\left(\left|\sum_{i=1}^m Y_i\right| > \alpha(\nu m + \sum_{i=1}^{2m} x_i)\right), \end{aligned}$$

because each swap in a randomly chosen $\sigma \in \binom{[2m]}{m}$ is independent. Since $\mathbf{E}(Y_i) = 0$ and $-|x_i - x_{m+i}| \leq Y_i \leq |x_i - x_{m+i}|$, we can apply Hoeffding's inequality (see e.g. [Pol84]) to bound the latter probability by

$$2e^{-\alpha^2(\nu m + \sum_{i=1}^{2m} x_i)^2 / 2 \sum_{i=1}^m (x_i - x_{m+i})^2}.$$

Let $\beta = \sum_{i=1}^{2m} x_i$. Since $0 \leq x_i \leq M$,

$$\sum_{i=1}^m (x_i - x_{m+i})^2 \leq \sum_{i=1}^m M|x_i - x_{m+i}| \leq \beta M.$$

Hence we have

$$2e^{-\alpha^2(\nu m + \sum_{i=1}^{2m} x_i)^2 / 2 \sum_{i=1}^m (x_i - x_{m+i})^2} \leq 2e^{-\alpha^2(\nu m + \beta)^2 / 2\beta M} = 2e^{-\frac{\alpha^2}{2M}(\frac{\nu m + \beta}{\beta})^2}.$$

The expression in parentheses is minimized, and therefore the whole expression is maximized, by setting $\beta = \nu m$, giving a value of $4\nu m$. Hence

$$\Pr(d_\nu(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) > \alpha) \leq 2e^{-2\alpha^2\nu m/M}.$$

□

For our next lemma we will need some notation to refer to the separate empirical estimates based on the first and second halves of an even length sample.

Definition 13 For all $m \geq 1$ and $\vec{z} \in Z^{2m}$, we let $\widehat{\mathbf{E}}'_z(f) = \frac{1}{m} \sum_{i=1}^m f(z_i)$ and $\widehat{\mathbf{E}}''_z(f) = \frac{1}{m} \sum_{i=m+1}^{2m} f(z_i)$.

10.4 Proof of Main Theorem on Uniform Convergence of Empirical Estimates

In this section we prove Theorem 3 from section 3. The proof is given in a series of lemmas. We first extend the metric d_ν to a pseudo-metric on vectors in $(\mathfrak{R}^+)^{2m}$ for $m \geq 1$. We will do so in a somewhat unusual manner that will be useful in what follows. In fact, this extension can be defined for any metric, so we will state it in its general form.

Definition 12 For each integer $m \geq 1$, let ${}_{,2m}$ denote the set of all permutations σ of $\{1, \dots, 2m\}$ such that for all i , $1 \leq i \leq m$, either $\sigma(i) = m + i$ and $\sigma(m + i) = i$, or $\sigma(i) = i$ and $\sigma(m + i) = m + i$. Thus the permutations in ${}_{,2m}$ swap selected indices in the first half of the sequence $\{1, \dots, 2m\}$ with corresponding indices in the second half. For any $\vec{x} = (x_1, \dots, x_{2m}) \in (\mathfrak{R}^+)^{2m}$ and $\sigma \in {}_{,2m}$, let $\mu_1(\vec{x}, \sigma) = \frac{1}{m} \sum_{i=1}^m x_{\sigma(i)}$ and $\mu_2(\vec{x}, \sigma) = \frac{1}{m} \sum_{i=m+1}^{2m} x_{\sigma(i)}$. For any metric d on \mathfrak{R}^+ and $\vec{x}, \vec{y} \in (\mathfrak{R}^+)^{2m}$, let $\vec{d}(\vec{x}, \vec{y}) = \max\{d(\mu_1(\vec{x}, \sigma), \mu_1(\vec{y}, \sigma)) + d(\mu_2(\vec{x}, \sigma), \mu_2(\vec{y}, \sigma)) : \sigma \in {}_{,2m}\}$.

It is easily verified that \vec{d} is a pseudo-metric on $(\mathfrak{R}^+)^{2m}$. Symmetry is obvious, and the triangle inequality follows easily from the triangle inequality for d on \mathfrak{R}^+ :

$$\begin{aligned} \vec{d}(\vec{x}, \vec{y}) + \vec{d}(\vec{y}, \vec{z}) &= \max\{d(\mu_1(\vec{x}, \sigma), \mu_1(\vec{y}, \sigma)) + d(\mu_2(\vec{x}, \sigma), \mu_2(\vec{y}, \sigma)) : \sigma \in {}_{,2m}\} \\ &\quad + \max\{d(\mu_1(\vec{y}, \sigma), \mu_1(\vec{z}, \sigma)) + d(\mu_2(\vec{y}, \sigma), \mu_2(\vec{z}, \sigma)) : \sigma \in {}_{,2m}\} \\ &\geq \max\{d(\mu_1(\vec{x}, \sigma), \mu_1(\vec{y}, \sigma)) + d(\mu_2(\vec{x}, \sigma), \mu_2(\vec{y}, \sigma)) \\ &\quad + d(\mu_1(\vec{y}, \sigma), \mu_1(\vec{z}, \sigma)) + d(\mu_2(\vec{y}, \sigma), \mu_2(\vec{z}, \sigma)) : \sigma \in {}_{,2m}\} \\ &\geq \max\{d(\mu_1(\vec{x}, \sigma), \mu_1(\vec{z}, \sigma)) + d(\mu_2(\vec{x}, \sigma), \mu_2(\vec{z}, \sigma)) : \sigma \in {}_{,2m}\} \\ &= \vec{d}(\vec{x}, \vec{z}) \end{aligned}$$

We note the following additional property of this extension.

Lemma 10 For all $\vec{x}, \vec{y} \in (\mathfrak{R}^+)^{2m}$ and $\sigma \in {}_{,2m}$, $d(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) \leq d(\mu_1(\vec{y}, \sigma), \mu_2(\vec{y}, \sigma)) + \vec{d}(\vec{x}, \vec{y})$.

Proof. We have

$$d(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) \leq d(\mu_1(\vec{y}, \sigma), \mu_2(\vec{y}, \sigma)) + d(\mu_1(\vec{x}, \sigma), \mu_1(\vec{y}, \sigma)) + d(\mu_2(\vec{x}, \sigma), \mu_2(\vec{y}, \sigma))$$

by the triangle inequality on d . The last two terms of this sum combined are at most $\vec{d}(\vec{x}, \vec{y})$ by definition. \square

We now restrict ourselves to the case that the metric d is the metric d_ν for some $\nu > 0$. The following lemma will play a key role in establishing our basic exponential inequality.

Lemma 11 Let $\vec{x} = (x_1, \dots, x_{2m})$ be a sequence of reals such that $0 \leq x_i \leq M$, $1 \leq i \leq 2m$. Assume $\nu > 0$ and $0 < \alpha < 1$. Then if a permutation $\sigma \in {}_{,2m}$ is chosen uniformly at random,

$$\Pr(d_\nu(\mu_1(\vec{x}, \sigma), \mu_2(\vec{x}, \sigma)) > \alpha) \leq 2e^{-2\alpha^2\nu m/M}$$

2.

$$\Pr\left(d_\nu\left(\frac{1}{n}\sum_{i=1}^n Z_i, \mu\right) > \alpha\right) \leq 2e^{\frac{-18\alpha^2\nu n}{(3+\alpha)^2 M}} < 2e^{-\frac{9}{8}\alpha^2\nu n/M} < 2e^{-\alpha^2\nu n/M}.$$

Proof. Let $Y_i = Z_i - \mu$, $1 \leq i \leq n$. Then

$$\begin{aligned} & \Pr\left(d_\nu\left(\frac{1}{n}\sum_{i=1}^n Z_i, \mu\right) > \alpha\right) \\ &= \Pr\left(\frac{|\sum_{i=1}^n Z_i - \mu n|}{\nu n + \sum_{i=1}^n Z_i + \mu n} > \alpha\right) \\ &\leq \Pr\left(\left|\sum_{i=1}^n Y_i\right| > \alpha n(\nu + \mu)\right), \end{aligned}$$

since $\sum_{i=1}^n Z_i \geq 0$.

To obtain the first bound, note that by Chebyshev's inequality,

$$\Pr\left(\left|\sum_{i=1}^n Y_i\right| > \alpha n(\nu + \mu)\right) \leq \sigma^2/(\alpha n(\nu + \mu))^2,$$

where $\sigma^2 = n \mathbf{Var}(Y_i)$. Since $0 \leq Z_i \leq M$, $\mathbf{Var}(Z_i) = \mathbf{Var}(Y_i) \leq \mu(M - \mu)$. Hence

$$\Pr\left(\left|\sum_{i=1}^n Y_i\right| > \alpha n(\nu + \mu)\right) \leq \frac{\mu(M - \mu)}{\alpha^2 n(\nu + \mu)^2}.$$

It is easily verified that the maximum value of this expression occurs at $\mu = (\nu M)/(2\nu + M)$, and that this gives an upper bound of

$$\frac{M^2}{4\alpha^2\nu n(M + \nu)}.$$

To obtain the second bound, we apply Bernstein's inequality (see e.g. [Pol84], Page 192), which states that

$$\Pr\left(\left|\sum_{i=1}^n Y_i\right| > \eta\right) \leq 2e^{-\eta^2/2(n\mathbf{Var}(Y_i) + \frac{1}{3}B\eta)}$$

for any zero mean i.i.d. random variables Y_1, \dots, Y_n bounded in absolute value by B . Substituting $\eta = \alpha n(\nu + \mu)$ and upper bounds $B \leq M$ and $\mathbf{Var}(Y_i) \leq \mu M$, this gives a bound of

$$2e^{-\alpha^2 n^2 (\nu + \mu)^2 / 2(n\mu M + \frac{1}{3}M\alpha n(\nu + \mu))} = 2e^{\frac{-3\alpha^2 n(\nu + \mu)^2}{2M(\alpha\nu + (3 + \alpha)\mu)}}.$$

Since $(\nu + \mu)^2/(\alpha\nu + (3 + \alpha)\mu)$ is minimized at $\mu = \frac{3-\alpha}{3+\alpha}\nu$, the latter expression is bounded by substituting this value of μ . This gives the first bound of part (2). The second bound follows from the fact that $\alpha < 1$. \square

$\mathcal{B}(T)$ of Borel sets on T is then the restriction to T of the σ -algebra \mathcal{B} of Borel sets on \bar{T} , i.e.

$$\mathcal{B}(T) = \{B \cap T : B \in \mathcal{B}\}.$$

Finally, if \mathcal{A} is any σ -algebra on Z and \mathcal{C} is any σ -algebra on T then $\mathcal{A} \times \mathcal{C}$ denotes the smallest σ algebra on $Z \times T$ that contains $\{A \times C : A \in \mathcal{A}, C \in \mathcal{C}\}$. We are now ready for our main definition.

Definition 11 *We say that the class \mathbf{F} is permissible if it can be indexed by a set T such that*

1. *T is a Borel subspace of a compact metric space \bar{T} and*
2. *the function $f : Z \times T \rightarrow \mathfrak{R}$ that indexes \mathbf{F} by T is measurable with respect to the σ -algebra $\mathcal{A} \times \mathcal{B}(T)$.*

Most uncountable classes of functions that come up in practice can be indexed by a finite number of real parameters (i.e. with $T = \mathfrak{R}^n$ for some $n \geq 1$) in such a way that condition (2) is satisfied. Condition (1) is satisfied as well in this case, since we can take \bar{T} to be the one-point compactification of T , obtained by adding a *point at infinity* to T (see e.g. [Sim63]).

Results given in Pollard ([Pol84], Appendix C) imply that the sets used in Lemmas 12 and 13 are measurable when \mathbf{F} is permissible. He also shows that the packing numbers $\mathcal{M}(\epsilon, \mathbf{F}|_Z, d_{L^1})$ are measurable functions of $\vec{z} \in Z^m$ for any $m \geq 1$. Since Theorem 12 relates these closely to the covering numbers $\mathcal{N}(\epsilon, \mathbf{F}|_Z, d_{L^1})$, this allows us to further formalize our usage of random covering numbers. A more formal treatment would either replace the covering numbers with the packing numbers in our upper bounds, or reword probabilistic bounds on the covering numbers to use outer measure arguments.

10.3 Measuring the accuracy of empirical estimates with the d_ν metric

In this section we give two bounds on the probability of large deviation of empirical estimates from true means, as measured by the d_ν metric. One is derived from Chebyshev's inequality and the other from Bernstein's inequality. The first bound is better for estimates obtained from small samples, the latter for estimates obtained from larger samples.

Lemma 9 *Let Z_1, \dots, Z_n be i.i.d. random variables with range $0 \leq Z_i \leq M$ and $\mathbf{E}(Z_i) = \mu$, $1 \leq i \leq n$. Assume $\nu > 0$ and $0 < \alpha < 1$. Then*

- 1.

$$\Pr \left(d_\nu \left(\frac{1}{n} \sum_{i=1}^n Z_i, \mu \right) > \alpha \right) \leq \frac{M^2}{4\alpha^2 \nu n (M + \nu)} < \frac{M}{4\alpha^2 \nu n}.$$

The following inequalities are easily verified (see e.g. [KT61]):

Theorem 12 *If T is a totally bounded subset of the (pseudo) metric space (S, ρ) then for any $\epsilon > 0$,*

$$\mathcal{M}(2\epsilon, T, \rho) \leq \mathcal{N}(\epsilon, T, \rho) \leq \mathcal{M}(\epsilon, T, \rho).$$

Hence both these measures of boundedness, by covering number and by packing number, are equivalent to within a factor of 2 of ϵ . Following [KT61] we define the *upper metric dimension* of a (pseudo) metric space (S, ρ) by

$$\overline{\mathbf{dim}}(S) = \mathbf{limsup}_{\epsilon \rightarrow 0} \frac{\log \mathcal{N}(\epsilon, S, \rho)}{\log(1/\epsilon)}.$$

The *lower metric dimension*, denoted by \mathbf{dim} , of a (pseudo) metric space (S, ρ) is defined similarly using \mathbf{liminf} . When $\overline{\mathbf{dim}}(S) = \mathbf{dim}(S)$, then this quantity is denoted $\mathbf{dim}(S)$, and referred to simply as the *metric dimension* of (S, ρ) . This quantity has also been called the *fractal dimension* [Far82] and the *capacity dimension* [FOY83]. A very lucid and intuitive treatment is given in [Man82].

10.2 Permissible classes of functions

In order to obtain the uniform convergence results given in Theorem 2, certain measurability assumptions have to be made concerning the class of functions \mathbf{F} when this class is uncountable. These we have indicated by saying that \mathbf{F} must be a *permissible* class [Pol84]. Here we give a definition of permissible that is a special case of that given by Pollard. This definition will be suitable for our purposes; we refer the reader to [Pol84] and [Dud84] for a more general treatment. See exercise 10, page 39 of [Pol84] for an indication of the kind of problems that can come up with non-permissible classes.

Throughout the paper we have assumed that \mathbf{F} is a class of real-valued functions on a set Z , and that P is a measure defined on some σ -algebra \mathcal{A} of subsets of Z such that each function in \mathbf{F} is measurable. We will need further conditions on \mathbf{F} when it is uncountable. Let us say that the class \mathbf{F} is *indexed* by the set T if

$$\mathbf{F} = \{f(\cdot, t) : t \in T\},$$

where f is a real-valued function on $Z \times T$ and $f(\cdot, t)$ denotes the real-valued function on Z obtained from f by fixing the second parameter to t . We will say that the function f *indexes \mathbf{F} by T* .

We will need some structure on T as well. If T is contained in a topological space \bar{T} then we will say that T is a *Borel subspace* if T is a Borel set with respect to the topology on \bar{T} , i.e. if T is in the smallest σ -algebra on \bar{T} containing the open sets. The σ -algebra

the problem of overfitting, including distribution specific bounds on sample complexity (Theorem 2 is actually distribution specific, since the random covering numbers are distribution specific, yet we only apply it here in a distribution independent setting), decision rule spaces with infinite pseudo and metric dimensions (these include various classes of “smooth” functions and their relatives, see [Dud84], Chapter 7 and [Qui89]) and non i.i.d. sources of examples (see [Whi90a,ND90]). Despite these shortcomings, we feel that the theory we give here provides useful insights into the nature of the problem of overfitting in learning, and because of its generality will be a useful starting point for further research in this area.

9 Acknowledgements

I would like to thank Dana Angluin, David Pollard and Phil Long for their careful criticisms of an earlier draft of this paper, and their numerous suggestions for improvements. I also thank Naoki Abe, Anselm Blumer, Richard Dudley, and Michael Kearns for helpful comments on earlier drafts. I would also like to thank Ron Rivest, David Rumelhart, Andrzej Ehrenfeucht and Nick Littlestone for stimulating discussions on these topics.

10 Appendix

10.1 Metric spaces, covering numbers and metric dimension

A *pseudo metric* on a set S is a function ρ from $S \times S$ into \mathbb{R}^+ such that for all $x, y, z \in S$, $x = y \Rightarrow \rho(x, y) = 0$, $\rho(x, y) = \rho(y, x)$ (symmetry), and $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ (triangle inequality). If in addition $\rho(x, y) = 0 \Rightarrow x = y$, then ρ is a *metric*. (S, ρ) is a (*pseudo*) *metric space*. (S, ρ) is *complete* if every Cauchy sequence of points in S converges to a point in S ; (S, ρ) is *separable* if it contains a countable dense subset, i.e. a countable subset A such that for every $x \in X$ and $\epsilon > 0$ there exists $a \in A$ with $\rho(x, a) < \epsilon$. If $\rho(x, y) = 1 \Leftrightarrow x \neq y$ then ρ is called the *discrete metric*.

The *diameter* of a set $T \subseteq S$ is $\sup\{\rho(x, y) : x, y \in T\}$. If the diameter of T is finite then we say that T is *bounded*. For any $\epsilon > 0$, an ϵ -*cover* for T is a finite set $N \subseteq S$ (not necessarily contained in T) such that for all $x \in T$ there is a $y \in N$ with $\rho(x, y) \leq \epsilon$. If T has a (finite) ϵ -cover for all $\epsilon > 0$ then T is *totally bounded*. (Note that this implies that (T, ρ) is separable and bounded.) In this case the function $\mathcal{N}(\epsilon, T, \rho)$ denotes the size of the smallest ϵ -cover for T (w.r.t. the space S and the (pseudo) metric ρ). We refer to $\mathcal{N}(\epsilon, T, \rho)$ as a *covering number*. A set $R \subseteq T$ is ϵ -*separated* if for all distinct $x, y \in R$, $\rho(x, y) > \epsilon$. We denote by $\mathcal{M}(\epsilon, T, \rho)$ the size of the largest ϵ -separated subset of T . We refer to $\mathcal{M}(\epsilon, T, \rho)$ as a *packing number*. The third argument to \mathcal{N} and \mathcal{M} will be omitted when the metric ρ is clear from the context.

similar to that obtained in [BH89].

Since W appears to be the dominant factor in these bounds, apart from the accuracy parameters α and ν , these bounds support the conventional wisdom that the training set size should be primarily related to the number of adjustable parameters in the net. They also support the notion that this relationship between appropriate training size and the number of parameters is nearly linear, at least in the worst case. Further work is needed to sharpen these relationships (see e.g. the lower bounds obtained in [BH89]).

8 Conclusion

We have extended the PAC learning model to a more general decision theoretic framework so that it addresses many of the concerns raised by machine learning practitioners, and also introduced a number of new theoretical tools. Here we concentrate on applications of the extended model to the problem of obtaining upper bounds on sufficient training sample size. Further work will be required to obtain lower bounds on sample size needed, and to determine the computational complexity of finding decision rules with near minimal empirical risk. Some promising results along these lines are given in [KS90]. However, even granting that such results can be obtained, the extended model still has a number of shortcomings in its present form. Some of these can be easily remedied, others may be more problematic.

First, we define the model only for a fixed decision rule space \mathcal{H} . The model should be extended to learning problems on a sequence of decision rule spaces $\{\mathcal{H}_n : n \geq 1\}$, where \mathcal{H}_n is a decision rule space on an n -attribute domain X_n (e.g. $[0, 1]^n$), and to families of decision rule spaces of different “complexities” on a fixed domain [KLPV87] [BEHW89] [HKLW91], so that tradeoffs between decision rule complexity and empirical risk can be addressed. The former extension is easy, the latter more involved. One approach to the latter problem is via Vapnik’s principle of structural risk minimization [Vap89] (see also [Dev88]). Other approaches include the MDL (see e.g. [BC90]), regularization (see e.g. [PG89]), and more general Bayesian methods (see e.g. [Ber85]).

Second, the constants in the upper bounds are still too large to give sample size estimates that are useful in practice. It may be difficult to improve them to the point where the results are directly usable in applied work. Thus even with matching asymptotic lower bounds, practitioners may still need to rely at least in part on empirically derived sample size bounds. It is possible that the Bayesian viewpoint may yield better tools for calculating sample complexities. Support for this belief is given in [CB90, CB, HKS91, OH91a]. However, necessary sample size estimates for decision rule spaces as general as those studied from the minimax perspective using uniform convergence have not yet been tackled from the Bayesian perspective.

Finally, many other issues would need to be considered in a complete treatment of

Despite the uncertainty about the need for the Lipschitz bounds, the result does give some indication of the maximum training sample size that will be needed for many popular network configurations. For example, if the squashing function is chosen as $\sigma(x) = 1/(1 + e^{-x/T})$ for some *temperature* $T > 0$ then it can be shown that the Lipschitz bound s for σ is $1/4T$. When the modifier $\mu \equiv 0$, then $r = 0$. Thus in this case the term $d \log(s(\beta W_{max} + r))$ in the bound of Corollary 3 becomes $d \log(\beta W_{max}/T)$. If the maximum weight β , the temperature T and the depth d are constants, along with $c_2 - c_1$, M , and c_l , then the asymptotic bound of the theorem becomes

$$O\left(\frac{1}{\alpha^{2\nu}} \left(W \log \frac{W_{max}}{\alpha\nu} + \log \frac{1}{\delta}\right)\right),$$

which is similar to the bound obtained in [BH89].

It should also be noted that Corollary 3 does have the feature that no Lipschitz bounds are required on the computation units at depth one. Thus if all computation units are at depth one, i.e. there are no hidden units, then no Lipschitz units are required at all. If the architecture has only one layer of hidden units at depth one and a single output unit at depth two, as is quite common, then Lipschitz bounds are required only on the output unit. This means that the weights and biases associated with the hidden units do not need to be bounded in order to get the rates of uniform convergence given by Corollary 3, as they would, for example, if the methods given in [Whi90a] were used to obtain a result of this type.

