# Sample-Efficient Strategies for Learning in the Presence of Noise

NICOLÒ CESA-BIANCHI

*University of Milan, Milan, Italy*

ELI DICHTERMAN

*IBM Haifa Research Laboratory, Haifa, Israel*

PAUL FISCHER

*University of Dortmund, Dortmund, Germany*

ELI SHAMIR

*Hebrew University, Jerusalem, Israel*

AND

HANS ULRICH SIMON

*Ruhr-Universität, Bochum, Germany*

---

Authors' present addresses: N. Cesa-Bianchi, DSI, University of Milan, Via Comelico 39, I-2035, Milano, Italy, e-mail: cesabian@dsi.unimi.it; E. Dichterman, IBM Haifa Research Laboratory, MATAM, Haifa, 31905, Israel, e-mail: eli@haifa.vnet.ibm.com; P. Fischer, Lehrstuhl Informatik II, University of Dortmund, D-44221 Dortmund, Germany, e-mail: paul@ls2.informatik.uni-dortmund.de; E. Shamir, Department of Computer Science, Hebrew University, Jerusalem, Israel, e-mail: shamir@cs.huji.ac.il; H. Ulrich Simon, Lehrstuhl Mathematik und Informatik, Fakultät für Mathematik, Ruhr-Universität Bochum, D-44780 Bochum, Germany, e-mail: simon@lmi.ruhr-uni-bochum.de.

Abstract. In this paper, we prove various results about PAC learning in the presence of malicious noise. Our main interest is the sample size behavior of learning algorithms. We prove the first nontrivial sample complexity lower bound in this model by showing that order of $\epsilon/\Delta^2 + d/\Delta$ (up to logarithmic factors) examples are necessary for PAC learning any target class of $\{0, 1\}$-valued functions of VC dimension $d$, where $\epsilon$ is the desired accuracy and $\eta = \epsilon/(1 + \epsilon) - \Delta$ the malicious noise rate (it is well known that any nontrivial target class cannot be PAC learned with accuracy $\epsilon$ and malicious noise rate $\eta \geq \epsilon/(1 + \epsilon)$, this irrespective to sample complexity). We also show that this result cannot be significantly improved in general by presenting efficient learning algorithms for the class of all subsets of $d$ elements and the class of unions of at most $d$ intervals on the real line. This is especially interesting as we can also show that the popular minimum disagreement strategy needs samples of size $d\epsilon/\Delta^2$, hence is not optimal with respect to sample size. We then discuss the use of randomized hypotheses. For these the bound $\epsilon/(1 + \epsilon)$ on the noise rate is no longer true and is replaced by $2\epsilon/(1 + 2\epsilon)$. In fact, we present a generic algorithm using randomized hypotheses that can tolerate noise rates slightly larger than $\epsilon/(1 + \epsilon)$ while using samples of size $d/\epsilon$ as in the noise-free case. Again one observes a quadratic powerlaw (in this case $d\epsilon/\Delta^2$, $\Delta = 2\epsilon/(1 + 2\epsilon) - \eta$) as $\Delta$ goes to zero. We show upper and lower bounds of this order.

Categories and Subject Descriptors: I.2.6 [**Artificial Intelligence**]: Learning–*concept learning*

General Terms: Theory

Additional Key Words and Phrases: Learning with malicious noise, PAC learning

## 1. *Introduction*

Any realistic learning algorithm should be able to cope with errors in the training data. A model of learning in the presence of *malicious noise* was introduced by Valiant [1984] as an extension of his basic PAC framework for learning classes of $\{0, 1\}$-valued functions. In this *malicious PAC model*, each training example given to the learner is independently replaced, with fixed probability $\eta$, by an adversarially chosen one (which may or may not be consistent with the $\{0, 1\}$-valued target function). In their comprehensive investigation of malicious PAC learning, Kearns and Li [1993] show that a malicious noise rate $\eta \geq \epsilon/(1 + \epsilon)$ can make statistically indistinguishable two target functions that differ on a subset of the domain whose probability measure is at least $\epsilon$. This implies that, with this noise rate, no learner can generate hypotheses that are $\epsilon$-good in the PAC sense, irrespective of the sample size (number of training examples) and to the learner's computational power. In their work, Kearns and Li also analyze the performance of the minimum disagreement strategy in presence of malicious noise. They show that, for a sample of size[1] $d/\epsilon$ (where $d$ is the VC dimension of the target class) and for a noise rate bounded by any constant fraction of $\epsilon/(1 + \epsilon)$, the hypothesis in the target class having the smallest sample error is $\epsilon$-good in the PAC sense.