For an example of the above, consider networks that implement generalized radial basis¹⁴ functions, as described in [PG89]. These networks have one layer of hidden units at depth 1 and one output unit at depth 2. The structure of the hidden units is as described in the example above: the input transformers are identity functions, the modifier is $\sum_{i=1}^n x_i^2$ and the squashing function is usually a smooth decreasing function. The output unit simply computes a weighted sum, so for this unit the modifier is the 0 function and the squashing function is the identity. Since this is the only unit at depth 2 and above, we require a Lipschitz bound only for this unit. If β is a bound on the maximum weight coming into the output unit, and W_{max} is the number of units in the hidden layer, then the term $d \log(s(\beta W_{max} + r))$ in the above bound becomes $\log(\beta W_{max})$. Again, fixing β , $c_2 - c_1$, M , and c_l gives the same sample size bound,

$$O\left(\frac{1}{\alpha^{2\nu}} \left(W \log \frac{W_{max}}{\alpha\nu} + \log \frac{1}{\delta}\right)\right),$$

¹⁴The computation units in the network of radial basis functions described here are quite primitive in that they have no adjustable multiplicative parameter included in their basic radial distance calculation. Such parameters would be needed to do any reasonable type of kernel based density estimation (see e.g. [DH73]). These parameters can be simulated by inserting another layer of computation units between the inputs and the layer described here. Alternately, the analysis can also be done directly for adjustable kernel units. This cleaner approach is detailed in [Pol86].

transformers are identity functions, the global modifier has Lipschitz bound at most r , and the squashing function has Lipschitz bound at most s , where $s(\beta W_{max} + r) \geq 1$. Then for any $0 < \delta < 1$, the probability in (1.) is less than δ for sample size

$$m = O\left(\frac{M}{\alpha^2\nu} \left(W \left(\log \frac{c_l(c_2 - c_1)}{\alpha\nu} + d \log(s(\beta W_{max} + r))\right) + \log \frac{1}{\delta}\right)\right).$$

Proof. Let $\epsilon = \alpha\nu/8 \leq c_2 - c_1$. It can be verified that $l_{\mathcal{H}}$ is permissible for the decision rule space \mathcal{H} . Hence, using Theorem 9 and Theorem 11,

$$\begin{aligned} \Pr(\exists f \in \mathcal{H} : d_{\nu}(\hat{\mathbf{r}}_{f,l}(\vec{z}), \mathbf{r}_{f,l}(P)) > \alpha) &\leq 4\mathcal{C}(\alpha\nu/8, \mathcal{H}, \rho_l) e^{-\alpha^2\nu m/8M} \\ &\leq 4\mathcal{C}(\alpha\nu/8c_l, \mathcal{H}, d_{L^1}) e^{-\alpha^2\nu m/8M} \\ &\leq 4 \left(\frac{c_l 16\epsilon(c_2 - c_1) d \prod_{l=2}^d b_l}{\alpha\nu}\right)^{2W} e^{-\alpha^2\nu m/8M}. \end{aligned}$$

For the second bound, it is readily verified that the L^1 Lipschitz bound for a linear function defined by W_{max} weights and a bias is no more than W_{max} times the largest absolute value of any weight. Furthermore, the Lipschitz bound for the sum of two functions is no more than the sum of the Lipschitz bounds on the individual functions, and the Lipschitz bound for the composition of two functions is no more than the product of the Lipschitz bounds on the individual functions. Thus if the input transformers are identity functions, the global modifier and squashing function have Lipschitz bounds r and s respectively and no weight is allowed to have absolute value greater than β , then the Lipschitz bound for a computation unit is at most $s(\beta W_{max} + r)$. If this holds for all units at depth 2 and above, may take $b_j = s(\beta W_{max} + r)$ for all $j \geq 2$ in the first bound. Solving for m , this gives the order-of-magnitude estimate of the second bound. \square

We give the constants in the upper bound of part (1.) the above theorem only to show that they are not outlandishly large. We do not mean to suggest that the bound is tight. At present we cannot even verify that the asymptotic bound of part (2.) is tight. In particular, we cannot show that the dependence on the Lipschitz bounds is necessary. Evidence that it may not be necessary comes from the analysis of the case where the squashing function σ is a sharp threshold function, i.e. $\sigma(x) = \text{sign}(x)$. Corollary 3 does not apply in this case, because the jump in σ prevents us from obtaining a Lipschitz bound on the computation units. As we let a smooth σ approach the sign function, its slope increases without limit, and thus the bound given in Corollary 3 degenerates. Nevertheless, using the techniques in [BH89] it can be shown that results similar to Corollary 3 hold in this case, except that no Lipschitz bounds are required, and a bound on the sample size is

$$O\left(\frac{1}{\alpha^2\nu} \left(W \log \frac{N}{\alpha\nu} + \log \frac{1}{\delta}\right)\right),$$

where N is the total number of computation units in the net. Details are given in Theorem 13 in the Appendix.

dimensional vector space of real-valued functions, summed with a fixed modifier and then composed with a non-increasing or non-decreasing squashing function. In the latter case, the dimension N of this vector space is the number of free parameters associated with the corresponding computation unit, i.e. the number of weights plus one (for the adjustable bias). Hence by Theorems 4 and 5 in section 4, the pseudo dimension $\mathbf{dim}_{\mathbf{P}}(\mathcal{F}) \leq N$. Thus by Theorem 10 above,

$$\mathcal{C}(\epsilon_j, \mathcal{F}, d_{L^1}) \leq 2 \left(\frac{2\epsilon(c_2 - c_1)}{\epsilon_j} \ln \frac{2\epsilon(c_2 - c_1)}{\epsilon_j} \right)^N \leq \left(\frac{2\epsilon(c_2 - c_1)}{\epsilon_j} \right)^{2N},$$

since $2 \ln x < x$ and $N \geq 1$. Since the capacity of a class with only one function is 1, it follows from Lemma 7 part (1) that

$$\mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1}) \leq \left(\frac{2\epsilon(c_2 - c_1)}{\epsilon_j} \right)^{2W_j},$$

where W_j is the total number of weights and biases of all computation nodes at depth j . Multiplying these bounds over all j , it follows that

$$\begin{aligned} \mathcal{C}(\epsilon, \mathcal{H}, d_{L^1}) &\leq \prod_{j=1}^d \left(\frac{2\epsilon(c_2 - c_1)}{\epsilon_j} \right)^{2W_j} \\ &= \prod_{j=1}^d \left(\frac{2\epsilon(c_2 - c_1) d \prod_{l=j+1}^d b_l}{\epsilon} \right)^{2W_j} \\ &\leq \left(\frac{2\epsilon(c_2 - c_1) d \prod_{l=2}^d b_l}{\epsilon} \right)^{2W}. \end{aligned}$$

□

Corollary 3 *Let n, k, \mathcal{H}, W, d and b_2, \dots, b_d be as in the previous theorem. Let X be the instance space $[c_1, c_2]^n$, A be the decision space $[c_1, c_2]^k$, and Y be any outcome space. Let $l : Y \times A \rightarrow [0, M]$ be a loss function and c_l be a constant such that $\rho_l(\vec{a}, \vec{b}) \leq c_l d_{L^1}(\vec{a}, \vec{b})$ for all $\vec{a}, \vec{b} \in A$. Let $m \geq 1$, $0 < \nu \leq 8(c_2 - c_1)$, $0 < \alpha < 1$, and P be any probability distribution on $Z = X \times Y$. Let \vec{z} be generated by m independent random draws from Z according to P .*

1. *Then*

$$\Pr(\exists f \in \mathcal{H} : d_\nu(\hat{\mathbf{r}}_{f,l}(\vec{z}), \mathbf{r}_{f,l}(P)) > \alpha) \leq 4 \left(\frac{c_l 16\epsilon(c_2 - c_1) d \prod_{l=2}^d b_l}{\alpha \nu} \right)^{2W} e^{-\alpha^2 \nu m / 8M}.$$

2. *Assume that for each computation unit at depth 2 and above the number of weights is at most W_{max} , no weight is allowed to have absolute value greater than β , the input*

biases in \mathcal{A} . Assume $b_j \geq 1$ for $2 \leq j \leq d$, and let \mathcal{H} be all functions from $[c_0, c_1]^n$ into $[c_0, c_1]^k$ representable on \mathcal{A} by setting the adjustable weights and biases such that for all j , $2 \leq j \leq d$, the average of the Lipschitz bounds of the functions computed by computation units at depth j is at most b_j . Then for all $0 < \epsilon \leq c_2 - c_1$,

$$\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1}) \leq \left(\frac{2\epsilon(c_2 - c_1)d \prod_{l=2}^d b_l}{\epsilon} \right)^{2W}.$$

□

Proof. For each j , $0 \leq j \leq d$, let n_j be the number of units at depth j in the architecture \mathcal{A} . For each j , $0 \leq j \leq d-1$, let $l_j = \sum_{i=0}^j n_i$, and let $l_d = n_d = k$. For each j , $1 \leq j \leq d+1$, let $X_j = [c_1, c_2]^{l_{j-1}}$. Then for each j , $1 \leq j \leq d$, we can define the family \mathcal{H}_j of functions from X_j into X_{j+1} in the following manner.

First assume $j < d$. Let u_1, \dots, u_{n_j} be an enumeration of the computation units at depth j and f_1, \dots, f_{n_j} be functions such that f_i can be represented by u_i , $1 \leq i \leq n_j$, and the average Lipschitz bound on the f_i s is at most b_j . Let h_j be the free product of f_1, \dots, f_{n_j} and l_{j-1} copies of the identity function on $[c_1, c_2]$. Thus $h_j : X_j \rightarrow X_{j+1}$. The function h_j represents a mapping from the sequence of all activations of units at depth at most $j-1$ to the sequence of all activations of units at depth at most j , where the activations at depth at most $j-1$ are unaltered, and the new activations, i.e. those at depth j , are calculated by f_1, \dots, f_{n_j} . The family \mathcal{H}_j consists of all functions h_j obtained in this manner, by varying the weights and biases in the units u_1, \dots, u_{n_j} at depth j in such a manner that the Lipschitz constraint is satisfied.

When $j = d$, no subsequent calculations will be performed so we no longer need to preserve the activations of shallower units. Hence, we omit the identity function components in each $h_d \in \mathcal{H}_d$. Otherwise the definition of \mathcal{H}_d is the same as that for \mathcal{H}_j , where $j < d$.

It is clear that the class \mathcal{H} in the statement of the theorem can be represented as the class of compositions of functions from classes $\mathcal{H}_1, \dots, \mathcal{H}_d$. Since the identity function has Lipschitz bound $1 \leq b_j$, the average Lipschitz bound on the components of each function $h_j \in \mathcal{H}_j$ is at most b_j . It is easily verified that a free product function is Lipschitz bounded by the average of the Lipschitz bounds on its component functions. Hence by assumption, b_j is a uniform Lipschitz bound on \mathcal{H}_j , $2 \leq j \leq d$. For each j , $1 \leq j \leq d$, let $a_j = \prod_{l=j+1}^d b_l$ and $\epsilon_j = \frac{\epsilon}{da_j}$. Since $\epsilon < c_2 - c_1$ and $a_j \geq 1$, $\epsilon_j \leq c_2 - c_1$. Let ρ_j be the d_{L^1} metric on X_j , $1 \leq j \leq d+1$. Then by Lemma 8, part (1),

$$\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1}) \leq \prod_{j=1}^d \mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1}).$$

For each j , \mathcal{H}_j is contained in the free product of l_j function classes. Each class \mathcal{F} in this product is either the trivial class containing only the identity function, or is a finite

computed by a given computation unit with n incoming edges has *Lipschitz bound* b if for any $\vec{x}, \vec{y} \in [c_0, c_1]^n$, $|f(\vec{x}) - f(\vec{y})| \leq bd_{L^1}(\vec{x}, \vec{y})$.

We give a few examples to illustrate the flexibility of this model at the level of the individual computation unit. First assume that $k = n$, i.e. the number of input transformers is the same as the number of inputs, and that each input transformer simply extracts a component of the input, i.e. $\phi_j(\vec{x}) = x_j$, $1 \leq j \leq n$. In this case, which is the standard case for most neural net research, the overall input transformation is just the identity map and can be ignored. In this standard case, if the global modifier $\mu = 0$ we get what is known as a quasi-linear unit [RM86]:

$$f(\vec{x}) = \sigma \left(\theta + \sum_{j=1}^k w_j x_j \right).$$

In the standard case, if $\mu(\vec{x}) = \sum_{j=1}^n x_j^2$ we get a unit that computes a function of the form

$$f(\vec{x}) = \sigma \left(\theta' + \sum_{j=1}^n (x_j - a_j)^2 \right),$$

where $a_j = -w_j/2$ and $\theta' = \theta - \sum_{j=1}^n a_j^2$. This is similar to what is called a *radial basis* unit in the neural net literature [PG89,MD89].

Now assume that $k = n$ but the input transformers take logs of the components of the inputs, i.e. $\phi_j(\vec{x}) = \log x_j$. (Here we assume $c_0 > 0$.) Let $\mu = 0$ and change the squashing function σ to σ' , where $\sigma'(x) = \sigma(e^x)$. Then

$$f(\vec{x}) = \sigma' \left(\theta + \sum_{j=1}^n w_j \log x_j \right) = \sigma \left(e^\theta \prod_{j=1}^n x_j^{w_j} \right),$$

giving what is commonly known as a *product unit* [DR89].

We define a feedforward *architecture* as a feedforward net with unspecified weights and biases, i.e. each computation unit has a fixed global modifier, a fixed squashing function and fixed input transformers, but it has variable weights and a variable bias. We say a unit is *at depth* j in an architecture if the longest (directed) path from an input unit to that unit has j edges. Thus all input units are at depth 0, all computation units that have incoming edges only from input units are at depth 1, all computation units that have incoming edges only from input units and computation units at depth 1 are at depth 2, etc. The *depth of the architecture* is the depth of the deepest unit in it.

We can bound the capacity of the decision rule space represented by a feedforward architecture as follows.

Theorem 11 *Let \mathcal{A} be a feedforward architecture as above with $n \geq 1$ input units, $k \geq 1$ output units, and depth $d \geq 1$. Let W be the total number of adjustable weights and*

2. $\overline{\mathbf{dim}}(\mathcal{H}) \leq \mathbf{dim}_{\mathbf{P}}(\mathcal{H})$.

Proof. Let P be any probability measure on X . Then by Theorems 12 and 6 in sections 10.1 and 4,

$$\mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P)}) \leq \mathcal{M}(\epsilon, \mathcal{H}, d_{L^1(P)}) < 2 \left(\frac{2eM}{\epsilon} \ln \frac{2eM}{\epsilon} \right)^d.$$

(Theorem 6 is applied with $Z = X$ and $\mathbf{F} = \mathcal{H}$.) This gives (1.), and (2.) follows easily from (1.). \square

In the general case, where $A \subset \mathfrak{R}^k$ for $k > 1$, we can apply the methods from the previous section, in addition to the pseudo dimension methods from section 4, to obtain bounds on $\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1})$. We illustrate this for the case when \mathcal{H} is the class of decision rules represented by a feedforward neural network.

A feedforward neural network is defined as a directed acyclic graph in which the incoming edges to each node (or *unit*) are ordered and each incoming edge can carry a real number representing the *activation* on that edge. We will assume that all activations are restricted to the interval $[c_0, c_1]$ for some constants $c_0 < c_1$. The units are divided into *input units*, which have no incoming edges from other units and serve as input ports for the network (their activations are determined by these external inputs), and *computation units*, which have incoming edges from other units and compute an activation based on the activations on these incoming edges. After an activation has been determined, this activation is placed on the outgoing edges of the unit. Computation units with no outgoing edges are called *output units* and serve as output ports for the network. Computation units that are not output units are called *hidden units*. The network as a whole computes a function that maps from vectors of activation values in its input units to vectors of activation values in its output units by composing the functions computed by its computation units in the obvious way.

The action of a computation unit with n incoming edges can be specified by a function f from $[c_0, c_1]^n$ into $[c_0, c_1]$, where $f(\vec{x})$ is the resulting activation of the unit when the activations of its incoming edges are given by the vector $\vec{x} \in [c_0, c_1]^n$. In the nets we consider, the function f is defined by

$$f(\vec{x}) = \sigma \left(\mu(\phi_1(\vec{x}), \dots, \phi_k(\vec{x})) + \theta + \sum_{j=1}^k w_j \phi_j(\vec{x}) \right),$$

where the w_j 's are adjustable real *weights*, θ is an adjustable real *bias*, ϕ_1, \dots, ϕ_k are fixed real-valued functions which we call the *input transformers*, $\mu : \mathfrak{R}^k \rightarrow \mathfrak{R}$ is a fixed function which we call the *global modifier*, and $\sigma : \mathfrak{R} \rightarrow [c_0, c_1]$ is a fixed non-increasing or non-decreasing function which we call the *squashing function*. Different units can have different modifiers, transformers and squashing functions. We say that the function f

Note that the above trick does not apply to the mean squared loss $l(\vec{y}, \vec{a}) = \frac{1}{k} \sum_{i=1}^k (y_i - a_i)^2$ since this loss does not satisfy the triangle inequality. However, in this case it is easy to show by direct calculation that if the outcome space Y is bounded, e.g. $Y \subset [0, M]^k$, then $\rho_l(\vec{a}, \vec{b}) \leq 2M d_{L^1}(\vec{a}, \vec{b})$, and hence we may take $c_l = 2M$.

For our final example, consider the case when $Y = \{0, 1\}^k$, $A \subset [0, 1]^k$ and l is the cross entropy loss

$$l(\vec{y}, \vec{a}) = - \sum_{i=1}^k (y_i \ln a_i + (1 - y_i) \ln(1 - a_i)).$$

As discussed in section 1.1.3, this is the log likelihood loss for the regression problem in which the action \vec{a} represents a vector of probabilities for independent Bernoulli variables, and the outcome \vec{y} gives the observed values of these variables. This loss is bounded if we restrict the probabilities in \vec{a} to be between B and $1 - B$ for some $0 < B \leq 1/2$. In this case

$$\begin{aligned} \rho_l(\vec{a}, \vec{b}) &= \sup_{\vec{y} \in Y} \left| \sum_{i=1}^k \left(y_i \ln \frac{b_i}{a_i} + (1 - y_i) \ln \frac{(1 - b_i)}{(1 - a_i)} \right) \right| \\ &\leq \sum_{i=1}^k \left| \ln \frac{b_i}{a_i} \right| + \sum_{i=1}^k \left| \ln \frac{(1 - b_i)}{(1 - a_i)} \right| \\ &\leq \frac{2}{B} \sum_{i=1}^k |a_i - b_i|. \end{aligned}$$

The latter inequality follows from the fact that for $x, y > 0$,

$$\left| \ln \frac{x}{y} \right| = \ln \frac{\max(x, y)}{\min(x, y)} \leq \frac{\max(x, y)}{\min(x, y)} - 1 = \frac{|x - y|}{\min(x, y)}.$$

Thus in this case we may take $c_l = 2k/B$.

We now turn to the task of obtaining an upper bound on the capacity $\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1})$ when the decision rules in \mathcal{H} map into a decision space $A \subset \mathfrak{R}^k$, and in particular, when these decision rules are represented by neural networks. When $k = 1$, i.e. the neural net has only one output, the decision rule space \mathcal{H} is a family of real valued functions and $d_{L^1}(a, b) = |a - b|$ for $a, b \in A$. In this case we can apply the methods and results of section 4. We must first find an upper bound on $\mathbf{dim}_{\mathbf{P}}(\mathcal{H})$, the pseudo dimension of \mathcal{H} . Then, when A is bounded, from the bound on $\mathbf{dim}_{\mathbf{P}}(\mathcal{H})$ we get a bound on the capacity $\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1})$ using Theorem 6. This also gives a bound on the metric dimension of \mathcal{H} .

Theorem 10 *Let \mathcal{H} be a family of functions from X into $A = [0, M]$. Assume $\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) = d$ for some $1 \leq d < \infty$.*

1. For all $0 < \epsilon \leq M$,

$$\mathcal{C}(\epsilon, \mathcal{H}, d_{L^1}) < 2 \left(\frac{2eM}{\epsilon} \ln \frac{2eM}{\epsilon} \right)^d.$$

$$\begin{aligned}
&= \sum_{j=1}^k \limsup_{\epsilon \rightarrow 0} \left(\frac{\log(\mathcal{C}(\epsilon/ka_j, \mathcal{H}_j, \rho_{j+1}))}{\log(1/\epsilon)} \right) \\
&= \sum_{j=1}^k \limsup_{\epsilon_j \rightarrow 0} \left(\frac{\log(\mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1}))}{\log(1/\epsilon_j) + \log(1/ka_j)} \right) \\
&= \sum_{j=1}^k \limsup_{\epsilon_j \rightarrow 0} \left(\frac{\log(\mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1}))}{\log(1/\epsilon_j)} \right) \\
&= \sum_{j=1}^k \overline{\dim}(\mathcal{H}_j).
\end{aligned}$$

□

7 Sample size bounds for learning with multi-layer neural nets

We now present some applications of the results of the previous section to learning with feedforward neural nets (see e.g. [RM86,PG89]). The decision rule space \mathcal{H} represented by a feedforward neural net consists of a family of functions from an instance space $X \subset \mathfrak{R}^n$ into a decision space $A \subset \mathfrak{R}^k$ for some $k, n \geq 1$. To apply Theorem 9 of the previous section, we will need to obtain an upper bound on the capacity $\mathcal{C}(\epsilon, \mathcal{H}, \rho_l)$ of such decision rule spaces for various loss functions l .

For many loss functions, the metric ρ_l on $A \subset \mathfrak{R}^k$ can be bounded in terms of the d_{L^1} metric, i.e. we can find a constant c_l such that for all $\vec{a} = (a_1, \dots, a_k)$ and $\vec{b} = (b_1, \dots, b_k)$ in A , $\rho_l(\vec{a}, \vec{b}) \leq c_l d_{L^1}(\vec{a}, \vec{b}) = \frac{c_l}{k} \sum_{i=1}^k |a_i - b_i|$. In this case it is clear that $\mathcal{C}(\epsilon, \mathcal{H}, \rho_l) \leq \mathcal{C}(\epsilon/c_l, \mathcal{H}, d_{L^1})$. Thus our problem is reduced to obtaining an upper bound on the capacity $\mathcal{C}(\epsilon/c_l, \mathcal{H}, d_{L^1})$.

We now give a few examples to illustrate this reduction. First consider the common case in which the outcome space Y is also contained in \mathfrak{R}^k , e.g. we receive explicit feedback on each coordinate of our action $\vec{a} \in A \subset \mathfrak{R}^k$. This occurs when each coordinate a_i of the action \vec{a} is a prediction of the corresponding coordinate of the outcome \vec{y} . Here the loss function l may itself be a metric on \mathfrak{R}^k which measures the distance between the predicted vector and the actual outcome vector. When l is a metric, we have for any actions $\vec{a}, \vec{b} \in A$

$$\rho_l(\vec{a}, \vec{b}) = \sup_{\vec{y} \in Y} |l(\vec{y}, \vec{a}) - l(\vec{y}, \vec{b})| \leq l(\vec{a}, \vec{b})$$

by the triangle inequality for l . Thus if the metric l is bounded with respect to d_{L^1} metric, i.e. $l(\vec{a}, \vec{b}) \leq c_l d_{L^1}(\vec{a}, \vec{b})$ for all $\vec{a}, \vec{b} \in A$, then we have $\rho_l(\vec{a}, \vec{b}) \leq c_l d_{L^1}(\vec{a}, \vec{b})$. For example, if $l(\vec{y}, \vec{a}) = d_{L^2}(\vec{y}, \vec{a}) = \frac{1}{k} \left(\sum_{i=1}^k (y_i - a_i)^2 \right)^{1/2}$ then we may take $c_l = 1$, and similarly for the other d_{L^q} metrics, for $q > 1$.

- $f_2 \in U_{f_1}$ such that $d_{L^1(P_{f_1, \rho_3})}(f_2, g_2) \leq \epsilon_2, \dots$, and
- $f_k \in U_{f_1, \dots, f_{k-1}}$ such that $d_{L^1(P_{f_{k-1} \circ f_{k-2} \circ \dots \circ f_1, \rho_{k+1}})}(f_k, g_k) \leq \epsilon_k$.

Let $f = f_k \circ f_{k-1} \circ \dots \circ f_1 \in V$. It suffices to show that $d_{L^1(P, \rho_{k+1})}(f, g) \leq \epsilon$. We prove that for all $h, 1 \leq h \leq k$,

$$d_{L^1(P, \rho_{h+1})}(f_h \circ \dots \circ f_1, g_h \circ \dots \circ g_1) \leq \sum_{j=1}^h \left(\prod_{l=j+1}^h b_l \right) \epsilon_j.$$

Since

$$\epsilon = \sum_{j=1}^k \left(\prod_{l=j+1}^k b_l \right) \epsilon_j,$$

part (1) of the result follows.

If $h = 1$ then the result follows directly from our definition of f . Otherwise

$$\begin{aligned} & d_{L^1(P, \rho_{h+1})}(f_h \circ \dots \circ f_1, g_h \circ \dots \circ g_1) \\ &= \int_{X_1} \rho_{h+1}(f_h \circ \dots \circ f_1(x), g_h \circ \dots \circ g_1(x)) dP(x) \\ &\leq \int_{X_1} \rho_{h+1}(g_h \circ f_{h-1} \circ \dots \circ f_1(x), g_h \circ g_{h-1} \circ \dots \circ g_1(x)) dP(x) \\ &\quad + \int_{X_1} \rho_{h+1}(f_h \circ f_{h-1} \circ \dots \circ f_1(x), g_h \circ f_{h-1} \circ \dots \circ f_1(x)) dP(x) \\ &\leq b_h \int_{X_1} \rho_h(f_{h-1} \circ \dots \circ f_1(x), g_{h-1} \circ \dots \circ g_1(x)) dP(x) \\ &\quad + \int_{X_h} \rho_{h+1}(f_h(y), g_h(y)) dP_{f_{h-1}, \dots, f_1}(y) \quad (\text{by Lip. assump.}) \\ &\leq b_h \sum_{j=1}^{h-1} \left(\prod_{l=j+1}^{h-1} b_l \right) \epsilon_j + \epsilon_h \quad (\text{by i.h. and def. of } f_h) \\ &= \sum_{j=1}^h \left(\prod_{l=j+1}^h b_l \right) \epsilon_j. \end{aligned}$$

To prove part (2), let $a_j = \prod_{l=j+1}^k b_l$ and set $\epsilon_j = \frac{\epsilon}{k a_j}$ for $1 \leq j \leq k$. By part (1), $\mathcal{C}(\epsilon, \mathcal{H}, \rho_{k+1}) \leq \prod_{j=1}^k \mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1})$. Thus

$$\begin{aligned} \overline{\dim}(\mathcal{H}) &= \limsup_{\epsilon \rightarrow 0} \left(\frac{\log(\mathcal{C}(\epsilon, \mathcal{H}, \rho_{k+1}))}{\log(1/\epsilon)} \right) \\ &\leq \limsup_{\epsilon \rightarrow 0} \left(\frac{\log(\prod_{j=1}^k \mathcal{C}(\epsilon/k a_j, \mathcal{H}_j, \rho_{j+1}))}{\log(1/\epsilon)} \right) \end{aligned}$$

Definition 10 Let f be a function from a metric space (X, ρ) into a metric space (Y, σ) . A Lipschitz bound on f is a real number $b > 0$ such that for all $x, y \in X$, $\sigma(f(x), f(y)) \leq b\rho(x, y)$. The Lipschitz bound on f is the smallest such b . If \mathcal{F} is a class of functions from (X, ρ) into (Y, σ) then b is a uniform Lipschitz bound on \mathcal{F} if b is a Lipschitz bound on f for all $f \in \mathcal{F}$.

Lemma 8 Let $(X_1, \rho_1), \dots, (X_{k+1}, \rho_{k+1})$ be metric spaces, where (X_j, ρ_j) is bounded, $2 \leq j \leq k$, and \mathcal{H}_j be a class of functions with $f : X_j \rightarrow X_{j+1}$ for all $f \in \mathcal{H}_j$, $1 \leq j \leq k$. Let b_j be a uniform Lipschitz bound on \mathcal{H}_j for all $2 \leq j \leq k$. Let \mathcal{H} denote the class of all functions from X_1 into X_{k+1} defined by compositions of functions in the \mathcal{H}_j 's, i.e.

$$\mathcal{H} = \{f_k \circ f_{k-1} \circ \dots \circ f_1 : f_j \in \mathcal{H}_j, 1 \leq j \leq k\}.$$

1. For any $\epsilon, \epsilon_1, \dots, \epsilon_k > 0$ such that

$$\epsilon = \sum_{j=1}^k \left(\prod_{l=j+1}^k b_l \right) \epsilon_j$$

we have

$$\mathcal{C}(\epsilon, \mathcal{H}, \rho_{k+1}) \leq \prod_{j=1}^k \mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1}).$$

2. $\overline{\dim}(\mathcal{H}) \leq \sum_{j=1}^k \overline{\dim}(\mathcal{H}_j)$, and similarly for $\underline{\dim}$ and \mathbf{dim} , when the latter is defined.

Proof. Fix a probability measure P on X_1 . We define a tree-structured family of covers for the \mathcal{H}_j 's by induction as follows. For the basis case, let U be a minimum-sized ϵ_1 -cover for \mathcal{H}_1 w.r.t. the measure P on X_1 , i.e. $|U| = \mathcal{N}(\epsilon_1, \mathcal{H}_1, d_{L^1(P, \rho_2)})$ and every function in \mathcal{H}_1 is $L^1(P, \rho_2)$ -approximated to within ϵ_1 by some function in U . Now for each j , $2 \leq j \leq k$, and for each sequence of functions f_1, \dots, f_{j-1} where $f_1 \in U$, $f_2 \in U_{f_1}$, $f_3 \in U_{f_1, f_2}$, ..., $f_{j-1} \in U_{f_1, \dots, f_{j-2}}$, let $U_{f_1, \dots, f_{j-1}}$ be a minimum-sized ϵ_j cover for \mathcal{H}_j w.r.t. the L^1 metric for the measure $P_{f_{j-1} \circ f_{j-2} \circ \dots \circ f_1}$ on X_j and the metric ρ_{j+1} on X_{j+1} .

Next we define a cover V for \mathcal{H} by composing functions in the covers for the \mathcal{H}_j 's. If $k = 1$ then $V = U$. Otherwise

$$V = \{f_k \circ f_{k-1} \circ \dots \circ f_1 : f_1 \in U, f_2 \in U_{f_1}, f_3 \in U_{f_1, f_2}, \dots, \text{ and } f_k \in U_{f_1, \dots, f_{k-1}}\}.$$

Since $\mathcal{N}(\epsilon_j, \mathcal{H}_j, d_{L^1(P_{f_{j-1} \circ f_{j-2} \circ \dots \circ f_1}, \rho_{j+1})}) \leq \mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1})$ for all $1 \leq j \leq k$ and all f_1, \dots, f_{j-1} , it is clear that $|V| \leq \prod_{j=1}^k \mathcal{C}(\epsilon_j, \mathcal{H}_j, \rho_{j+1})$. Hence it remains to show that V is an ϵ -cover for \mathcal{H} .

Suppose that $g = g_k \circ g_{k-1} \circ \dots \circ g_1 \in \mathcal{H}$. Find

- $f_1 \in U$ such that $d_{L^1(P, \rho_2)}(f_1, g_1) \leq \epsilon_1$,

Proof. We begin with the second inequality of (1). For each $1 \leq j \leq k$ let U_j be an ϵ -cover for \mathcal{H}_j . Let

$$U = \prod_{j=1}^k U_j = \{(f_1, \dots, f_k) : f_j \in U_j, 1 \leq j \leq k\}.$$

It suffices to show that U is an ϵ -cover for \mathcal{H} . Let $g = (g_1, \dots, g_k)$ be any function in \mathcal{H} . For each j , $1 \leq j \leq k$, find $f_j \in U_j$ such that $d_{L^1(P, \rho_j)}(f_j, g_j) \leq \epsilon$. Let $f = (f_1, \dots, f_k)$. Then

$$\begin{aligned} d_{L^1(P, \rho)}(f, g) &= \int_X \frac{1}{k} \sum_{j=1}^k \rho_j(f_j(x), g_j(x)) dP(x) \\ &= \frac{1}{k} \sum_{j=1}^k \int_X \rho_j(f_j(x), g_j(x)) dP(x) \\ &= \frac{1}{k} \sum_{j=1}^k d_{L^1(P, \rho_j)}(f_j, g_j) \\ &\leq \epsilon. \end{aligned}$$

Hence U is an ϵ -cover for \mathcal{H} .

The first inequality of (1) is verified similarly. For each $1 \leq j \leq k$ let V_j be an $k\epsilon$ -separated subset of \mathcal{H}_j . Let $V = \prod_{j=1}^k V_j$. Let $f = (f_1, \dots, f_k)$ and $g = (g_1, \dots, g_k)$ be distinct functions in V . Then

$$d_{L^1(P, \rho)}(f, g) = \frac{1}{k} \sum_{j=1}^k d_{L^1(P, \rho_j)}(f_j, g_j) > \epsilon.$$

Hence V is an ϵ -separated subset of \mathcal{H} . It follows that $\prod_{j=1}^k \mathcal{M}(k\epsilon, \mathcal{H}_j, d_{L^1(P, \rho_j)}) \leq \mathcal{M}(\epsilon, \mathcal{H}, d_{L^1(P, \rho)})$. The first inequality of (1) then follows using Theorem 12.

From (1) we have

$$\prod_{j=1}^k \mathcal{C}(2k\epsilon, \mathcal{H}_j, \rho_j) \leq \mathcal{C}(\epsilon, \mathcal{H}, \rho) \leq \prod_{j=1}^k \mathcal{C}(\epsilon, \mathcal{H}_j, \rho_j).$$

Part (2) follows easily from this. \square

Definition 9 Let P be a probability measure on X and f be a measurable function from X into Y . Then P_f denotes the probability measure on Y induced by f , i.e.

$$P_f(S) = P(f^{-1}(S)) \text{ for all measurable } S \subset Y.$$

Theorem 9 Assume that the decision rule space \mathcal{H} and the loss function l are such that $l_{\mathcal{H}}$ is permissible. Let P be any probability distribution on $Z = X \times Y$. Assume $m \geq 1$, $\nu > 0$ and $0 < \alpha < 1$. Let \vec{z} be generated by m independent draws from Z according to P . Then

$$\Pr(\exists h \in \mathcal{H} : d_{\nu}(\hat{\mathbf{r}}_{h,l}(\vec{z}), \mathbf{r}_{h,l}(P)) > \alpha) \leq 4\mathcal{C}(\alpha\nu/8, \mathcal{H}, \rho_l)e^{-\alpha^2\nu m/8M}.$$

Proof. Let $\mathbf{F} = l_{\mathcal{H}}$. For any sequence \vec{z} of points in Z there is a trivial isometry between $(\mathbf{F}|_{\vec{z}}, d_{L^1})$ and $(\mathbf{F}, d_{L^1(P_{\vec{z}})})$, where $P_{\vec{z}}$ is the empirical measure induced by \vec{z} , in which each set has measure equal to the fraction of the points in \vec{z} it contains. Thus by Lemma 6 above, we have

$$\mathcal{N}(\epsilon, \mathbf{F}|_{\vec{z}}, d_{L^1}) = \mathcal{N}(\epsilon, \mathbf{F}, d_{L^1(P_{\vec{z}})}) \leq \mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P_{\vec{z}|_X}, \rho_l)}) \leq \mathcal{C}(\epsilon, \mathcal{H}, \rho_l)$$

for all $\vec{z} \in Z^{2m}$. Hence, setting $\epsilon = \alpha\nu/8$, the given probability is at most $4\mathcal{C}(\alpha\nu/8, \mathcal{H}, \rho_l)e^{-\alpha^2\nu m/8M}$ by Theorem 3. \square

In order to apply the above theorem, we need tools for bounding the capacity of various decision rule spaces. Along these lines, we close this section by proving two basic lemmas, one about the capacity of the free product of a set of function classes, and the other about the capacity of compositions of functions classes.

Definition 8 Let $(A_1, \rho_1), \dots, (A_k, \rho_k)$ be bounded metric spaces. Let $A = A_1 \times \dots \times A_k$ and ρ be the metric on A defined by

$$\rho(\vec{u}, \vec{v}) = \frac{1}{k} \sum_{j=1}^k \rho_j(u_j, v_j)$$

for any $\vec{u} = (u_1, \dots, u_k)$ and $\vec{v} = (v_1, \dots, v_k) \in A$. For each j , $1 \leq j \leq k$, let \mathcal{H}_j be a family of functions from X into A_j . The free product of \mathcal{H}_1 through \mathcal{H}_k is the class of functions

$$\mathcal{H} = \{(f_1, \dots, f_k) : f_j \in \mathcal{H}_j, 1 \leq j \leq k\},$$

where $(f_1, \dots, f_k) : X \rightarrow A$ is the function defined by

$$(f_1, \dots, f_k)(x) = (f_1(x), \dots, f_k(x)).$$

Lemma 7 If $\mathcal{H}, \mathcal{H}_1, \dots, \mathcal{H}_k$ are defined as above then

1. for any probability measure P on X and $\epsilon > 0$,

$$\prod_{j=1}^k \mathcal{N}(2k\epsilon, \mathcal{H}_j, d_{L^1(P, \rho_j)}) \leq \mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P, \rho)}) \leq \prod_{j=1}^k \mathcal{N}(\epsilon, \mathcal{H}_j, d_{L^1(P, \rho_j)}),$$

2. $\overline{\dim}(\mathcal{H}) = \sum_{j=1}^k \overline{\dim}(\mathcal{H}_j)$, and similarly for $\underline{\dim}$ and \mathbf{dim} , when the latter is defined.

If $\mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P, \rho)})$ is infinite for some measure P , or if the set in this supremum is unbounded, then $\mathcal{C}(\epsilon, \mathcal{H}, \rho) = \infty$. We call¹³ $\mathcal{C}(\epsilon, \mathcal{H}, \rho)$ the capacity of \mathcal{H} . In analogy with the definition of metric dimension, we define the upper metric dimension of \mathcal{H} by

$$\overline{\mathbf{dim}}(\mathcal{H}) = \limsup_{\epsilon \rightarrow 0} \frac{\log \mathcal{C}(\epsilon, \mathcal{H}, \rho)}{\log(1/\epsilon)},$$

and the lower metric dimension, denoted by $\underline{\mathbf{dim}}(\mathcal{H})$, is defined similarly using \liminf . When $\overline{\mathbf{dim}}(\mathcal{H}) = \underline{\mathbf{dim}}(\mathcal{H})$, then this quantity is denoted $\mathbf{dim}(\mathcal{H})$, and referred to simply as the metric dimension of \mathcal{H} . If $\mathcal{C}(\epsilon, \mathcal{H}, \rho) = \infty$ for some $\epsilon > 0$ then $\mathbf{dim}(\mathcal{H}) = \infty$.

We now show how bounds on the capacity of \mathcal{H} lead to distribution-independent bounds on the rate of uniform convergence of empirical risk estimates for functions in \mathcal{H} with respect to the loss function l . As before, let $Z = X \times Y$, P be a probability distribution on Z , and $l_{\mathcal{H}}$ be the family of functions on Z defined by $l_{\mathcal{H}} = \{l_h : h \in \mathcal{H}\}$, where $l_h(x, y) = l(y, h(x))$. Let $P_{|X}$ be the marginal on X of the joint distribution P on $X \times Y$ (see section 1.5).

Lemma 6 For all $\epsilon > 0$,

$$\mathcal{N}(\epsilon, l_{\mathcal{H}}, d_{L^1(P)}) \leq \mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P_{|X}, \rho_l)}).$$

Proof. For every $h \in \mathcal{H}$ let $\psi(h) = l_h$. Hence ψ maps from \mathcal{H} onto $l_{\mathcal{H}}$. It suffices to show that ψ is a contraction, i.e. that

$$\forall f, g \in \mathcal{H}, d_{L^1(P)}(\psi(f), \psi(g)) \leq d_{L^1(P_{|X}, \rho_l)}(f, g).$$

Let f and g be any two functions in \mathcal{H} . Then

$$\begin{aligned} d_{L^1(P)}(\psi(f), \psi(g)) &= \int_Z |l(y, f(x)) - l(y, g(x))| dP(x, y) \\ &\leq \int_Z \rho_l(f(x), g(x)) dP(x, y) \\ &= \int_X \rho_l(f(x), g(x)) dP_{|X}(x) \\ &= d_{L^1(P_{|X}, \rho_l)}(f, g). \end{aligned}$$

□

This gives the following theorem about distribution-independent uniform convergence of risk estimates for learning.

¹³The term *metric entropy* is often used for the quantities $\log \mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P, \rho)})$ and $\log \mathcal{C}(\epsilon, \mathcal{H}, \rho)$ [Dud87, Qui89]. It is also used for an analogous, but fundamentally distinct, concept in the dynamical systems literature (e.g. [Far82]). The term *capacity* has also been used with many other related meanings [Man82, Vap82, KT61, FOY83, BH89]. Our usage here is taken from [Dud84].

6 Capacity and Metric Dimension of Function Classes

In sections 4 and 5 we showed how the pseudo dimension can be used to obtain distribution independent bounds on the random covering numbers needed for Theorem 2, thereby obtaining bounds on the sample size needed for uniform convergence and learning results. In this section we develop an alternate way of obtaining distribution independent bounds on random covering numbers. This method can sometimes be used in conjunction with the method given in the previous sections to extend that method to cover cases where the decision space A is not contained in \mathfrak{R} . We will demonstrate this in our analysis of the sample size needed for learning in feedforward neural networks in the following section.

The key idea is to introduce a pseudo metric (see section 10.1) on the decision space A . The distance between two actions is the maximum difference in loss for these actions, over all possible outcomes.

Definition 6 For every loss function $l : Y \times A \rightarrow [0, M]$, by ρ_l we denote the pseudo metric on A defined by $\rho_l(a, b) = \sup_{y \in Y} |l(y, a) - l(y, b)|$ for all $a, b \in A$.

Note that (A, ρ_l) is a bounded pseudo metric space: no two actions in A are more than M apart.

Since decision rules in \mathcal{H} map from the instance space X into A , the pseudo metric ρ_l on A can be used to induce a pseudo metric on \mathcal{H} in which two decision rules differ only to the extent that the actions that they proscribe differ with respect to ρ_l . There are several ways to do this. The easiest is to use an L^∞ function distance on \mathcal{H} , defining the distance between decision rules f and g as the supremum of $\rho_l(f(x), g(x))$ over all $x \in X$. This works, and is a useful method of obtaining uniform convergence and learning results (see related techniques used in [Whi90a]). However, as we will see, the crucial issue is the size of the smallest ϵ -cover of the resulting pseudo metric space \mathcal{H} . In some cases we can get smaller covers, and hence better results, by using an L^1 function distance instead. Since the L^1 distance is never more than the L^∞ distance, the results are never worse. Thus we present this more powerful method here.

Definition 7 Let \mathcal{H} be a family of functions from a set X into a bounded pseudo metric space (A, ρ) . Let P be a probability measure on X . Then $d_{L^1(P, \rho)}$ is the pseudo metric on \mathcal{H} defined by

$$d_{L^1(P, \rho)}(f, g) = \mathbf{E}(\rho(f(x), g(x))) = \int_X \rho(f(x), g(x)) dP(x)$$

for all $f, g \in \mathcal{H}$. For every $\epsilon > 0$ let

$$\mathcal{C}(\epsilon, \mathcal{H}, \rho) = \sup\{\mathcal{N}(\epsilon, \mathcal{H}, d_{L^1(P, \rho)})\} \text{ over all probability measures } P \text{ on } X.$$

We now give several examples to illustrate the use of this theorem. First, consider the standard PAC model in which the outcome space Y and the decision space A are both $\{0, 1\}$. In this case any loss function such that $l(0, 1) \neq l(0, 0)$ and $l(1, 1) \neq l(1, 0)$ is monotone over A , and in particular, the standard discrete loss function, $l(y, a) = 1$ if $y \neq a$, else $l(y, a) = 0$, is monotone. Clearly $M = 1$ in this case. Thus from equation (6) above we get a sample complexity bound of

$$O\left(\frac{1}{\epsilon} \left(\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) \log \frac{1}{\epsilon} + \log \frac{1}{\delta} \right)\right).$$

As mentioned in the previous section, $\mathbf{dim}_{\mathbf{P}}(\mathcal{H})$ is the same as the Vapnik-Chervonenkis dimension of \mathbf{H} in this case, hence this bound is, up to constants, the same as that given in Theorem 2.1 of [BEHW89]. A number of applications of this result are outlined in [BEHW89]. Further applications, specifically for learning problems that have been studied recently in the mainstream artificial intelligence work, are given in [Hau88,Hau89].

For our second example, consider the case that the outcome space Y is $\{0, 1\}$ but the decision space A is $[0, 1]$. Assume that the loss function is $l(y, a) = |a \Leftrightarrow y|^q$ for some $q > 0$. This case was examined in [KS90] with $q = 2$ by Kearns and Schapire in their investigation into the learnability of p -concepts. They showed that $\mathbf{dim}_{\mathbf{P}}(l_{\mathcal{H}}) = \mathbf{dim}_{\mathbf{P}}(\mathcal{H})$ in this case. Since the loss l is monotone in A for all $q > 0$, Lemma 5 shows that this result holds for other values of q as well. Hence the conditions of Theorem 8 are met. Some applications are given in [KS90]. (It should be noted that it is important that $A = [0, 1]$ in this case. The result does not hold in general for larger A .)

For our third example, consider the problem of logistic regression, as described in section 1.1.3. In the simplest case the outcome space again has only two values, denoted y_1 and y_2 , where y_1 indicates that some event has taken place and y_2 indicates that it has not, and an action a represents an estimate of the log odds ratio $\ln(P(y_1)/P(y_2))$, where the probability P is conditioned on the observed instance x . Here $A = \mathfrak{R}$ and the log likelihood loss function is the logistic loss function, defined by $l(y_1, a) = \ln(1 + e^a) \Leftrightarrow a$ and $l(y_2, a) = \ln(1 + e^{-a})$. Again, it is easily verified that l is monotone in A . In logistic regression the standard assumption is that the instance space X is contained in \mathfrak{R}^n for some $n \geq 1$ and \mathcal{H} is contained in the family of all linear functions on X (see e.g. [MN89]). In this case, $\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) \leq n + 1$ by Theorem 4. By restricting A to a bounded range, we can then apply Theorem 8 to obtain sample complexity bounds that are linear in n .

As a last example, consider the problem of density estimation, as described in section 1.1.4. Here there is only one outcome in Y , $A \subset \mathfrak{R}^+$, and $l(y, a) = l(a) = \Leftrightarrow \log a$. Thus clearly l is monotone in A . Thus we can apply Theorem 8 to the problem of density estimation as well, whenever the family of densities \mathcal{H} is uniformly bounded (away from zero) and has finite pseudo dimension.

notation $\hat{\mathbf{r}}_{h,l}(\vec{z})$ and $\mathbf{r}_{h,l}(P)$ introduced in section 2 for the empirical risk estimate and true risk of a decision rule h , respectively.

Theorem 8 *Assume the decision space $A \subset \mathfrak{R}$, the loss function l is monotone over A and bounded between 0 and M , the decision rule space \mathcal{H} is such that $l_{\mathcal{H}}$ is permissible, and $1 \leq d = \mathbf{dim}_{\mathbf{P}}(\mathcal{H}) < \infty$. Assume $m \geq 1$, $0 < \nu < 8M$ and $0 < \alpha < 1$. Let P be any probability distribution on Z . Let \vec{z} be generated by m independent draws from Z according to P . Then*

$$\Pr(\exists h \in \mathcal{H} : d_{\nu}(\hat{\mathbf{r}}_{h,l}(\vec{z}), \mathbf{r}_{h,l}(P)) > \alpha) \leq 8 \left(\frac{16\epsilon M}{\alpha\nu} \ln \frac{16\epsilon M}{\alpha\nu} \right)^d e^{-\alpha^2\nu m/8M}.$$

Moreover, for $m \geq \frac{8M}{\alpha^2\nu} \left(2d \ln \frac{8\epsilon M}{\alpha\nu} + \ln \frac{8}{\delta} \right)$ this probability is at most δ .

Proof. By Lemma 5, $\mathbf{dim}_{\mathbf{P}}(l_{\mathcal{H}}) = \mathbf{dim}_{\mathbf{P}}(\mathcal{H})$ in this case. Thus, since $\hat{\mathbf{r}}_{h,l}(\vec{z}) = \hat{\mathbf{E}}_{\vec{z}}(l_h)$ and $\mathbf{r}_{h,l}(P) = \mathbf{E}(l_h)$, the result follows directly from Theorem 7. \square

When its conditions are satisfied, this theorem, combined with Lemma 1 from section 2.5, gives us a means of bounding in terms of $\mathbf{dim}_{\mathbf{P}}(\mathcal{H})$ the sample complexity $m(\alpha, \nu, \delta)$ of any algorithm that solves the basic learning problem by returning (with high probability) a decision rule with near minimal empirical risk on the training sample. The resulting bound is

$$m(\alpha, \nu, \delta) = O \left(\frac{M}{\alpha^2\nu} \left(\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) \log \frac{M}{\alpha\nu} + \log \frac{1}{\delta} \right) \right). \quad (5)$$

This is similar to the sample complexity

$$m(\alpha, \nu, \delta) = O \left(\frac{M}{\alpha^2\nu} \left(\log |\mathcal{H}| + \log \frac{1}{\delta} \right) \right)$$

that can be obtained by using Theorem 1, when \mathcal{H} (and hence $l_{\mathcal{H}}$) is finite. The term $\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) \log \frac{M}{\alpha\nu}$ replaces the term $\log |\mathcal{H}|$. In particular, for the ‘‘PAC settings’’ $\alpha = 1/2$ and $\nu = \epsilon$ we get the sample complexity

$$m(\epsilon, \delta) = O \left(\frac{M}{\epsilon} \left(\mathbf{dim}_{\mathbf{P}}(\mathcal{H}) \log \frac{M}{\epsilon} + \log \frac{1}{\delta} \right) \right), \quad (6)$$

in place of the sample complexity

$$m(\epsilon, \delta) = O \left(\frac{M}{\epsilon} \left(\log |\mathcal{H}| + \log \frac{1}{\delta} \right) \right)$$

derived from Theorem 1. Moreover, these bounds are distribution-independent, so \mathcal{P} can be taken to be the class of all probability distributions on Z .

Definition 5 Let $l : Y \times A \rightarrow \mathfrak{R}$ be a loss function, where $A \subset \mathfrak{R}$ and Y is an arbitrary set. For each $y \in Y$ define the function $f_y : A \rightarrow \mathfrak{R}$ by letting $f_y(a) = l(y, a)$ for each $a \in A$, i.e. f_y is the restriction of l obtained by fixing its first argument to y . We say that l is monotone over A if for every $y \in Y$, either f_y is strictly increasing on A , or f_y is strictly decreasing on A . Thus f_y may be increasing for some $y \in Y$ and decreasing for others.

Lemma 5 If $A \subset \mathfrak{R}$ and l is a loss function on $A \times Y$ that is monotone over A , then $\mathbf{dimp}(l_{\mathcal{H}}) = \mathbf{dimp}(\mathcal{H})$.

Proof. Suppose that $\vec{x} = (x_1, \dots, x_d)$ is shattered by \mathcal{H} , where $x_i \in X$, $1 \leq i \leq d$. Then there is some real vector $\vec{r} = (r_1, \dots, r_d)$ such that $\mathcal{H}_{|\vec{x}} \Leftrightarrow \vec{r}$ intersects all 2^d orthants of \mathfrak{R}^d . Hence for every Boolean vector $\vec{b} \in \{0, 1\}^d$ there exists a function $h_{\vec{b}} \in \mathcal{H}$ such that for every i , $1 \leq i \leq d$, we have $h_{\vec{b}}(x_i) > r_i$ if and only if the i^{th} bit of \vec{b} is 1. Fix an outcome $y \in Y$. Let $\vec{z} = (z_1, \dots, z_d)$, where $z_i = (x_i, y)$ for all i , $1 \leq i \leq d$. Note that if f_y is strictly increasing, then for any $h \in \mathcal{H}$ and $1 \leq i \leq d$, $h(x_i) > r_i \Leftrightarrow l_h(z_i) = l(y, h(x_i)) = f_y(h(x_i)) > f_y(r_i)$. Hence, for every Boolean vector $\vec{b} \in \{0, 1\}^d$ there exists a function $h_{\vec{b}} \in \mathcal{H}$ such that for every i , $1 \leq i \leq d$, we have $l_{h_{\vec{b}}}(z_i) > f_y(r_i)$ if and only if the i^{th} bit of \vec{b} is 1. A similar result holds if f_y is strictly decreasing. Thus \vec{z} is shattered by $l_{\mathcal{H}}$. It follows that $\mathbf{dimp}(l_{\mathcal{H}}) \geq \mathbf{dimp}(\mathcal{H})$.

For the other direction, assume $\vec{z} = (z_1, \dots, z_d)$ is shattered by $l_{\mathcal{H}}$, where $z_i = (x_i, y_i)$ for all i , $1 \leq i \leq d$. We will show that $\vec{x} = (x_1, \dots, x_d)$ is shattered by \mathcal{H} . Since \vec{z} is shattered by $l_{\mathcal{H}}$, there is some real vector $\vec{r} = (r_1, \dots, r_d)$ such that for every Boolean vector $\vec{b} \in \{0, 1\}^d$ there exists a function $h_{\vec{b}} \in \mathcal{H}$ such that for every i , $1 \leq i \leq d$, we have $l_{h_{\vec{b}}}(z_i) > r_i$ if and only if the i^{th} bit of \vec{b} is 1. Let $A_0 = \{h_{\vec{b}}(x_i) : 1 \leq i \leq d \text{ and } \vec{b} \in \{0, 1\}^d\}$. For each outcome $y \in Y$ define the function $g_y : \mathfrak{R} \rightarrow A_0$ as follows. If f_y is increasing then $g_y(r) = \max\{a \in A_0 : f_y(a) \leq r\}$ and if f_y is decreasing then $g_y(r) = \max\{a \in A_0 : f_y(a) > r\}$. Then for each i , $1 \leq i \leq d$, we either have

1. for all $h \in \mathcal{H}$, $l_h(z_i) = f_{y_i}(h(x_i)) > r_i \Leftrightarrow h(x_i) > g_{y_i}(r_i)$ or
2. for all $h \in \mathcal{H}$, $l_h(z_i) = f_{y_i}(h(x_i)) > r_i \Leftrightarrow h(x_i) \leq g_{y_i}(r_i)$

Hence, for every Boolean vector $\vec{b} \in \{0, 1\}^d$ there exists a function $h'_{\vec{b}} \in \mathcal{H}$ such that for every i , $1 \leq i \leq d$, we have $h'_{\vec{b}}(x_i) > g_{y_i}(r_i)$ if and only if the i^{th} bit of \vec{b} is 1. (To see this, let \vec{c} be the Boolean vector derived from \vec{b} by complementing the bit in each position i for which f_{y_i} is decreasing, and then let $h'_{\vec{b}} = h_{\vec{c}}$.) Thus $\vec{x} = (x_1, \dots, x_d)$ is shattered by \mathcal{H} . It follows that $\mathbf{dimp}(\mathcal{H}) \geq \mathbf{dimp}(l_{\mathcal{H}})$, and combined with the above inequality, this gives the result. \square

Combined with Theorem 7, this gives the following result on the uniform convergence of empirical risk estimates for the basic learning problem. Here and below we use the

for all $\vec{z} \in Z^{2m}$. Hence the given probability is at most

$$8 \left(\frac{2eM}{\epsilon} \ln \frac{2eM}{\epsilon} \right)^d e^{-\alpha^2 \nu m / 8M} = 8 \left(\frac{16eM}{\alpha \nu} \ln \frac{16eM}{\alpha \nu} \right)^d e^{-\alpha^2 \nu m / 8M}$$

by Theorem 3.

For the second result, setting the bound above equal to δ and solving for m gives

$$m \geq \frac{8M}{\alpha^2 \nu} \left(d \ln \left(\frac{16eM}{\alpha \nu} \ln \frac{16eM}{\alpha \nu} \right) + \ln \frac{8}{\delta} \right).$$

It is easily verified that $\ln(a \ln a) < 2 \ln(a/2)$ when $a \geq 5$, and from this the bound given in the second result follows. \square

Corollary 2 *Under the same assumptions as above, for all $0 < \epsilon \leq M$,*

$$\Pr \left(\exists f \in \mathbf{F} : |\widehat{\mathbf{E}}_{\vec{z}}(f) - \mathbf{E}(f)| > \epsilon \right) \leq 8 \left(\frac{32eM}{\epsilon} \ln \frac{32eM}{\epsilon} \right)^d e^{-\epsilon^2 m / 64M^2}.$$

Moreover, for $m \geq \frac{64M^2}{\epsilon^2} \left(2d \ln \frac{16eM}{\epsilon} + \ln \frac{8}{\delta} \right)$ this probability is at most δ .

Proof. This follows directly from the above result by setting $\nu = 2M$, $\alpha = \epsilon/4M$, and using property (3) of the d_ν metric, as in the proof of Corollary 1 in section 3.2. \square

5 Some Applications of Pseudo Dimension in Learning

We now look at how the theoretical results obtained in the previous two sections can be applied to certain types of learning problems. Suppose that we have a basic learning problem defined by $X, Y, A, \mathcal{H}, \mathcal{P}$, and \mathcal{L} , where \mathcal{L} is the family of $L_{\alpha, \nu}$ regret functions for an underlying loss function l , as in section 2. As before, let $Z = X \times Y$, $l_h : Z \rightarrow [0, M]$ be defined by $l_h(x, y) = l(y, h(x))$ for all $h \in \mathcal{H}$, and $l_{\mathcal{H}} = \{l_h : h \in \mathcal{H}\}$. In this section we will show how to obtain sample complexity bounds on algorithms for this basic learning problem using Theorem 7 above.

To obtain these bounds, we will need bounds on $\mathbf{dim}_{\mathbf{P}}(l_{\mathcal{H}})$. In this section we will look at some useful tricks for computing $\mathbf{dim}_{\mathbf{P}}(l_{\mathcal{H}})$ in the important case $A \subset \mathfrak{R}$, i.e. when each decision is represented by a real number. In the following section we discuss more general decision spaces.

When $A \subset \mathfrak{R}$, the functions in \mathcal{H} are themselves real-valued, so we can talk about the pseudo dimension of \mathcal{H} itself, without reference to any particular loss function. What makes this useful is that in many important cases the pseudo dimension of \mathcal{H} is the same as the pseudo dimension of $l_{\mathcal{H}}$. Thus we can factor out the effects of the loss function in deriving our sample size bounds, and concentrate on the pseudo dimension of the decision rule space \mathcal{H} .

Since this holds for every distinct f and g in \mathbf{G} , combining this with the inequality above, we have

$$\mathbf{E}(|\text{sign}(\mathbf{F}_{|\vec{z}} - \vec{r}^\dagger)|) \geq |\mathbf{G}|(1 - |\mathbf{G}|e^{-\epsilon m/M}).$$

Since $|\mathbf{G}| = \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})$, this gives the result. \square

Proof of Theorem 6. Since $\text{dimp}(\mathbf{F}) = d$, it follows from Sauer's lemma that $|\text{sign}(\mathbf{F}_{|\vec{z}} - \vec{r}^\dagger)| \leq (\epsilon m/d)^d$ for all $m \geq d$, $\vec{z} \in Z^m$ and $\vec{r} \in [0, M]^m$. Hence the above lemma implies that

$$(\epsilon m/d)^d \geq \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}) \left(1 - \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})e^{-\epsilon m/M}\right) \quad (4)$$

for all probability measures P on Z and $m \geq d$. It is easily verified that if

$$\frac{M}{\epsilon} \ln(2\mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})) < d$$

then the upper bound given in Theorem 6 follows trivially using the fact that $\epsilon \leq M$. Thus we may assume that $\frac{M}{\epsilon} \ln(2\mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})) \geq d$. Hence, if $m \geq \frac{M}{\epsilon} \ln(2\mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}))$, then $m \geq d$ and

$$(1 - \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})e^{-\epsilon m/M}) \geq 1/2.$$

Thus from (4) we obtain

$$\left(\frac{eM \ln(2\mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}))}{\epsilon d}\right)^d \geq \frac{1}{2} \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}).$$

With some simple calculations, this gives the bound of Theorem 6. \square

Using our results on uniform convergence from sections 3.2 and 10.4, we can now show the following.

Theorem 7 *Let \mathbf{F} be a permissible family of functions from a set Z into $[0, M]$ with $\text{dimp}(\mathbf{F}) = d$ for some $1 \leq d < \infty$. Assume $m \geq 1$, $0 < \nu \leq 8M$ and $0 < \alpha < 1$. Let \vec{z} be generated by m independent draws according to any distribution on Z . Then*

$$\Pr\left(\exists f \in \mathbf{F} : d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha\right) \leq 8 \left(\frac{16eM}{\alpha\nu} \ln \frac{16eM}{\alpha\nu}\right)^d e^{-\alpha^2 \nu m / 8M}.$$

Moreover, for $m \geq \frac{8M}{\alpha^2 \nu} \left(2d \ln \frac{8eM}{\alpha\nu} + \ln \frac{8}{\delta}\right)$ this probability is at most δ .

Proof. Let $\epsilon = \alpha\nu/8$. Since $\alpha < 1$ and $\nu \leq 8M$, $\epsilon \leq M$. For any sequence \vec{z} of points in Z there is a trivial isometry between $(\mathbf{F}_{|\vec{z}}, d_{L^1})$ and $(\mathbf{F}, d_{L^1(P_{\vec{z}})})$, where $P_{\vec{z}}$ is the empirical measure induced by \vec{z} , in which each set has measure equal to the fraction of the points in \vec{z} it contains. Thus by Theorem 12 of section 10.1 and Theorem 6, we have

$$\mathcal{N}(\epsilon, \mathbf{F}_{|\vec{z}}, d_{L^1}) \leq \mathcal{M}(\epsilon, \mathbf{F}_{|\vec{z}}, d_{L^1}) \leq 2 \left(\frac{2eM}{\epsilon} \ln \frac{2eM}{\epsilon}\right)^d,$$

$\mathbf{F}|_{\mathcal{Z}} = \{0, 1\}^d$. Then ¹²

$$|\mathbf{F}| \leq \sum_{i=0}^d \binom{m}{i} \leq (em/d)^d$$

where e is the base of the natural logarithm. \square

In the next lemma we bound the packing numbers of $(\mathbf{F}, d_{L^1(P)})$ in terms of the expected number of orthants intersected by a random translation of a random restriction of \mathbf{F} . This is the key lemma of the proof.

Lemma 4 *Let \mathbf{F} be a family of functions from a set Z into $[0, M]$ and let P be a probability measure on Z . Let $\vec{r} = (r_1, \dots, r_m)$ be a random vector in $[0, M]^m$ where each r_i is drawn independently at random from the uniform distribution on $[0, M]$. Let $\vec{z} = (z_1, \dots, z_m)$ be a random vector in Z^m where each z_i is drawn independently at random from P . Then for all $\epsilon > 0$*

$$\mathbf{E}(|\text{sign}(\mathbf{F}|_{\mathcal{Z}} - \vec{r})|) \geq \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}) \left(1 - \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}) e^{-\epsilon m/M}\right).$$

Proof. For all $f \in \mathbf{F}$ we will denote $(f(z_1), \dots, f(z_m))$ by $f|_{\mathcal{Z}}$. Choose $\epsilon > 0$. Let \mathbf{G} be an ϵ -separated subset of \mathbf{F} (w.r.t. $d_{L^1(P)}$), with $|\mathbf{G}| = \mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)})$. Then

$$\begin{aligned} \mathbf{E}(|\text{sign}(\mathbf{F}|_{\mathcal{Z}} - \vec{r})|) &\geq \mathbf{E}(|\text{sign}(\mathbf{G}|_{\mathcal{Z}} - \vec{r})|) \\ &\geq \mathbf{E}(|\{f \in \mathbf{G} : \text{sign}(f|_{\mathcal{Z}} - \vec{r}) \neq \text{sign}(g|_{\mathcal{Z}} - \vec{r}) \text{ for all } g \in \mathbf{G}, g \neq f\}|) \\ &= \sum_{f \in \mathbf{G}} \Pr(\text{sign}(f|_{\mathcal{Z}} - \vec{r}) \neq \text{sign}(g|_{\mathcal{Z}} - \vec{r}) \text{ for all } g \in \mathbf{G}, g \neq f) \\ &= \sum_{f \in \mathbf{G}} \left(1 - \Pr(\exists g \in \mathbf{G}, g \neq f : \text{sign}(f|_{\mathcal{Z}} - \vec{r}) = \text{sign}(g|_{\mathcal{Z}} - \vec{r}))\right) \\ &\geq \sum_{f \in \mathbf{G}} \left(1 - |\mathbf{G}| \max_{g \in \mathbf{G}, g \neq f} \Pr(\text{sign}(f|_{\mathcal{Z}} - \vec{r}) = \text{sign}(g|_{\mathcal{Z}} - \vec{r}))\right). \end{aligned}$$

Let f and g be distinct functions in \mathbf{G} . Since \mathbf{G} is ϵ -separated,

$$\int_Z |f(z) - g(z)| dP(z) > \epsilon.$$

In addition, the range of f and g is $[0, M]$. Hence if z_i is drawn at random from P and r_i drawn at random from the uniform distribution on $[0, M]$, then the probability that r_i lies between $f(z_i)$ and $g(z_i)$ is at least ϵ/M . And $\text{sign}(f|_{\mathcal{Z}} - \vec{r}) = \text{sign}(g|_{\mathcal{Z}} - \vec{r})$ only if this fails to occur for each i , $1 \leq i \leq m$. Thus

$$\Pr(\text{sign}(f|_{\mathcal{Z}} - \vec{r}) = \text{sign}(g|_{\mathcal{Z}} - \vec{r})) \leq \left(1 - \frac{\epsilon}{M}\right)^m \leq e^{-\epsilon m/M}.$$

¹²See e.g. [Dud84], Prop. 2.2.9, or [BEHW89], Appendix, for a proof of the second inequality. Note also that [VC71] actually contains a slightly weaker result.

the absolute value of their difference, i.e. the L^1 distance, relative to the given measure. To make this work, we need to make some assumptions about the integrability of the functions in \mathbf{F} under the given measure. Since we will be concerned only with families of functions taking values in a bounded range in this paper, this will cause no problems for us. For convenience, we choose this range to be $[0, M]$. For a more general treatment, see [Pol84][Dud84].

Definition 4 *Let \mathbf{F} be a class of functions from Z into $[0, M]$, where $M > 0$, and P be a probability measure on Z . Then $d_{L^1(P)}$ is the pseudo metric on \mathbf{F} defined by*

$$d_{L^1(P)}(f, g) = \mathbf{E}(|f - g|) = \int_Z |f(z) - g(z)| dP(z) \text{ for all } f, g \in \mathbf{F}.$$

The notions of ϵ -cover and metric dimension used in the previous section can be generalized to arbitrary pseudo metric spaces. This generalization is given in section 10.1 of the appendix. In the remainder of the paper we will use the concepts and notation given there without further special reference.

Using techniques that go back to Dudley [Dud78], Pollard has obtained a beautiful theorem bounding the metric dimension of $(\mathbf{F}, d_{L^1(P)})$ by $\mathbf{dim}_P(\mathbf{F})$ for any probability measure P on Z . Actually this result is much stronger in that it gives explicit bounds on the packing numbers for \mathbf{F} using $d_{L^1(P)}$ balls of radius ϵ . Since packing numbers are closely related to covering numbers (Theorem 12 in section 10.1), these bounds can then be used with Theorem 2 to obtain uniform convergence results for empirical estimates of functions in \mathbf{F} . We now state and prove a version of Pollard's result ([Pol84], Lemma 25, p. 27) for the special case when \mathbf{F} is a class of functions taking values in the interval $[0, M]$ with somewhat better bounds on the packing numbers.

Theorem 6 (Pollard) *Let \mathbf{F} be a family of functions from a set Z into $[0, M]$, where $\mathbf{dim}_P(\mathbf{F}) = d$ for some $1 \leq d < \infty$. Let P be a probability measure on Z . Then for all $0 < \epsilon \leq M$,*

$$\mathcal{M}(\epsilon, \mathbf{F}, d_{L^1(P)}) < 2 \left(\frac{2\epsilon M}{\epsilon} \ln \frac{2\epsilon M}{\epsilon} \right)^d.$$

The proof we give uses essentially the same techniques as Pollard's, with some minor modifications. It relies on a few lemmas, which we give now. The first, which we give without proof, was discovered independently by a number of people (see [Ass83]), including Vapnik and Chervonenkis [VC71], but is most often attributed to Sauer [Sau72] in the computer science literature.

Lemma 3 (Sauer) *Let \mathbf{F} be a class of functions from $S = \{1, 2, \dots, m\}$ into $\{0, 1\}$ with $|\mathbf{F}| > 1$ and let d be the length of the longest sequence of points \vec{z} from S such that*

On the other hand, if \mathbf{F} is a d -dimensional vector space of real-valued functions on Z , then there exists a sequence \vec{z} of d points in Z such that $\mathbf{F}|_{\vec{z}} = \mathbb{R}^d$. Hence \vec{z} is shattered, implying that $\mathbf{dim}_P(\mathbf{F}) \geq d$. \square

There are many other ways that the VC dimension can be generalized to real-valued functions [Nat89b] [Nat89a] [Pol84] [Vap89] [Dud87]. Dudley [Dud87] compares several such generalizations, albeit in a different context. The generalization we have proposed here, the pseudo dimension, is a minor variant of the notion used by Pollard in [Pol84] to define classes of real-valued functions of polynomial discrimination, called VC-subgraph classes in [Dud87]. The pseudo dimension will be used in the form defined above in Pollard's new book [Pol90].

The pseudo dimension has a few invariance properties that are useful (see [Pol90] for further results of this type).

Theorem 5 *Let \mathbf{F} be a family of functions from Z into \mathbb{R} . Fix any function g from Z into \mathbb{R} and let $\mathbf{G} = \{g + f : f \in \mathbf{F}\}$. Let I be a real interval (possibly all of \mathbb{R}) such that every function in \mathbf{F} takes values only in I . Fix any nondecreasing (resp. nonincreasing) function $h : I \rightarrow \mathbb{R}$ and let $\mathbf{H} = \{h \circ f : f \in \mathbf{F}\}$, where \circ indicates function composition. Then*

1. ([WD81]) $\mathbf{dim}_P(\mathbf{G}) = \mathbf{dim}_P(\mathbf{F})$ and
2. ([NP87,Dud87]) $\mathbf{dim}_P(\mathbf{H}) \leq \mathbf{dim}_P(\mathbf{F})$, with equality if h is continuous and strictly increasing (resp. continuous and strictly decreasing).

Proof. Part (1) follows directly from the fact that the notion of a set of points being full is invariant under translation. For part (2) it suffices to prove the results for h nondecreasing and h continuous and strictly increasing. Let $\vec{z} = (z_1, \dots, z_d)$ be such that $\mathbf{H}|_{\vec{z}}$ is full, i.e. such that $\mathbf{H}|_{\vec{z}} - \vec{x}$ intersects all 2^d orthants of \mathbb{R}^d for some vector $\vec{x} = (x_1, \dots, x_d)$ in \mathbb{R}^d . Then for every Boolean vector $\vec{b} \in \{0, 1\}^d$ there exists a function $f_{\vec{b}} \in \mathbf{F}$ such that for every i , $1 \leq i \leq d$, we have $h \circ f_{\vec{b}}(z_i) > x_i$ if and only if the i^{th} bit of \vec{b} is 1. For each i , $1 \leq i \leq d$, let

$$u_i = \min\{f_{\vec{b}}(z_i) : \text{the } i^{\text{th}} \text{ bit of } \vec{b} \text{ is } 1\}$$

and

$$l_i = \max\{f_{\vec{b}}(z_i) : \text{the } i^{\text{th}} \text{ bit of } \vec{b} \text{ is } 0\}.$$

Since h is nondecreasing, we have $u_i > l_i$ for each i . Let $r_i = (u_i + l_i)/2$ for each i and $\vec{r} = (r_1, \dots, r_d)$. Let $T = \{f_{\vec{b}} : \vec{b} \in \{0, 1\}^d\}$. Then clearly $T - \vec{r}$ intersects every orthant of \mathbb{R}^d , so T is full. Since $T \subset \mathbf{F}$, this implies that $\mathbf{F}|_{\vec{z}}$ is full, and hence $\mathbf{dim}_P(\mathbf{H}) \leq \mathbf{dim}_P(\mathbf{F})$. Equality follows when h is continuous and strictly increasing since we obtain the class \mathbf{F} from \mathbf{H} by composing with h^{-1} . \square

By putting a probability measure on Z , we can view a class \mathbf{F} of real-valued functions on Z as a pseudo metric space. The distance between two functions is the integral of

Proof. Let T be a hyperplane in \mathfrak{R}^d . Choose a vector $\vec{x} \in \mathfrak{R}^d$ as follows. If T includes the origin, then let \vec{x} be any vector that is orthogonal to T and has at least one strictly negative coordinate. (For any nonzero orthogonal vector \vec{x} , if \vec{x} doesn't have a negative coordinate then $-\vec{x}$ does.) Otherwise let \vec{x} be the (nonzero) vector in T on the line perpendicular to T that passes through the origin. To complete the proof, we show that for all $\vec{y} \in T$, $\text{sign}(\vec{y}) \neq \vec{1} - \text{sign}(\vec{x})$, where $\vec{1}$ denotes the all 1's vector.

Suppose to the contrary that $\text{sign}(\vec{y}) = \vec{1} - \text{sign}(\vec{x})$ for some $\vec{y} \in T$. This implies that the inner product $\sum_{i=1}^d x_i y_i$ is non positive, and is in fact strictly negative if either \vec{x} or \vec{y} contain a strictly negative coordinate. However, by our choice of \vec{x} , either \vec{x} is orthogonal to \vec{y} and contains a strictly negative coordinate, giving an immediate contradiction, or \vec{x} is non-zero and \vec{x} is orthogonal to $\vec{y} - \vec{x}$. In this last case,

$$\sum_{i=1}^d x_i y_i = \sum_{i=1}^d x_i^2,$$

which is again a contradiction, since the left side is non-positive while the right side is strictly positive. \square

It follows from this lemma that if T is contained in a hyperplane of \mathfrak{R}^d then T is not full.

Definition 3 Let \mathbf{F} be a family of functions from a set Z into \mathfrak{R} . For any sequence $\vec{z} = (z_1, \dots, z_d)$ of points in Z , let $\mathbf{F}_{|\vec{z}} = \{(f(z_1), \dots, f(z_d)) : f \in \mathbf{F}\}$. If $\mathbf{F}_{|\vec{z}}$ is full then we say that \vec{z} is shattered by \mathbf{F} . The pseudo dimension of \mathbf{F} , denoted $\mathbf{dimp}(\mathbf{F})$, is the largest d such that there exists a sequence of d points in Z that is shattered by \mathbf{F} . If arbitrarily long finite sequences are shattered, then $\mathbf{dimp}(\mathbf{F})$ is infinite.

It is clear that when \mathbf{F} is a set of $\{0, 1\}$ -valued functions then for any sequence \vec{z} of d points in Z , $\mathbf{F}_{|\vec{z}}$ is full if and only if $\mathbf{F}_{|\vec{z}} = \{0, 1\}^d$. Thus in this case $\mathbf{dimp}(\mathbf{F})$ is the length d of the longest sequence of points \vec{z} such that $\mathbf{F}_{|\vec{z}} = \{0, 1\}^d$. This is the definition of the Vapnik-Chervonenkis dimension of a class \mathbf{F} of $\{0, 1\}$ -valued functions [Vap82][HW87][BEHW89]. Thus the pseudo dimension generalizes the Vapnik-Chervonenkis dimension to arbitrary classes of real-valued functions.

The pseudo dimension also generalizes the algebraic notion of the dimension of a vector space of real-valued functions [Dud78].

Theorem 4 (Dudley) Let \mathbf{F} be a d -dimensional vector space of functions from a set Z into \mathfrak{R} . Then $\mathbf{dimp}(\mathbf{F}) = d$.

Proof. Fix any sequence $\vec{z} = (z_1, \dots, z_{d+1})$ of points in Z . For any $f \in \mathbf{F}$ let $\Psi(f) = (f(z_1), \dots, f(z_{d+1}))$. Then Ψ is a linear mapping from \mathbf{F} into \mathfrak{R}^{d+1} , and the image of Ψ is $\mathbf{F}_{|\vec{z}}$. Since \mathbf{F} is a vector space of dimension d , this implies that $\mathbf{F}_{|\vec{z}}$ is a subspace of \mathfrak{R}^{d+1} of dimension at most d . Hence by Lemma 2, $\mathbf{F}_{|\vec{z}}$ is not full. This implies $\mathbf{dimp}(\mathbf{F}) \leq d$.

Theorem 3 *Let \mathbf{F} be a permissible set of functions on Z with $0 \leq f(z) \leq M$ for all $f \in \mathbf{F}$ and $z \in Z$. Assume $\nu > 0$, $0 < \alpha < 1$ and $m \geq 1$. Suppose that \vec{z} is generated by m independent random draws according to any probability measure on Z . Let*

$$p(\alpha, \nu, m) = \Pr \left\{ \vec{z} \in Z^m : \exists f \in \mathbf{F} \text{ with } d_\nu(\hat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha \right\}.$$

Then

$$p(\alpha, \nu, m) \leq 2\mathbf{E}(\min(2\mathcal{N}(\alpha\nu/8, \mathbf{F}_{|\vec{z}}, d_{L^1})e^{-\alpha^2\nu m/8M}, 1)),$$

where the expectation is over \vec{z} drawn randomly from Z^{2m} . If in addition $\mathbf{F}_{|\vec{z}}$ is finite for all $\vec{z} \in Z^{2m}$ then

$$p(\alpha, \nu, m) \leq 2\mathbf{E}(\min(2|\mathbf{F}_{|\vec{z}}|e^{-\alpha^2\nu m/2M}, 1)).$$

Theorem 2 is obtained as a corollary of this result by substituting $m/2$ for m and not taking the minimum with 1 in the left hand side of the first bound for $p(\alpha, \nu, m)$. We will use Theorem 3 to obtain slightly better constants in some of the results in the sequel.

4 Pseudo dimension of classes of real-valued functions

In this section we will look at one way that bounds on the covering numbers appearing in Theorem 2 can be obtained. This technique, due to Pollard [Pol84], who extended methods from [Dud78], is based on certain intuitions from combinatorial geometry. It generalizes the techniques based on the Vapnik-Chervonenkis dimension used in [BEHW89], which apply only to $\{0, 1\}$ -valued functions. We begin by establishing some basic notation.

Definition 2 *For $x \in \mathbb{R}$, let $\text{sign}(x) = 1$ if $x > 0$ else $\text{sign}(x) = 0$. For $\vec{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, let $\text{sign}(\vec{x}) = (\text{sign}(x_1), \dots, \text{sign}(x_d))$ and for $T \subset \mathbb{R}^d$ let $\text{sign}(T) = \{\text{sign}(\vec{x}) : \vec{x} \in T\}$. For any Boolean vector $\vec{b} = (b_1, \dots, b_d)$, $\{\vec{x} \in \mathbb{R}^d : \text{sign}(\vec{x}) = \vec{b}\}$ is called the \vec{b} -orthant of \mathbb{R}^d , where we have, somewhat arbitrarily, included points with value zero for a particular coordinate in the associated lower orthant. Thus $\text{sign}(T)$ denotes the set of orthants intersected by T . For any $T \subset \mathbb{R}^d$, and $\vec{x} \in \mathbb{R}^d$, let $T + \vec{x} = \{\vec{y} + \vec{x} : \vec{y} \in T\}$, i.e. the translation of T obtained by adding the vector \vec{x} . We say that T is full if there exists $\vec{x} \in \mathbb{R}^d$ such that $\text{sign}(T + \vec{x}) = \{0, 1\}^d$, i.e. if there exists some translation of T that intersects all 2^d orthants of \mathbb{R}^d .*

The following result is well known and can be proved in a variety of ways. For example, it follows easily from well known bounds on the number of cells in arrangements of hyperplanes (see e.g. [Ede87]). We give an elementary proof using a technique from [Dud78].

Lemma 2 *No hyperplane in \mathbb{R}^d intersects all orthants of \mathbb{R}^d .*

$|r - s| \leq \epsilon$ whenever $d_\nu(r, s) \leq \alpha$ for all $0 \leq r, s \leq M$ when this setting of ν and α is used. \square

The constants in these results are only crude estimates. No serious attempt has been made to minimize them. (See the recent results of Talagrand [Tal91] for much better constants for Corollary 1).

The bound in this latter result depends critically on the relative magnitudes of the negative exponent in $e^{-\epsilon^2 m/128M^2}$ and the exponent in the expectation of the covering number $\mathcal{N}(\epsilon/16, \mathbf{F}_{|\vec{z}})$, which reflects the extent to which $\mathbf{F}_{|\vec{z}}$ “fills up” the m -cube $[0, M]^m$. For example, if $\mathbf{F}_{|\vec{z}}$ has metric dimension at most n for all m and all \vec{z} , then there is a constant c_0 such that for any $\eta > 0$, $\mathcal{N}(\epsilon/16, \mathbf{F}_{|\vec{z}}) \leq (c_0 M/\epsilon)^{n+\eta}$ for suitably small ϵ . In this case the negative exponential term eventually dominates the expected covering number, and beyond a critical sample size

$$m_0 = O\left(\frac{nM^2}{\epsilon^2} \log \frac{M}{\epsilon}\right),$$

the bound goes to zero exponentially fast. We will see examples of this in the following section, where we give bounds on the metric dimension of $\mathbf{F}_{|\vec{z}}$ in terms of a combinatorial parameter called the pseudo dimension of \mathbf{F} . The theorem actually shows that this exponential drop off occurs even if this metric dimension bound holds only for “most” \vec{z} .

On the other hand, if with high probability $\mathbf{F}_{|\vec{z}}$ “fills up” the m -cube $[0, M]^m$ to the extent that $\mathcal{N}(\epsilon/16, \mathbf{F}_{|\vec{z}}) \approx (c_0/\epsilon)^m$, which is as large as possible, then the covering number dominates, and the bound is trivial. Results in [Vap82] (Theorem A.2, page 220) indicate that uniform convergence does not take place in this case. Similar remarks apply to the bound given in Theorem 2, which uses the d_ν metric.

The proof of Theorem 2 follows the proof of Pollard’s Theorem 24 ([Pol84], p. 25) in general outline. However, the use of the d_ν metric necessitates a number of substantial modifications. The approach taken here is different from that taken (independently, but prior to this work) by Pollard in [Pol86]. Still different, and more involved, techniques are used in the more general theory of weighted empirical processes developed by Alexander [Ale85, Ale87]. While the proof of Theorem 2 we give is simpler, it is still somewhat lengthy, so it is given in Appendix 10.4.

Actually, we can prove a slightly stronger result than Theorem 2. This result is obtained by bounding the probability of uniform convergence on a sample of length m in terms of the expected covering numbers associated with a sample of length $2m$, and by expanding the expectation to include the negative exponential term with a “truncation” at 1. It turns out that this saves us a factor of 1/2 in the negative exponential term. We also include special bounds for the case that $\mathbf{F}_{|\vec{z}}$ is always finite. This case comes up, for example, when $\mathbf{F} = l_{\mathcal{H}}$ and we use the discrete loss function l , as in the PAC learning model.

Following [KT61] we define the *upper metric dimension* of the set T of points by

$$\overline{\mathbf{dim}}(T) = \limsup_{\epsilon \rightarrow 0} \frac{\log \mathcal{N}(\epsilon, T)}{\log(1/\epsilon)}.$$

The *lower metric dimension*, denoted by \mathbf{dim} , is defined similarly using \liminf . When $\overline{\mathbf{dim}}(T) = \mathbf{dim}(T)$, then this quantity is denoted $\mathbf{dim}(T)$, and referred to simply as the *metric dimension* of T . Note that if $\mathcal{N}(\epsilon, T) = (g(\epsilon)/\epsilon)^n$, where $g(\epsilon)$ is polylogarithmic in $1/\epsilon$, then $\mathbf{dim}(T) = n$. Hence the metric dimension essentially picks out the exponent in the rate of growth of the covering number as a function of $1/\epsilon$.

Assume all functions in \mathbf{F} map from Z into $[0, M]$. For any sample $\vec{z} = (z_1, \dots, z_m)$, with $z_i \in Z$, let

$$\mathbf{F}_{|\vec{z}} = \{(f(z_1), \dots, f(z_m)) : f \in \mathbf{F}\}.$$

We will call $\mathbf{F}_{|\vec{z}}$ the *restriction of \mathbf{F} to \vec{z}* . Note that $\mathbf{F}_{|\vec{z}}$ is a set of points in the m -cube $[0, M]^m$. We can consider the size of the covering number $\mathcal{N}(\epsilon, \mathbf{F}_{|\vec{z}})$ as giving some indication of the “richness at scale $\approx \epsilon$ ” of the class \mathbf{F} of functions, restricted to the domain z_1, \dots, z_m . The metric dimension of $\mathbf{F}_{|\vec{z}}$ gives some indication of the “number of essential degrees of freedom” in this restriction of \mathbf{F} .

When z_1, \dots, z_m are drawn independently at random from Z , the *random covering number* $\mathbf{E}(\mathcal{N}(\epsilon, \mathbf{F}_{|\vec{z}}))$ gives some indication of the “richness” of \mathbf{F} on a “typical” set of m points in the domain Z . Note that for finite \mathbf{F} , we have $\mathcal{N}(\epsilon, \mathbf{F}_{|\vec{z}}) \leq |\mathbf{F}|$ for all ϵ and all samples \vec{z} , and hence the random covering number $\mathbf{E}(\mathcal{N}(\epsilon, \mathbf{F}_{|\vec{z}})) \leq |\mathbf{F}|$ for all ϵ , all sample sizes m , and all distributions on Z . The main result about uniform empirical estimates for infinite classes of functions is similar to Theorem 1 except that the random covering numbers are used in place of $|\mathbf{F}|$.

Theorem 2 ([Pol86]) *Let \mathbf{F} be a permissible¹¹ set of functions on Z with $0 \leq f(z) \leq M$ for all $f \in \mathbf{F}$ and $z \in Z$. Let $\vec{z} = (z_1, \dots, z_m)$ be a sequence of m examples drawn independently from Z according to any distribution on Z . Then for any $\nu > 0$ and $0 < \alpha < 1$,*

$$\Pr(\exists f \in \mathbf{F} : d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha) \leq 4\mathbf{E}(\mathcal{N}(\alpha\nu/8, \mathbf{F}_{|\vec{z}})) e^{-\alpha^2\nu m/16M}. \quad \square$$

Corollary 1 ([Pol84]) *Under the same assumptions as above, for all $\epsilon > 0$,*

$$\Pr(\exists f \in \mathbf{F} : |\widehat{\mathbf{E}}_{\vec{z}}(f) - \mathbf{E}(f)| > \epsilon) \leq 4\mathbf{E}(\mathcal{N}(\epsilon/16, \mathbf{F}_{|\vec{z}})) e^{-\epsilon^2 m/128M^2}.$$

Proof of corollary. This follows directly from the above result by setting $\nu = 2M$, and $\alpha = \epsilon/4M$. To see this, note that property (3) of the d_ν metric (section 2.2) implies that

¹¹This is a measurability condition defined in [Pol84] which need not concern us in practice. Further details are given in section 10.2 of the appendix.

In the latter case we get a sample complexity

$$m(\alpha, \nu, \delta) = O\left(\frac{M}{\alpha^2\nu} \left(\log |l_{\mathcal{H}}| + \log \frac{1}{\delta}\right)\right). \quad (2)$$

As shown in the previous section, a generalization of the PAC learning model can be obtained by using either the L_ϵ or $L_{\alpha,\nu}$ regret functions, in the latter case by setting $\alpha = 1/2$ and $\nu = \epsilon$. Note that plugging this latter setting into (2) gives a sample complexity

$$m(\epsilon, \delta) = O\left(\frac{M}{\epsilon} \left(\log |l_{\mathcal{H}}| + \log \frac{1}{\delta}\right)\right), \quad (3)$$

a significant improvement over (1), which is quadratic in M/ϵ . Thus the generalization of the PAC model using the d_ν metric to measure distance from optimality, and the resulting $L_{\alpha,\nu}$ family of regret functions, offers new insight in this regard. (Vapnik’s use of the relative difference between empirical estimates and true expectations [Vap82] also has this advantage; see [AST90] also the appendix of [BEHW89].)

3.2 The general case

The main task of this section is to generalize Theorem 1 to infinite collections of uniformly bounded functions. The basic idea is simple: we replace the infinite class \mathbf{F} of functions with a finite class \mathbf{F}_0 that “approximates” it, in the sense that each function in \mathbf{F} is close to some function in \mathbf{F}_0 , and argue that some type of uniform convergence of empirical estimates for \mathbf{F}_0 implies uniform convergence for \mathbf{F} . In the simplest version of this technique, the choice of \mathbf{F}_0 depends only on \mathbf{F} and the distribution P , as in the “direct method” discussed in section II.2 of [Pol84] (see also [Vap82] section 6.6, [Dud84] chapter 6, [BI88], [Whi90a]). However, more general results (apart from certain measurability constraints) are obtained by allowing \mathbf{F}_0 to depend on the particular random sample \vec{z} (e.g. [Pol84], chapter 2). Here \mathbf{F}_0 is called a “random cover”, and its size is called a “random covering number”. It is this type of result that we derive here.

We will need a few preliminary definitions to introduce the notion of ϵ -covers and metric dimension. A more general treatment of these ideas is given in the appendix, section 10.1. This more general treatment will be used later, but the following definitions suffice for this section.

For any real vectors $\vec{x} = (x_1, \dots, x_m)$ and $\vec{y} = (y_1, \dots, y_m)$ in \mathfrak{R}^m , let $d_{L^1}(\vec{x}, \vec{y}) = \frac{1}{m} \sum_{i=1}^m |x_i - y_i|$. Thus d_{L^1} is the L^1 distance metric. Let T be a set of points that lie in a bounded region of \mathfrak{R}^m . For any $\epsilon > 0$, an ϵ -cover for T is a finite set $N \subset \mathfrak{R}^m$ (not necessarily contained in T) such that for all $\vec{x} \in T$ there is a $\vec{y} \in N$ with $d_{L^1}(\vec{x}, \vec{y}) \leq \epsilon$. The function $\mathcal{N}(\epsilon, T)$ denotes the size of the smallest ϵ -cover for T . We refer to $\mathcal{N}(\epsilon, T)$ as a *covering number*.

Theorem 1 Let \mathbf{F} be a finite set of functions on Z with $0 \leq f(z) \leq M$ for all $f \in \mathbf{F}$ and $z \in Z$. Let $\vec{z} = (z_1, \dots, z_m)$ be a sequence of m examples drawn independently from Z according to any distribution on Z , and let $\epsilon > 0$. Then

$$\Pr \left(\exists f \in \mathbf{F} : |\widehat{\mathbf{E}}_{\vec{z}}(f) - \mathbf{E}(f)| > \epsilon \right) \leq 2|\mathbf{F}|e^{-2\epsilon^2 m/M^2}.$$

For $0 < \delta \leq 1$ and sample size

$$m \geq \frac{M^2}{2\epsilon^2} \left(\ln |\mathbf{F}| + \ln \frac{2}{\delta} \right)$$

this probability is at most δ . Further, for any $\nu > 0$ and $0 < \alpha < 1$,

$$\Pr \left(\exists f \in \mathbf{F} : d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha \right) \leq 2|\mathbf{F}|e^{-\alpha^2 \nu m/M}.$$

For $0 < \delta \leq 1$ and sample size

$$m \geq \frac{M}{\alpha^2 \nu} \left(\ln |\mathbf{F}| + \ln \frac{2}{\delta} \right)$$

this probability is at most δ . \square

Proof. For the second part of the theorem, using Bernstein's inequality (see e.g. [Pol84]) it is easy to show that for any single function f with $0 \leq f \leq M$,

$$\Pr \left(d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha \right) < 2e^{-\alpha^2 \nu m/M}.$$

Details are given in Lemma 9, part (2) in the Appendix. It follows that the probability that there is any $f \in \mathbf{F}$ with $d_\nu(\widehat{\mathbf{E}}_{\vec{z}}(f), \mathbf{E}(f)) > \alpha$ is at most $2|\mathbf{F}|e^{-\alpha^2 \nu m/M}$. Setting this bound to δ and solving for m gives the result on the sample size. The proof of the first part of the lemma is similar, except we use Hoeffding's inequality (see e.g. [Pol84]), which implies that for any single f ,

$$\Pr \left(|\widehat{\mathbf{E}}_{\vec{z}}(f) - \mathbf{E}(f)| > \epsilon \right) \leq 2e^{-2\epsilon^2 m/M^2}.$$

\square

By letting $\mathbf{F} = l_{\mathcal{H}}$, this theorem can be used in conjunction with Lemma 1 from the previous section to obtain bounds on the sample complexity of learning algorithms that minimize empirical risk. Here we can use either the L_ϵ or $L_{\alpha, \nu}$ family of regret functions. In the former case we get a sample complexity

$$m(\epsilon, \delta) = O \left(\frac{M^2}{\epsilon^2} \left(\log |l_{\mathcal{H}}| + \log \frac{1}{\delta} \right) \right). \quad (1)$$

size needed so that

$$\Pr \left(\exists f \in \mathbf{F} : |\hat{\mathbf{E}}_z(f) - \mathbf{E}(f)| > \epsilon \right) < \delta$$

for $\epsilon, \delta > 0$ [Pol84] [Dud84] [Vap82]. Vapnik also obtains better bounds in some important cases by considering the relative deviation of empirical estimates from true expectations. He looks at bounds on the sample size needed so that

$$\Pr \left(\exists f \in \mathbf{F} : \frac{\mathbf{E}(f) - \hat{\mathbf{E}}_z(f)}{\mathbf{E}(f)} > \epsilon \right) < \delta$$

and also bounds on the sample size needed so that

$$\Pr \left(\exists f \in \mathbf{F} : \frac{\mathbf{E}(f) - \hat{\mathbf{E}}_z(f)}{\sqrt{\mathbf{E}(f)}} > \epsilon \right) < \delta.$$

(Anthony and Shawe-Taylor also obtain bounds of the latter form [AST90].) Note that these are one-sided bounds, in that they only bound the probability that the empirical mean is significantly smaller than the true mean. While extremely useful, as we mentioned in the previous section, these measures of deviation suffer from a discontinuity at $\mathbf{E}(f) = 0$, and lack of convenient metric properties. As in [Pol86], we will give bounds on the sample size needed so that

$$\Pr \left(\exists f \in \mathbf{F} : d_\nu(\hat{\mathbf{E}}_z(f), \mathbf{E}(f)) > \alpha \right) = \Pr \left(\exists f \in \mathbf{F} : \frac{|\hat{\mathbf{E}}_z(f) - \mathbf{E}(f)|}{\nu + \hat{\mathbf{E}}_z(f) + \mathbf{E}(f)} > \alpha \right) < \delta,$$

i.e. the deviation measured using the d_ν metric¹⁰. By setting ν and α appropriately, we obtain results similar to those of [Pol84] and [Vap82] as special cases of our main theorem. However, our results are restricted to the case that all functions in \mathbf{F} are positive and uniformly bounded.

3.1 The case of finite \mathbf{F}

Before considering the general case, it is useful to see what bounds we can get in the case that \mathbf{F} is a finite set of functions. Here we can easily prove the following.

¹⁰In [Pol86], Pollard also gives results that can be used to bound the sample size needed so that

$$\Pr \left(\exists f \in \mathbf{F} : \frac{|\hat{\mathbf{E}}_z(f) - \mathbf{E}(f)|}{\nu + \sqrt{\hat{\mathbf{E}}_z(f)} + \sqrt{\mathbf{E}(f)}} > \alpha \right) < \delta,$$

in analogy with the second type of bound given by Vapnik, except that these bounds are two-sided. We do not pursue these further here.

3 Uniformly good empirical estimates of means

In this section we concentrate on the problem of bounding the number of random examples needed to get good empirical estimates of the risk of each of the decision rules in a decision rule space \mathcal{H} . For each decision rule $h \in \mathcal{H}$ and example $z = (x, y) \in Z$, let $l_h(z) = l(y, h(x))$. As in the previous section, we assume that l is a non-negative bounded loss function taking values in the interval $[0, M]$, thus for each decision rule h , l_h defines a random variable taking values in $[0, M]$. The value of l_h on an example (x, y) is the loss incurred when you use h to determine the action to take for instance x , and the outcome is y . The risk of h is just the expectation of l_h , i.e.

$$\mathbf{r}_h(P) = \mathbf{E}(l_h) = \int_Z l_h(z) dP(z).$$

Furthermore, if $\vec{z} = (z_1, \dots, z_m)$ is a sequence of examples from Z , then the empirical risk of h on \vec{z} is the empirical estimate of the mean of l_h based on the sample \vec{z} , which we denote by $\hat{\mathbf{E}}_{\vec{z}}(l_h)$, i.e.

$$\hat{\mathbf{r}}_h(\vec{z}) = \hat{\mathbf{E}}_{\vec{z}}(l_h) = \frac{1}{m} \sum_{i=1}^m l_h(z_i).$$

Let $l_{\mathcal{H}} = \{l_h : h \in \mathcal{H}\}$. We need to draw enough random examples to get a uniformly good empirical estimate of the expectation of every random variable in $l_{\mathcal{H}}$.

The general problem of obtaining a uniformly good estimate of the expectation of every function in a class \mathbf{F} of real-valued functions has been widely studied (see e.g. [Vap82, Pol84, Dud84] and their references). If no assumptions at all are made about the functions in \mathbf{F} , we immediately run into the problem that some functions in \mathbf{F} could take on arbitrarily large values with arbitrarily small probabilities, making it impossible to obtain uniformly good empirical estimates of all expectations with any finite sample size. This problem can be avoided by making assumptions about the moments of the functions in \mathbf{F} , as in [Vap82], or by assuming that there exists a single non-negative function with a finite expectation (called an *envelope*) that lies above the absolute value of every function in \mathbf{F} , as in [Pol84, Dud84]. In our case, when the loss takes only values in $[0, M]$, then the constant function M serves as an envelope. This case is especially nice since this same envelope works for all distributions on the domain Z of the functions in \mathbf{F} .

The usual measure of deviation of empirical estimates from true means is simply the absolute value of the difference. Thus we would say that the empirical estimates for the expectations of the functions in \mathbf{F} converge uniformly to the true expectations if as the size of the random sample \vec{z} grows,

$$\Pr \left(\exists f \in \mathbf{F} : |\hat{\mathbf{E}}_{\vec{z}}(f) - \mathbf{E}(f)| > \epsilon \right)$$

goes to zero for any $\epsilon > 0$. (This is called (*uniform*) *convergence in probability*, see e.g. [Bil86]). Vapnik, Dudley, Pollard and others have obtained general bounds on the sample

the empirical risk estimates of decision rules in \mathcal{H} converge *uniformly* to the true risk if for all ϵ and $\delta > 0$ there exists a sample size m such that when the $z_i \in \vec{z}$, $1 \leq i \leq m$, are drawn independently at random from Z according to the distribution P , with probability at least $1 - \delta$, we have $\rho(\hat{\mathbf{r}}_h(\vec{z}), \mathbf{r}_h(P)) \leq \epsilon$ for all $h \in \mathcal{H}$. Here ρ is some metric on \mathbb{R}^+ , e.g. either the absolute difference or the d_ν metric.

The following result shows that uniform convergence of the empirical risk estimates, along with a learning method \mathcal{A} that gives a randomized solution to the optimization problem on the estimates, gives a solution to the basic learning problem. We state it for the d_ν metric, but the same argument works also for the absolute difference metric.

Lemma 1 *Let $\nu > 0$ and $0 < \alpha, \delta < 1$. Suppose the sample size $m = m(\alpha, \nu, \delta)$ is such that for all probability distributions $P \in \mathcal{P}$*

$$\Pr(\exists h \in \mathcal{H} : d_\nu(\hat{\mathbf{r}}_h(\vec{z}), \mathbf{r}_h(P)) > \alpha/3) \leq \delta/2,$$

where the $z_i \in \vec{z}$, $1 \leq i \leq m$, are drawn independently at random from Z according to the distribution P . Suppose also that the algorithm \mathcal{A} is such that for all $P \in \mathcal{P}$

$$\Pr(d_\nu(\hat{\mathbf{r}}_{\mathcal{A}(\vec{z})}(\vec{z}), \hat{\mathbf{r}}^*(\vec{z})) > \alpha/3) \leq \delta/2,$$

where \vec{z} is drawn randomly by P as above. Then for all $P \in \mathcal{P}$

$$\Pr(d_\nu(\mathbf{r}_{\mathcal{A}(\vec{z})}(P), \mathbf{r}^*(P)) > \alpha) \leq \delta,$$

i.e. \mathcal{A} solves the basic learning problem for the family of $L_{\alpha, \nu}$ regret functions and has sample complexity at most $m(\alpha, \nu, \delta)$.

Proof. By the triangle inequality for d_ν , if

1. $d_\nu(\mathbf{r}_{\mathcal{A}(\vec{z})}(P), \hat{\mathbf{r}}_{\mathcal{A}(\vec{z})}(\vec{z})) \leq \alpha/3$,
2. $d_\nu(\hat{\mathbf{r}}_{\mathcal{A}(\vec{z})}(\vec{z}), \hat{\mathbf{r}}^*(\vec{z})) \leq \alpha/3$, and
3. $d_\nu(\hat{\mathbf{r}}^*(\vec{z}), \mathbf{r}^*(P)) \leq \alpha/3$,

then

$$d_\nu(\mathbf{r}_{\mathcal{A}(\vec{z})}(P), \mathbf{r}^*(P)) \leq \alpha.$$

The second assumption of the lemma states that (2) holds with probability at least $1 - \delta/2$. The first assumption implies that both (1) and (3) hold with probability at least $1 - \delta/2$. (If (3) fails then we can find a decision rule $h \in \mathcal{H}$ such that $d_\nu(\hat{\mathbf{r}}_h(\vec{z}), \mathbf{r}_h(P)) > \alpha/3$. Here we use the compatibility of d_ν with the ordering on the reals.) Hence with probability at least $1 - \delta$ all of (1) - (3) hold. The result follows. \square

In statistics, this type of result is called a *consistency theorem* about the “statistic” (i.e. the decision rule) computed by the learning method \mathcal{A} . This use of the term “consistency” differs sharply from that common in PAC learning research.

for $L = L_\epsilon$ and $L = L_{\alpha,\nu}$. When $L = L_\epsilon$, this condition means that given m random training examples drawn according to P , with probability at least $1 - \delta$, the decision rule \hat{h} produced by the algorithm \mathcal{A} satisfies

$$\mathbf{r}_{\hat{h},l}(P) \leq \mathbf{r}_l^*(P) + \epsilon,$$

i.e. the risk of \hat{h} is at most ϵ greater than that of the optimal decision rule in \mathcal{H} . When $L = L_{\alpha,\nu}$, this condition is the same, except that we require

$$\mathbf{r}_{\hat{h},l}(P) \leq \frac{1 + \alpha}{1 - \alpha} \mathbf{r}_l^*(P) + \frac{\alpha\nu}{1 - \alpha}.$$

Thus in the former case, the sample complexity is defined in terms of small additive deviations from optimality, and in the latter, we allow both additive and multiplicative deviations. These deviations are controlled by the parameters α and ν .

For example, when $\mathbf{r}_l^*(P) = 0$ as in the standard PAC model, then setting $\alpha = 1/2$ and $\nu = \epsilon$ makes the L_ϵ and $L_{\alpha,\nu}$ conditions equivalent; each reduces to the PAC condition

$$\mathbf{r}_{\hat{h},l}(P) \leq \epsilon.$$

When $\mathbf{r}_l^*(P) > 0$, then $L_{\alpha,\nu}$ condition approximates the L_ϵ condition when α is small and $\nu \approx \epsilon/\alpha$. In particular, since we are assuming that the underlying loss function l is bounded between 0 and M , we have $0 \leq \mathbf{r}_{\hat{h},l}(P), \mathbf{r}_l^*(P) \leq M$ and property (3) of the d_ν metric shows that the $L_{\alpha,\nu}$ condition with $\nu = 2M$ and $\alpha = \epsilon/4M$ implies the L_ϵ condition. This shows how the two parameter $L_{\alpha,\nu}$ condition is generally more flexible than the single parameter L_ϵ condition.

2.5 Relation between learning and optimization

Let us assume that the underlying loss function l is fixed, and we are using either the L_ϵ or $L_{\alpha,\nu}$ regret functions derived from l . In order to solve a basic learning problem, we must find, with high probability, a decision rule \hat{h} with risk close to optimal. As the true distribution P is unknown, to do this we must rely on estimates of $\mathbf{r}_h(P)$ for the various $h \in \mathcal{H}$ which are derived from the given random training sample. For a given $h \in \mathcal{H}$ and training sample $\vec{z} = (z_1, \dots, z_m)$, where $z_i = (x_i, y_i) \in Z$, let $\hat{\mathbf{r}}_h(\vec{z})$ denote the empirical risk on \vec{z} , i.e. $\hat{\mathbf{r}}_h(\vec{z}) = \frac{1}{m} \sum_{i=1}^m l(y_i, h(x_i))$. Let $\hat{\mathbf{r}}^*(\vec{z}) = \inf\{\hat{\mathbf{r}}_h(\vec{z}) : h \in \mathcal{H}\}$. We can then define a natural optimization problem associated with the basic learning problem: given the training sample \vec{z} , find a decision rule $\hat{h} \in \mathcal{H}$ such that $\hat{\mathbf{r}}_{\hat{h}}(\vec{z})$ is close to $\hat{\mathbf{r}}^*(\vec{z})$, i.e. a decision rule whose empirical risk on the training sample is close to minimal.

Solving the optimization problem does not automatically solve the learning problem. We need to have good empirical risk estimates as well. Since l is bounded, for every $h \in \mathcal{H}$, as the sample size $m \rightarrow \infty$, $\hat{\mathbf{r}}_h(\vec{z}) \rightarrow \mathbf{r}_h(P)$ with probability 1. We will say that

ble states of nature P in \mathcal{P} , or do we want to assume a *prior distribution* on possible distributions in \mathcal{P} , so that we can define a notion of “average case” big ‘L’ risk to be minimized. The former goal is known as *minimax optimality*, and has been used in the PAC model. The latter is the *Bayesian* notion of optimality [Ber85,Kie87], and has been used in several approaches to learning in neural nets based on statistical mechanics [DSW⁺87,TLS89,GT90,STS90,OH91a,OH91b]. Unfortunately this last question has no clear cut answer, and leads us directly into a longstanding unresolved debate in statistics (see e.g. [Lin90] and following discussion.). Since we have set out to generalize the PAC model, and since our results are best illustrated in the minimax setting, we will formalize the notion of a basic learning problem using the minimax criterion. In subsequent work we hope to further explore this Bayesian setting. (For recent work in Bayesian approaches to neural network learning see [Mac92,BW91], and for Bayesian versions of the PAC model see [HKS91,Bun90].)

We can now define exactly what we mean by a basic learning problem, and what it means for a learning method to solve this problem in this minimax setting.

Definition 1 *A basic learning problem is defined by six components $X, Y, A, \mathcal{H}, \mathcal{P}$, and \mathcal{L} , where the first five components are as defined in section 2.1, and the last component, \mathcal{L} , is a family of regret functions as defined in section 2.3 (e.g. $\mathcal{L} = \{L_{\alpha,\nu} : \nu > 0 \text{ and } 0 < \alpha < 1\}$, or $\mathcal{L} = \{L_\epsilon : \epsilon > 0\}$ for some loss function l). Let \mathcal{A} be a learning method as defined in section 2.3. We say that \mathcal{A} solves the basic learning problem if for all $L \in \mathcal{L}$ and all $0 < \delta < 1$ there exists a finite sample size $m = m(L, \delta)$ such that*

$$\text{for all } P \in \mathcal{P}, R_{L,\mathcal{A},m}(P) \leq \delta.$$

The sample complexity of the learning method \mathcal{A} is the smallest such integer valued function $m(L, \delta)$. When $\mathcal{L} = \{L_{\alpha,\nu} : \nu > 0 \text{ and } 0 < \alpha < 1\}$ we will denote $m(L_{\alpha,\nu}, \delta)$ by $m(\alpha, \nu, \delta)$ and when $\mathcal{L} = \{L_\epsilon : \epsilon > 0\}$ we will denote $m(L_\epsilon, \delta)$ by $m(\epsilon, \delta)$.

As discussed above, this definition generalizes the PAC criterion, and several others as well. In fact, this definition is quite generous, in that sample size needed to get the big ‘L’ risk less than δ is only required to be finite for each $\delta > 0$. In particular, using property (3) of the d_ν metric from section 2.2, when the underlying loss function l is bounded, as we assume here, any algorithm \mathcal{A} solves the basic learning problem using the $L_{\alpha,\nu}$ class of regret functions if and only if it solves it using the L_ϵ class. Thus it doesn’t matter which of these two classes of regret functions we use. However, in practice it is the sample complexity of \mathcal{A} that is critical, and this will depend on which class of regret functions are used.

The nature of this dependence is seen more clearly when we expand the condition

$$R_{L,\mathcal{A},m}(P) \leq \delta$$

Hence demanding big ‘L’ risk at most δ gives the usual PAC criterion that the risk (or “error”) of the decision rule (or “hypothesis”) produced by \mathcal{A} be greater than ϵ with probability at most δ .

The regret function L can also be defined similarly, but using the d_ν metric to measure distance from optimality, instead of the absolute difference. Specifically, for every $\nu > 0$ and $0 < \alpha < 1$, we can define the regret function $L_{\alpha,\nu}$ by letting $L_{\alpha,\nu}(P, h) = 1$ if $d_\nu(\mathbf{r}_{h,i}(P), \mathbf{r}_i^*(P)) > \alpha$, and $L_{\alpha,\nu}(P, h) = 0$ otherwise. In this case the big ‘L’ risk $R_{L_{\alpha,\nu}, \mathcal{A}, m}(P)$ measures the probability that the risk of the decision rule produced by the algorithm \mathcal{A} has distance more than α from optimal in the d_ν metric, when the algorithm is given m random training examples drawn according to P . We will see in sections 2.4 and 3.1 why this sometimes gives a more useful and flexible definition of regret. In fact, in this paper, we will give our main results in terms of the family $\{L_{\alpha,\nu} : \nu > 0 \text{ and } 0 < \alpha < 1\}$ of regret functions, and show how corresponding results may be derived as corollaries for the family $\{L_\epsilon : \epsilon > 0\}$ of regret functions.

Other regret functions are also possible and lead to different learning criteria. For example, another, perhaps simpler, way to define regret is to let $L(P, h) = \mathbf{r}_{h,i}(P) \Leftrightarrow \mathbf{r}_i^*(P)$. When $\mathbf{r}_i^*(P) = 0$, as it does in the standard noise-free PAC model, this definition makes the regret L equal to the risk $\mathbf{r}_i(P)$, i.e. the expectation of the underlying loss l . In this case the big ‘L’ risk $R_{L, \mathcal{A}, m}(P)$ measures the expectation of the loss incurred by the learning algorithm \mathcal{A} when it is given m random training examples drawn according to P , forms a decision rule h , and then uses h to determine the action on one further independent random example drawn according to P . This gives a generalization of the learning criterion studied in [HLW90]. When $\mathbf{r}_i^*(P) \neq 0$, then the big ‘L’ risk gives the expectation of the amount of such loss above and beyond the expected loss that would be suffered if the optimal decision rule were used. In particular, in density estimation, where P and h are both densities on the instance space X , if $P \in \mathcal{H}$ then defining the regret by $L(P, h) = \mathbf{r}_{h,i}(P) \Leftrightarrow \mathbf{r}_i^*(P)$ makes it equal to the Kullback-Leibler divergence from P to h . Hence the big ‘L’ risk is the expected Kullback-Leibler divergence of the decision rule h returned by the algorithm from the true density (see sections 1.1.3 and 1.1.4).

It is also possible to define the regret function L directly, without using an underlying loss function l . For example, in density estimation it is possible to use other measures of the distance between two densities, e.g. the Hellinger distance or the total variational distance, as in [BC90, Yam90]. The criterion from [KS90] for inferring a good *model of probability* can also be defined using an appropriate regret function, without defining an underlying loss l .

2.4 Full formalization of the basic learning problem

Having defined the regret function, and thereby the big ‘L’ risk function, we still face one last issue: do we want to minimize big ‘L’ risk in the worst case over all possi-

2.3 The regret function L and the big ‘L’ risk R

Once we have specified how we measure closeness to optimality, we still need to specify our criteria for a successful learning algorithm. Do we need to have the risk of the decision rule found close to the optimum $\mathbf{r}_l^*(P)$ with high probability, or should its average distance from $\mathbf{r}_l^*(P)$ be small? Do we measure success in terms of the performance of the algorithm on the worst case distribution in \mathcal{P} , or do we use some average case analysis over distributions in \mathcal{P} ? These questions lead us right back to decision theory again, but this time at a higher level in the analysis of learning.

To see this, consider the structure of a learning algorithm \mathcal{A} . For any sample size m , the algorithm \mathcal{A} may be given a sample $\vec{z} = ((x_1, y_1), \dots, (x_m, y_m))$ drawn at random from Z^m according an unknown product distribution P^m , where $P \in \mathcal{P}$. For any such \vec{z} it will choose a decision rule $\mathcal{A}(\vec{z}) \in \mathcal{H}$. Thus abstractly, the algorithm defines a function \mathcal{A} from the set of all samples over Z into \mathcal{H} , i.e. $\mathcal{A} : \bigcup_{m \geq 1} Z^m \rightarrow \mathcal{H}$. Since we are not requiring computability here, we will call such \mathcal{A} a *learning method*. When $P \in \mathcal{P}$ is the actual “state of nature” governing the generation of examples, and the algorithm produces the decision rule $h \in \mathcal{H}$, let us say that we suffer a nonnegative real-valued *regret* $L(P, h)$. Thus, formally $L : \mathcal{P} \times \mathcal{H} \rightarrow \mathbb{R}^+$. In our treatment here, the regret function L will be derived from the loss function l , and will measure the extent to which we have failed to produce a near optimal decision rule, assuming P is the true state of nature (i.e. the amount of “regret” we feel for not having produced the optimal decision rule). Finally, for each possible state of nature P , the average regret suffered by the algorithm, over all possible training samples $\vec{z} \in Z^m$, is the *big ‘L’ risk* of that algorithm under P for sample size m . This big ‘L’ risk is defined formally by

$$R_{L, \mathcal{A}, m}(P) = \int_{\vec{z} \in Z^m} L(P, \mathcal{A}(\vec{z})) dP^m(\vec{z}).$$

The goal of learning is to minimize big ‘L’ risk.

We illustrate these definitions with a few examples. First suppose we want to capture the notion of successful learning that is used in the PAC model. Then one possibility is to introduce an *accuracy parameter* $\epsilon > 0$ and define the regret function $L = L_\epsilon$ by letting $L_\epsilon(P, h) = 1$ if $\mathbf{r}_{h,l}(P) \Leftrightarrow \mathbf{r}_l^*(P) > \epsilon$, and $L_\epsilon(P, h) = 0$ otherwise. This we suffer regret only when the decision rule h produced by the learning algorithm has risk that is more than ϵ from optimal, measured by the absolute difference metric. For this definition of regret, the big ‘L’ risk $R_{L_\epsilon, \mathcal{A}, m}(P)$ measures the probability that the decision rule produced by \mathcal{A} has risk more than ϵ from optimal, when \mathcal{A} is given m random training examples drawn according to P . We then demand that this big ‘L’ risk be small, i.e. smaller than some given *confidence parameter* $\delta > 0$.

In the PAC model it is commonly assumed that the examples given to the algorithm \mathcal{A} are noise-free examples of some underlying target function $f \in \mathcal{H}$. In this case the risk $\mathbf{r}^*(P)$ of the optimal decision rule in \mathcal{H} is zero, and hence $L_\epsilon(P, h) = 1 \Leftrightarrow \mathbf{r}_{h,l}(P) > \epsilon$.

according to P . It is defined by

$$\mathbf{r}_{h,l}(P) = \mathbf{r}_h(P) = \mathbf{E}(l(y, h(x))) = \int_Z l(y, h(x)) dP(x, y)$$

(the subscript l will be omitted when the loss function is clear from the context.) Since l is bounded this expectation is finite for every distribution P . In decision theory the expected loss $\mathbf{r}_h(P)$ is called the *risk* of h when P is the true underlying distribution. This quantity generalizes the notion of the *error* of h used in computational learning theory.

In section 1.1 we stated the goal of learning quite informally: Given examples chosen independently at random from some unknown probability distribution $P \in \mathcal{P}$, find a decision rule \hat{h} in \mathcal{H} that comes “close to” minimizing the risk $\mathbf{r}_h(P)$ over all $h \in \mathcal{H}$. Let $\mathbf{r}_l^*(P)$ (or $\mathbf{r}^*(P)$ when l is clear from the context) denote the infimum of $\mathbf{r}_h(P)$ over all h in the decision rule space \mathcal{H} . To formalize our notion of a basic learning problem, we first need to say what we mean that $\mathbf{r}_{\hat{h}}(P)$ is “close to” $\mathbf{r}^*(P)$.

Let $r = \mathbf{r}_{\hat{h}}(P)$ and $s = \mathbf{r}^*(P)$. One natural interpretation is to demand that $|r \leftrightarrow s| \leq \epsilon$ for some small $\epsilon > 0$. However, we will see in section 3.1 that sometimes it is better to use a relative measure of distance. For any real $\nu > 0$, let d_ν be the function defined by

$$d_\nu(r, s) = \frac{|r \leftrightarrow s|}{\nu + r + s}$$

for any non-negative reals r and s . It is straightforward but tedious to verify that d_ν is a metric on \mathfrak{R}^+ . The d_ν metric is similar to the standard function

$$\frac{|r \leftrightarrow s|}{s}$$

used to measure the difference between the quality r of a given solution and the quality s of an optimal solution in combinatorial optimization. However, our measure has been modified to be well-behaved when one or both of its arguments are zero, and to be symmetric in its arguments (so that it is a metric). Three other properties of d_ν are also useful.

1. For all non-negative reals r and s , $0 \leq d_\nu(r, s) < 1$.
2. For all non-negative $r \leq s \leq t$, $d_\nu(r, s) \leq d_\nu(r, t)$ and $d_\nu(s, t) \leq d_\nu(r, t)$.
3. For $0 \leq r, s \leq M$, $\frac{|r-s|}{\nu+2M} \leq d_\nu(r, s) \leq \frac{|r-s|}{\nu}$.

We will refer to the second property by saying that d_ν is *compatible with the ordering on the reals*.

The fifth component, \mathcal{P} , is a family of joint probability distributions on $X \times Y$. These represent the possible “states of nature” that might be governing the generation of examples. The set $Z = X \times Y$ will be called the *sample space*. We assume that examples are drawn independently at random according to some probability distribution $P \in \mathcal{P}$ on the sample space Z . A sequence of examples will be called a *sample*. In what follows⁸ we will usually assume that \mathcal{P} includes all probability distributions on Z . Hence our results will be distribution independent.

The last component, the loss function l , is a mapping from $Y \times A$ into \mathfrak{R} . In this paper we will assume that l is bounded and nonnegative, i.e. $0 \leq l \leq M$ for some real M . When Y and A are finite it is always possible to enforce this condition by simply adding a constant to l , which doesn’t change the learning problem in any essential way. When either Y or A is infinite, the learning problem sometimes needs to be restricted to meet this condition. For example, in regression⁹ we might restrict the possible parameter vectors in A and/or the possible outcomes in Y such that for every $y \in Y$ and $a \in A$, $\hat{P}(y; a) \geq b$ for some constant b . We can then take $M = -\log b$. In density estimation, the same thing can be accomplished by restricting the instance space X to a bounded subset of \mathfrak{R}^n on which all densities in \mathcal{H} have values uniformly greater than b and less than B for constants $0 < b < B$. We can then add $\log B$ to the loss function to make it positive. The same method works for estimating distributions on discrete spaces: we restrict ourselves to a finite instance space X and demand that for all $x \in X$ and all probability distributions $h \in \mathcal{H}$, $h(x) \geq b > 0$ (see e.g. [AW90, Yam90]). These restrictions are often reasonable in practice, e.g. most measurements naturally have bounded ranges, but they can be annoying (see [Vap89], [Pol84], and [Pol90] for alternative approaches for unbounded loss functions).

2.2 Measuring distance from optimality with the d_ν metric

For a given decision rule $h \in \mathcal{H}$ and distribution P on the sample space Z , the expected loss of h is the average value of $l(y, h(x))$, when the example (x, y) is drawn at random

⁸It is, however, possible and in fact common to assume that \mathcal{P} is a very specific class of probability distributions on Z . For example, let $X = \mathfrak{R}^n$. Then if we are doing classification learning and Y is discrete we may assume that y is selected according to an arbitrary distribution on Y , and for each y , $P(x|y)$ is a multi-variate Gaussian distribution on X [DH73]. On the other hand, if we are doing linear regression, then Y is real-valued and we might assume that x is selected according to an arbitrary distribution on X , and y is a linear function of x with additive Gaussian noise. In PAC learning theory we have a discrete analog of the latter case. Here we usually have $X = \{0, 1\}^n$, $Y = \{0, 1\}$, and y a Boolean function of x of a particular type (e.g. defined by a small disjunctive normal form formula), possibly plus random noise.

⁹Note that to get bounded loss in linear regression, X must a bounded subset of \mathfrak{R}^n as well, since we can’t bound Y without bounding X . The coefficients of the functions in \mathcal{H} must also be bounded.

X, Y, A, \mathcal{H} and l	sections 1.1 and 2.1
\mathcal{P}	section 2.1
$\mathbf{r}_{h,l}(P), \mathbf{r}_h(P)$ (true risk)	section 2.2
$\mathbf{r}_l^*(P), \mathbf{r}^*(P)$ (optimal risk)	section 2.2
$\hat{\mathbf{r}}_h(\vec{z})$ (empirical risk)	section 2.2
$\hat{\mathbf{r}}^*(\vec{z})$ (optimal empirical risk)	section 2.2
d_ν	section 2.2
$L, L_\epsilon, L_{\alpha,\nu}$ (regret functions)	section 2.3
R (big ‘L’ risk)	section 2.3
$m(\epsilon, \delta), m(\alpha, \nu, \delta)$ (sample complexity)	section 2.4
\mathcal{N} (covering number)	sections 10.1 and 3.2
\mathcal{M} (packing number)	section 10.1
dim (metric dimension)	section 10.1
dim_P (pseudo dimension)	section 4
\mathcal{C} (capacity)	section 6
ρ_l	section 6
$l_{\mathcal{H}}$	section 3
$\mathbf{E}_{ \mathcal{E}}$	section 3
$\hat{\mathbf{E}}$ (empirical expectation)	section 3
d_{L^1} (L^1 distance for vectors)	section 3.2
$d_{L^1(P)}$ (L^1 distance for functions)	section 4
$d_{L^1(P,\rho)}$ (L^1 distance for functions)	section 6

2 Learning and optimization

We now further formalize the basic problem of learning, as introduced in section 1.1. We will introduce a formal notion of a learning algorithm, and a higher level loss function, which we will call a *regret function*, that measures how well the learning algorithm performs. The regret function will be defined in terms of the low level loss function l discussed in the previous section. Finally, we will show how an algorithm can solve the learning problem by solving a related optimization problem.

2.1 The basic components $X, Y, A, \mathcal{H}, \mathcal{P}$ and l

We first review and further formalize the six components of the basic learning problem introduced in the previous section: $X, Y, A, \mathcal{H}, \mathcal{P}$ and l . The first four components are the instance, outcome, decision and decision rule spaces, respectively. The first three of these are arbitrary sets, and the fourth, \mathcal{H} is a family of functions from X into A . These have been discussed extensively in the previous section.

1.5 Notational conventions

We denote the real numbers by \mathfrak{R} and the non-negative real numbers by \mathfrak{R}^+ . By \log and \ln we denote the logarithm base 2 and the natural logarithm, respectively. We use $\mathbf{E}(\cdot)$ to denote the expectation of a random variable, and $\mathbf{Var}(\cdot)$ to denote the variance of a random variable. When the probability space is defined implicitly from the context, we use $\mathbf{Pr}(\cdot)$ to denote the probability of a set. However, usually the measure on the underlying probability space will be defined explicitly using the symbol P .

Here, P will usually denote a probability measure on some appropriate⁶ σ -algebra over the set $Z = X \times Y$, where X is the instance space and Y is the outcome space. We use P^m to denote the m -fold product measure on Z^m . Functions on Z and subsets of Z mentioned in what follows will be assumed to be measurable without explicit reference. Alternately, we will also view X and Y as random variables on some other, unspecified, probability space, e.g. when they are viewed as real valued measurements. In this case P is viewed as a joint distribution on X and Y . In either case, the probability of a set $T \subset Z$ is defined by

$$P(T) = \int_T dP(z)$$

(where $z = (x, y)$ with $x \in X$ and $y \in Y$) and the expectation of function f on Z is denoted by

$$\mathbf{E}(f) = \int_Z f(z) dP(z).$$

When Z is countable we will, with some abuse of notation, also use P for the probability mass function, i.e. for $z \in Z$, $P(z)$ denotes $P(\{z\})$. Hence $P(T) = \sum_{z \in T} P(z)$ and $\mathbf{E}(f) = \sum_{z \in Z} f(z)P(z)$ in this case. When Z is continuous, a density associated with P (if it exists) is denoted by p .

When Z is countable we use $P(y|x)$ to denote the probability that $Y = y$ given that $X = x$ (viewing X and Y as random variables) and similarly for $P(x|y)$. Hence $P(\cdot|x)$ denotes the conditional distribution on Y , given $X = x$. The marginal distribution in X is defined by⁷ $P_{|X}(x) = \sum_{y \in Y} P(x, y)$. Here and elsewhere, we abbreviate $P((x, y))$ by $P(x, y)$.

Finally, we list some other notation that is used several places in the text, indicating which section it is defined in.

⁶If Z is countable then we assume this σ -algebra contains all subsets of Z , otherwise we assume that Z is a complete, separable metric space (see section 10.1) and that this σ -algebra is the smallest σ -algebra that contains the open sets of Z (i.e. the σ -algebra of Borel sets).

⁷When Z is uncountable, the marginal and conditional distributions are defined so that

$$\int_Z f(x, y) dP(x, y) = \int_X \left(\int_Y f(x, y) dP(y|x) \right) dP_{|X}(x)$$

for every bounded measurable function f .

is not only fascinating from a purely mathematical standpoint, but also potentially very useful in machine learning and other applied fields.

1.4 Organization of the paper

The remainder of the paper is organized as follows. The learning framework we have described above in section 1.1 is defined more formally in section 2. There we also look at the question of evaluating the performance of learning algorithms in terms of the number of training examples they use. This question is also formalized from a decision theory perspective. We then provide a lemma (Lemma 1) that can be used to evaluate the performance of learning algorithms that work by minimizing empirical loss. To use this lemma, we need bounds on the rate of uniform convergence of empirical loss estimates to true expected losses. These are given in section 3. The key bound is given in Theorem 2 in section 3, and in a more general version in Theorem 3.

To use the bound from Theorem 2 we need bounds on the “random covering numbers” associated with the decision rule space \mathcal{H} , the loss function l and the distribution P . These are related to the idea of an ϵ -cover described above. In section 4 we introduce Pollard’s notion of the pseudo dimension as a means of bounding the random covering numbers. Applications of this method to several learning problems are described in section 5.

The techniques of sections 4 and 5 only apply to the case when the action set A is real-valued. Tools for bounding the random covering numbers that apply in more general cases are developed in section 6. Here we introduce the notion of the capacity of the decision rule space \mathcal{H} (for a particular loss function l), and the related notion of the metric dimension of \mathcal{H} . In section 7 we use these notions to obtain bounds on the performance (in terms of the number of training examples used) of learning algorithms that use multilayer feedforward neural networks, and work by minimizing empirical loss (Corollary 3). Finally, some further discussion of our results is given in the conclusion, section 8.

Many of the more technical proofs and definitions have been moved into the appendix to make the paper more readable. The appendix has several sections. Section 10.1 contains a brief overview of the theory of metric spaces, ϵ -covers and metric dimension. Notation from this section is used in several places in the paper. Section 10.2 deals with certain technical measurability requirements. Section 10.3 gives an analogue of Chernoff and Hoeffding bounds using the d_v metric. Section 10.4 contains the proof of Theorem 2. Finally, Section 10.5 contains a result on feedforward neural networks of linear threshold functions that is similar to that given in [BH89], and provides a counterpart to Corollary 3 in section 7.

(MDL) approaches [BC90,Ris86], try to find a decision rule that minimizes some function of empirical loss and decision rule complexity. These can also achieve expected loss approaching that of Bayes optimal decision rule in the limit, and may be more effective in practice. Although uniform convergence results such as those we develop here are also used in the analysis of such methods [Vap82] (and in the analysis of cross-validation methods [NP87]), the full treatment of such approaches is beyond the scope of the present paper. It should also be noted that Bayesian methods and structural risk minimization can be applied even when the decision rule space includes only neural networks of a fixed size. An example is the recent work using weight penalty functions in neural net training [WHR90,LDS90,NH91,Mac92,BW91]. Such approaches may significantly reduce the training sample size needed to avoid overfitting in practice.

1.3 Overview of methods used

We now briefly discuss the methodology and previous work used in obtaining our results. Our work builds directly on the work of Vapnik and Chervonenkis, Pollard, and Dudley on the uniform convergence of empirical estimates [Vap82][Pol84][Dud84] and its application to pattern recognition [Vap82,Vap89] [Dev88]. It also builds on the work of Benedek and Itai on PAC learnability with respect to specific probability distributions [BI88], and is related to the work of Natarajan and Tadepalli on extensions of the VC dimension to multi-valued functions [NT88] [Nat89b] and PAC learnability with respect to classes of probability distributions [Nat88] [Nat89a]. In addition, Quiroz and Kulkarni have each independently generalized the PAC model in a related manner [Qui89,Kul89].

One of the key ideas we use is the notion of an ϵ -cover of a metric space [Dud84] [Pol84] [BI88] [Nat89a] [Qui89] and the associated idea of *metric dimension* [KT61] (also called the *fractal* dimension [Far82]). This notion of dimension has played an important role in the now very active study of fractals in nature [Man82], especially in connection with chaos in dynamical systems [Far82][FOY83]. Here we build further on the beautiful results of Vapnik and Chervonenkis [Vap82], Dudley [Dud78] and Pollard [Pol84], which relate a type of generalized VC dimension for a decision rule space to the number of balls of radius ϵ required to cover the space, with respect to certain metrics. The sizes of the smallest such covers determine the *metric dimension* of the space. Our treatment closely parallels the approach given in [Pol90]. It is interesting to note that related results connecting ϵ -covers with the VC dimension have also been independently developed in [BI88] and in recent computational geometry work [Wel88]⁵. This work seems to lead to a potentially rich area of investigation that combines elements of combinatorics, topology and geometry, and probability and measure in a novel framework. We feel that this area

⁵Specifically, Lemma 7.13 of [Dud78] is nearly equivalent to Lemma 4.1 of [Wel88] (using the primal space instead of the dual). This result also gives a stronger version of Theorem 4, part(3) of [BI88]. We give a still stronger version of this result in Theorem 6 below.

to avoid overfitting when learning with the decision rule space of feedforward neural nets [RM86], extending previous work in [BH89] and [Whi90a] (see also related work in [AST90]). These are the nets most widely used in current neural net learning research. Our model for feedforward neural nets is quite general in that it allows many types of units in the nets, including quasi-linear units [RM86], radial basis units [PG89], and product units [DR89].

In our general setting, successful learning means finding a decision rule with average loss close to minimal over all decision rules in the given decision rule space, rather than loss close to zero as in the PAC model. In addition to using an additive model as in [LMR88], we also define “close to” using a measure of relative difference (the d_v metric) similar to the standard multiplicative measure of approximation used in combinatorial optimization. This allows us to state the relevant uniform convergence bounds as generalized “Chernoff-style” [AV79] bounds, as in [Pol86],[BFOS84] (chapter 12), rather than “Hoeffding-style” bounds (as in Pollard’s results [Pol84]), giving better bounds on sufficient training sample size in some important cases. These two types of bounds are analogous to the two types of bounds that Vapnik gives in his book [Vap82] in that one uses a measure of absolute difference and the other a measure of relative difference. However, both of our bounds are “two-sided”, i.e. they bound deviations both above and below the mean.

We give these upper bounds on required sample size only to give some indication of the order-of-magnitude dependence of sample size on certain critical parameters of the learning problem, and to illustrate the theory. They are still too crude to be used directly in practice, e.g. as explicit formulae for choosing an appropriate sample size. Cross validation techniques, in which some of the training examples are held in reserve and used instead to test the performance of the decision rules produced by the learning algorithm, are likely to perform better for this task in practice (see e.g. [Whi90a,WK91]). Nevertheless, cross validation is only a means of estimating the amount of overfitting in the learning method in particular cases, i.e. it is only an engineering trick and provides no scientific explanation of the phenomenon. Our goal is to understand and explain overfitting in general decision rule spaces, from a scientific rather than an engineering viewpoint.

Finally, we should note that in practice, many learning algorithms do more than just search for a decision rule in a fixed decision rule space that minimizes empirical loss. For example, it is common to let the decision rule space depend on the number of training examples available, using richer and richer decision rule spaces as more examples become available (see e.g. [Whi90a,BEHW89]). This can allow the learning algorithm to produce a sequence of decision rules with expected losses that approach the loss of Bayes optimal decision rule in the limit of infinite training sample size for a large class of possible joint distributions. The results given here can be used to estimate the appropriate rate at which the decision rule space should grow relative to the sample size to avoid overfitting. Other approaches, e.g. the method of *structural risk minimization* introduced by Vapnik [Vap82], and the *Bayesian* [Ber85,Mac92,BW91] and *minimum description length*

outcome spaces are used. In pattern recognition and statistics, the instance space X is usually a finite dimensional real vector space, i.e. each instance consists of a vector of real valued measurements of some attributes. In density estimation, a decision rule represents a density on X , and many choices are possible. One common choice is a mixture of Gaussian densities (e.g. [DH73][Now90]). In standard regression, the outcome and decision spaces Y and A are identical and real valued, and linear functions are most often used as decision rules. For more complex outcome spaces such as those in the medical diagnosis example given above, the decision rule space for regression is usually defined using a *generalized linear model* [MN89]. Similarly, in binary classification, where there are only two possible outcomes in Y as in the PAC model, linear threshold functions are most often used as decision rules, and there are straightforward generalizations for the case of k -ary classification (see e.g. [DH73]). This “linear bias” in pattern recognition and statistics is in contrast to that in the PAC model and other AI areas, including work in neural networks, in which a rich variety of decision rule spaces are used (see e.g. [Tou89,Tou90,Hau88,Hau89]). Our main goal here is to develop analytic tools to help understand the problem of overfitting in these more complex decision rule spaces.

In order to focus on the problem of overfitting, we take a simplified view of learning, in which the learner chooses a decision rule space \mathcal{H} , and then tries to find a decision rule in \mathcal{H} with near minimal expected loss. To do this, the learner looks for a decision rule that minimizes the observed average loss on the training examples, which is called *empirical loss* or *empirical risk*. For example, in standard linear regression⁴ the learning algorithm is the method of least squares, i.e. we find the linear function h that minimizes the average of $l(y, h(x)) = (h(x) - y)^2$ over all examples (x, y) in our training set. It is well known that if we have too few training examples, then we tend to overfit them, and the function we find does not come close to minimizing the actual expected quadratic loss, which would be obtained by integrating over all possible (mostly unseen) examples with respect to the unknown joint distribution on them. This same situation occurs with all nontrivial decision rule spaces, including the nonlinear regression models defined by feedforward neural nets.

Using certain measures of the “dimension” or “capacity” of the decision rule space \mathcal{H} and classes derived from \mathcal{H} (see below), we obtain general upper bounds on the number of random training examples needed so that with high probability, any decision rule in \mathcal{H} that has small empirical loss on the training examples will have small actual expected loss, i.e. we get uniform convergence results for empirical estimates like those in [Vap82][Dud84][Pol84,Pol90]. We show how these give upper bounds on sufficient training sample size like those derived in [BEHW89] and elsewhere using the notion of the VC dimension, and generalize those results.

As an application, we give specific bounds on the number of training examples needed

⁴For general regression with the negative log likelihood loss function, the principle of minimizing empirical loss is the same as the principle of *maximum likelihood* [Ber85,Kie87].

drawn randomly from some density $p(x)$ on X . Let the decision set A be the positive real numbers and each decision rule h in \mathcal{H} be a density on X . Then, as above, information theoretic considerations suggest the the loss function $l(y, a) = l(a) = -\log a$. Again, as above, the expected loss of h is minimized when h is the true density p . Further, if p is not a member of \mathcal{H} , then the best decision rule in \mathcal{H} , in terms of minimizing the expected loss, is the one with the smallest Kullback-Leibler divergence from the true density p [Kul59].

When the instance space X is discrete, we are not estimating a density on X but rather a probability distribution. The same ideas as above carry over, except that we let the decision space $A = (0, 1)$ and each decision rule h in \mathcal{H} represent a probability distribution on X . Here we can also use the same loss function, and it has the same properties.

These examples illustrate the diversity of the learning problems that can be cast in the proposed decision theoretic framework, even under the restrictive assumptions we make here, i.e. that the outcome y does not depend on the action a , and that the learner always observes both the outcome and the loss. By weakening these assumptions, we can model other types of learning as well, including *associative reinforcement learning* [BA85,Gul90] and the theory of *learning automata* (with static environment) [NT89]. However, we will not pursue this here.

1.2 Summary and discussion of the results presented here

There are three major practical issues in this decision theoretic view of learning. The first is the number of random examples needed in order to be able to produce a good decision rule in the decision rule space \mathcal{H} , i.e. a decision rule whose expected loss is near the minimum of all decision rules in \mathcal{H} . If too few examples are used, we run into the problem of *overfitting*, where the decision rule produced performs well on the training data, but not on further random examples drawn from the same joint distribution that generated this training data. The second is the adequacy of the decision rule space \mathcal{H} . If \mathcal{H} does not contain any decision rule with expected loss close to that of Bayes optimal decision rule for the particular joint distribution we are dealing with, then we can never hope to achieve near optimal performance using this decision rule space. Choosing the right decision rule space often requires considerable insight into the particular problem domain. Finally, the third practical problem is the computational complexity of the method we use to produce our decision rule from the training examples. This issue has been addressed extensively in the PAC literature, and is also addressed in [KS90,AW90]. Of these three important issues, here we examine only the first. This issue is referred to as the problem of estimating the “sample complexity” of the learning problem in the PAC literature [EHKV89].

The number of random training examples needed to avoid overfitting depends critically on the nature of the decision rule space used. Different kinds of decision rule spaces are used in different areas of learning research, partly because different kinds of instance and

and the log likelihood loss is

$$l(y, a) = - \sum_{i=1}^k (y_i \log a_i + (1 - y_i) \log(1 - a_i)),$$

which we will call the *cross entropy loss*.

In the medical diagnosis example, the outcome space Y is discrete. However, in most uses of regression Y is real valued, e.g. the outcome y is the measurement of some real valued quantity, and the instance x represents the experimental conditions under which this quantity was measured. In this case regression is usually defined as estimating the conditional expectation of Y given the instance x . Thus $A \subset \mathfrak{R}$, and the action $a \in A$ for a given instance x consists of an estimate of the mean of the various outcomes y that would typically be observed for that instance x . It is easy to show that by using the *quadratic* loss function $l(y, a) = (a - y)^2$, the expected loss is minimized when a is the true mean, and hence this version of regression also fits naturally³ into the decision theoretic framework. An alternate approach is to use the L_1 loss function $l(y, a) = |a - y|$, in which case the expected loss is minimized when a is the median of the conditional distribution Y given the instance x . (See e.g. [Whi90b], [Hau90].)

1.1.4 Density and parameter estimation

Finally, the problems of parameter estimation and density estimation can also be viewed as special cases of this decision theoretic framework. For parameter estimation, note that when the instance space X has only one element then the particular instance x can be ignored entirely. Thus the regression problem reduces to the problem of estimating the parameters of a single distribution on the outcome space Y from a sample of random outcomes y from Y , i.e. to the simpler problem of parameter estimation. Here the decision rule is not a function but merely a single vector of parameters, and the decision rule space \mathcal{H} is the same as the decision space A .

For density estimation, we can consider the dual case in which the outcome space Y has only one element, and hence can be ignored. Thus examples are unlabeled instances x

³In fact, the standard version of regression, defined as estimating the conditional mean of Y given instance x using the quadratic loss function, is actually a special case of the general version of regression defined above, where for continuous outcome spaces Y , the object is to estimate the parameters specifying the conditional density of Y given instance x , using the log likelihood loss function. To see this, assume that we represent the conditional density on Y with a Gaussian density $\hat{p}(y; \mu, \sigma) = (2\pi\sigma^2)^{-1/2} e^{-(\mu-y)^2/2\sigma^2}$, where μ is the mean and σ^2 the variance. Let the variance be fixed, independent of x , so that the estimate $\hat{p}(y; \mu, \sigma)$, of the conditional density on Y given x is completely determined by the mean μ . Thus the decision space $A \subset \mathfrak{R}$, and each action a in A is interpreted as specifying the mean of a Gaussian density. Substituting $\mu = a$ and evaluating $-\log \hat{p}(y; \mu, \sigma)$, the log likelihood loss is seen to be $l(y, a) = \frac{1}{2\sigma^2}(a - y)^2 + \frac{1}{2} \log(2\pi\sigma^2)$. For fixed variance σ^2 , this is equivalent, for learning, to the quadratic loss $(a - y)^2$, since additive and multiplicative constants in the definition of l only rescale it without changing the value of a that minimizes its expectation.

theoretic interpretation²: it is the *Kullback-Leibler divergence* [Kul59] (or *information gain* [Ren70]) from the actual conditional probability distribution P to the estimated conditional distribution \hat{P} , plus the entropy of P .

For a given x , the entropy of the true conditional distribution P is a constant, independent of the action a . Thus choosing the action a for each instance x that minimizes the expected log likelihood loss is equivalent to choosing the action a that gives the closest estimate \hat{P} to the true conditional distribution P over possible outcomes in Y as measured by the Kullback-Leibler divergence, given that instance x . It is well known that the Kullback-Leibler divergence is minimized when $\hat{P} = P$. This is Bayes optimal decision rule in regression.

In the regression version of our medical diagnosis situation, the definition of the log likelihood loss function depends on the interpretation of the components of the parameter vector a . If there are k possible diseases and the patient can have at most one of these, then we might have $k + 1$ possible mutually exclusive disease states y_1, \dots, y_{k+1} , where y_{k+1} means healthy. Hence $Y = \{y_1, \dots, y_{k+1}\}$. Then we might specify that an action a takes the form

$$a = (a_1, \dots, a_{k+1}),$$

where $a_i = \hat{P}(y_i; a)$, the estimated probability of disease state y_i . Here the components of the vector a must be positive and sum to one. In this case the log likelihood loss would be $l(y_i, a) = -\log a_i = -\log \hat{P}(y_i; a)$.

Often the constraints on the components of a are a nuisance, so other interpretations of a are used, e.g. that $a_i = \log \hat{P}(y_i; a) - \log \hat{P}(y_{k+1}; a)$ for each i , $1 \leq i \leq k + 1$. In this case the a_1, \dots, a_k are arbitrary real numbers and $a_{k+1} = 0$, and hence can be ignored. Since $\hat{P}(y_i; a) = e^{a_i} / \sum_{j=1}^{k+1} e^{a_j}$, the log likelihood loss is $l(y_i, a) = -a_i + \log \sum_{j=1}^{k+1} e^{a_j} = -a_i + \log(1 + \sum_{j=1}^k e^{a_j})$. This is known as the *logistic loss* [MN89, Bar89]. A third interpretation would be to allow the possibility that the patient may have more than one disease, and assume, for the purposes of estimation, that diseases occur independently. Then the disease state y might be defined as a binary vector of length k , where the i^{th} bit y_i is 1 if and only if the i^{th} disease is present. Hence $Y = \{0, 1\}^k$. Similarly, the vector a would be a vector of independent probabilities (a_1, \dots, a_k) , where a_i is the estimated probability of the patient having the i^{th} disease. In this case

$$\hat{P}(y; a) = \prod_{i=1}^k a_i^{y_i} (1 - a_i)^{(1-y_i)}$$

²The Kullback-Leibler divergence from P to \hat{P} , denoted $I(P||\hat{P})$, is defined as $\sum_{y \in Y} P(y) \log \frac{P(y)}{\hat{P}(y; a)}$ for countable Y . The entropy of P , denoted $H(P)$, is $-\sum_{y \in Y} P(y) \log P(y)$. Thus $I(P||\hat{P}) + H(P) = -\sum_{y \in Y} P(y) \log \hat{P}(y; a)$, which is the expectation of the (negative) log likelihood loss. Analogous results hold for densities when the relevant quantities are finite [Kul59].

prediction of the outcome y . Hence, a decision rule h maps from the instance space X into the outcome space Y , just as the target function does. In much of AI, and in PAC learning in particular, it is common to refer to h as a *hypothesis* in this case, and to \mathcal{H} as the *hypothesis space*.

This same setup, where the outcome y is a function of the instance x , can be applied to any function learning problem by letting X and Y be arbitrary sets. In the general function learning problem, the loss function $l(y, a)$ usually measures the distance between the prediction a and the actual value y in some metric. In the PAC model, l is the discrete metric: $l(y, a) = 0$ if $a = y$, else $l(y, a) = 1$. Thus the expected loss of the decision rule (or hypothesis) is just the probability that it predicts incorrectly, the usual PAC notion of the *error* of the hypothesis. In general, Y may be a set of strings, graphs, real vectors, etc., in which case other distance metrics or more general kinds of loss functions may be more appropriate.

1.1.3 Regression

The general problem of regression has a different character from that of classification learning, but can also be addressed in the decision theoretic learning framework. To illustrate this, as a third example consider a variant of the medical diagnosis situation in which the doctor provides an estimate of the probability that the patient has each of several diseases, rather than predicting that he has one specific disease or asserting that he is healthy. (Here we assume that the actual disease state includes at most one disease.) For example, the doctor may say “Given these test results x , I would say you have disease 1 with probability 55%, disease 2 with probability 5%, and no disease at all with probability 40%.” Here the doctor is actually trying to estimate the conditional distribution on disease states Y given the test results x . Her action a entails providing a vector of parameters that determine that estimated distribution, e.g. $(0.55, 0.05, 0.4)$. The decision space A is the set of all such parameter vectors.

Now let Y be an arbitrary discrete outcome space. Keeping the instance x fixed, for each parameter vector a in A and outcome y in Y let $\hat{P}(y; a)$ denote the probability of outcome y with respect to the distribution on Y defined by the parameter vector a . Thus when we take action a on instance x , we are asserting that, given the instance x , we estimate the conditional probability of outcome y to be $\hat{P}(y; a)$ for each outcome y in Y . Let $P(y)$ denote the actual conditional probability of outcome y , given the instance x , with respect to the unknown joint distribution on $X \times Y$. (The distributions P and \hat{P} can be replaced by densities when Y is continuous.) Let us define¹ the loss function l by setting $l(y, a) = -\log \hat{P}(y; a)$. This is called the (negative) *log likelihood* loss function. If we define loss in this way, then the expected loss resulting from action a has a natural information

¹We assume $\hat{P}(y; a) > 0$ for all y in Y .

1.1.1 Betting example

For our first example, consider the problem of learning to maximize profit (or minimize loss!) at the horse races. Here an instance x in X is a race, an action a in A consists of placing or not placing a certain bet, and an outcome y in Y is determined by the winner and the second and third place finishers. The loss $l(y, a)$ is the amount of money lost when bet a is placed and the outcome of the race is y . A negative loss is interpreted as gain. The joint distribution on $X \times Y$ represents the probability of various races and outcomes. This joint distribution is unknown to the learner; he only has random examples $(x_1, y_1), \dots, (x_m, y_m)$, each consisting of a race/outcome pair generated from this distribution. From these examples, the learner develops a deterministic betting strategy (decision rule). The best decision rule h is one that specifies a bet a for each race x that minimizes the expectation of the loss $l(y, a)$, when y is chosen randomly from the unknown conditional distribution on Y given x , which is determined by the underlying joint distribution on $X \times Y$. This (not necessarily unique) best decision rule minimizes the expected loss on a random example (x, y) . It is known as *Bayes optimal decision rule*. The learner tries to approximate Bayes optimal decision rule as best he can using decision rules from a given decision rule space \mathcal{H} (e.g. “simple” or “easy to compute” decision rules, or perhaps decision rules that can be represented by a particular kind of neural network).

1.1.2 Classification

As a second example, consider the problem of medical diagnosis. Here an instance x is a vector of measurements from medical tests conducted on the patient, an action a is a diagnosis of the patient’s disease state, and an outcome y may be defined as the actual disease state of the patient. Here $A = Y$, i.e. the possible diagnoses are the same as the possible disease states. To specify the loss function l , we may stipulate that there is zero loss for the correct diagnosis $a = y$, but for each pair (y, a) with diagnosis a differing from disease state y there is some positive real loss $l(y, a)$, depending on the severity of the consequences of that particular misdiagnosis. Here a decision rule is a diagnostic method, and Bayes optimal decision rule is the one that minimizes the expected loss from misdiagnosis when examples (x, y) of test results and associated disease states occur randomly according to some unknown “natural” joint distribution.

This medical diagnosis situation is a typical example of a *classification learning* problem in the field of pattern recognition (see e.g. [DH73]). The problem of learning a Boolean function from noise-free examples, as investigated in the PAC model, is a special case of classification learning. Here the outcome space Y is $\{0, 1\}$ and only the instance x in an example (x, y) is drawn at random. The outcome y is $f(x)$ for some unknown Boolean *target function* f , rather than being determined stochastically. As above, the decision space A is the same as the outcome space Y , and the action a can be interpreted as a

1.1 Overview of the proposed framework

To extend the PAC model, we propose a more general framework based on statistical decision theory (see e.g. Ferguson [Fer67], Kiefer [Kie87] or Berger [Ber85]). In this general framework we assume the learner receives randomly drawn training examples, each example consisting of an instance $x \in X$ and an outcome $y \in Y$, where X and Y are arbitrary sets called *instance* and *outcome spaces*, respectively. These examples are generated according to a joint distribution on $X \times Y$, unknown to the learner. This distribution comes from a (known) class \mathcal{P} of joint distributions on $X \times Y$, representing possible “states of nature.” After training, the learner will receive further random examples drawn from this same joint distribution. For each example (x, y) , the learner will be shown only the instance x . Then he will be asked to choose an action a from a set of possible actions A , called the *decision space*. Following this, the outcome y will be revealed to the learner. In the case that we examine here, the outcome y depends only on the instance x and not on the action a chosen by the learner. For each action a and outcome y , the learner will suffer a loss, which is measured by a fixed real-valued *loss function* l on $Y \times A$. We assume that the loss function is known to the learner. The learner tries to choose his actions so as to minimize his loss.

Here we look at the case in which, based on the training examples, the learner develops a deterministic strategy that specifies what he believes is the appropriate action a for each instance x in X . He then uses this strategy on all future examples. Thus we look at “batch” learning rather than “incremental” or “on-line” learning [Lit88]. The learner’s strategy, which is a function from the instance space X into the decision space A , will be called a *decision rule*. We assume that the decision rule is chosen from a fixed *decision rule space* \mathcal{H} of functions from X into A . For example, instances in X may be encoded as inputs to a neural network, and outputs of the network may be interpreted as actions in A . In this case the network represents a decision rule, and the decision rule space \mathcal{H} may be all functions represented by networks obtained by varying the parameters of a fixed underlying network. The goal of learning is to find a decision rule in \mathcal{H} that minimizes the expected loss, when examples are drawn at random from the unknown joint distribution on $X \times Y$.

This learning framework can be applied in a variety of situations. We now give several illustrations. For further discussion, we refer the reader to the excellent surveys of White [Whi90b], Barron [Bar89], Devroye [Dev88], and Vapnik [Vap89], to which we are greatly indebted. We also recommend the text by Kiefer [Kie87] for a general introduction to statistical inference and decision theory.

1 Introduction

The introduction of the Probably Approximately Correct (PAC) model [Val84] [Ang88] of learning from examples has done an admirable job of drawing together practitioners of machine learning with theoretically oriented computer scientists in the pursuit of a solid and useful mathematical foundation for applied machine learning work. These practitioners include both those in mainstream artificial intelligence and in neural net research. However, in attempting to address the issues that are relevant to this applied work in machine learning, a number of shortcomings of the model have cropped up repeatedly. Among these are the following:

1. The model is defined only for $\{0, 1\}$ -valued functions. Practitioners would like to learn functions on an instance space X that take values in an arbitrary set Y , e.g. multi-valued discrete functions, real-valued functions and vector-valued functions.
2. Some practitioners are wary of the assumption that the examples are generated from an underlying “target function”, and are not satisfied with the noise models that have been proposed to weaken this assumption (e.g. [AL88] [Slo88] [SV88]). They would like to see more general regression models investigated in which the y component in a training example $(x, y) \in X \times Y$ is randomly specified according to a conditional distribution on Y , given x . Here the general goal is to approximate this conditional distribution for each instance $x \in X$. In the computational learning theory literature, a model of this type is investigated in [KS90], with $Y = \{0, 1\}$, and in a more general case in [Yam90].
3. Many learning problems are unsupervised, i.e. the learner has access only to randomly drawn, unlabeled examples from an instance space X . Here learning can often be viewed as some form of approximation of the distribution that is generating these examples. This is usually called *density estimation* when the instance space X is continuous and no specific parametric form for the underlying distribution on X is assumed. It is often called *parameter estimation* when specific parametric probability models are used. One example of this in the computational learning theory literature is the recent investigation of Abe and Warmuth into the complexity of learning the parameters in a hidden Markov model [AW90].

Our purpose here is twofold. First, we propose an extension of the PAC model, based on the work of Vapnik and Chervonenkis [Vap89] and Pollard [Pol84, Pol90], that addresses these and other issues. Second, we use this extension to obtain distribution-independent upper bounds on the size of the training set needed for learning with various kinds of feedforward neural networks. [RM86] [PG89], a popular learning method that is not covered by the basic PAC model.

Decision Theoretic Generalizations of the PAC Model for Neural Net and Other Learning Applications

David Haussler

September, 1989

Revised: December, 1990

Revised again: January, 1992

haussler@saturn.ucsc.edu

Baskin Center for Computer Engineering and Information Sciences

University of California, Santa Cruz, CA 95064

Running head: Generalizations of the PAC Model

Abstract: We describe a generalization of the PAC learning model that is based on statistical decision theory. In this model the learner receives randomly drawn examples, each example consisting of an instance $x \in X$ and an outcome $y \in Y$, and tries to find a decision rule $h : X \rightarrow A$, where $h \in \mathcal{H}$, that specifies the appropriate action $a \in A$ to take for each instance x , in order to minimize the expectation of a loss $l(y, a)$. Here X , Y , and A are arbitrary sets, l is a real-valued function, and examples are generated according to an arbitrary joint distribution on $X \times Y$. Special cases include the problem of learning a function from X into Y , the problem of learning the conditional probability distribution on Y given X (regression), and the problem of learning a distribution on X (density estimation).

We give theorems on the uniform convergence of empirical loss estimates to true expected loss rates for certain decision rule spaces \mathcal{H} , and show how this implies learnability with bounded sample size, disregarding computational complexity. As an application, we give distribution-independent upper bounds on the sample size needed for learning with feedforward neural networks. Our theorems use a generalized notion of VC dimension that applies to classes of real-valued functions, adapted from Vapnik and Pollard's work, and a notion of *capacity* and *metric dimension* for classes of functions that map into a bounded metric space.