We begin this paper by studying the behaviour of sample complexity in the case of deterministic hypotheses and a high malicious noise rate, that is, a malicious noise rate arbitrarily close to the information-theoretic upper bound $\eta_{\text{det}} := \epsilon/(1 + \epsilon)$. We prove the first nontrivial lower bound in this model (Theorem 3.4) by showing that at least order of $\epsilon/\Delta^2 + d/\Delta$ examples are needed to PAC learn, with accuracy $\epsilon$ and tolerating a malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, every class of $\{0, 1\}$-valued functions of VC dimension $d$. Our proof combines, in an original way, techniques from Ehrenfeucht et al. [1989], Kearns and Li [1993] and Simon [1996a] and uses some new estimates of the tails of the

---

[1] All sample size orders in this section are given up to logarithmic factors.

binomial distribution that may be of independent interest. We then prove that this lower bound cannot be improved in general. Namely, we show (Theorem 3.10 and Corollary 3.15) that there is an algorithm RMD (for Randomized Minimum Disagreement) that, for each $d$, learns both the class $\mathscr{C}_d$ of all subsets of $d$ elements and the class $\mathscr{I}_d$ of unions of at most $d$ intervals on the real line using a noisy sample whose size is of the same order as the size of our lower bound. Algorithm RMD makes essential use of the fact that the learning domain is small. Hence, for the class $\mathscr{I}_d$ we first discretize the universe in a suitable way. Then Algorithm RMD uses a majority vote to decide the classification of those domain points which have a clear majority of one label, and tosses a fair coin to decide the classification of the remaining points.

We also show a lower bound of order $d\epsilon/\Delta^2$ for the sample size of the popular strategy of choosing any hypothesis that minimizes disagreements on the sample (Theorem 3.9). This bound holds for any class of VC dimension $d \geq 3$ and for every noise rate $\eta$ such that $\epsilon/(1 + \epsilon) - \eta = \Delta = o(\epsilon)$. This implies that, for high noise rate $\eta$ and for every target class of VC dimension $d$ large enough, there are distributions over the target domain where the minimum disagreement strategy is outperformed by algorithm RMD. To our knowledge, this is the first example of a natural PAC learning problem for which choosing any minimum disagreement hypothesis from a fixed hypothesis class is probably worse, in terms of sample complexity, than a different learning strategy.

In the second part of the paper, we consider the use of randomized hypotheses for learning with small sample sizes and high malicious noise rates. An easy modification of Kearns and Li's argument (Proposition 4.1) shows that no learner can output $\epsilon$-good randomized hypotheses with a noise rate larger or equal to $\eta_{\text{rand}} := 2\epsilon/(1 + 2\epsilon)$. Given the gap between this bound and the corresponding bound $\epsilon/(1 + \epsilon)$ for learners using deterministic hypotheses, we address the problem whether allowing randomized hypotheses helps in this setting. In fact, we present an algorithm (Theorem 4.2) that PAC learns any target class of VC dimension $d$ using randomized hypotheses and $d/\epsilon$ training examples while tolerating any noise rate bounded by a constant fraction of $(7/6)\epsilon/(1 + (7/6)\epsilon)$. The algorithm works by finding up to three functions in the target class that satisfy a certain independence condition defined on the sample. The value of the final hypothesis on a domain point is then computed by taking a majority vote over these functions (or by tossing a coin in case only two functions are found). Our investigation then moves on to consider the case of a noise rate close to the information-theoretic bound $2\epsilon/(1 + 2\epsilon)$ for randomized hypotheses. We show a strategy (Theorem 4.4) for learning the powerset of $d$ elements using $d\epsilon/\Delta^2$ training examples and tolerating a malicious noise rate of $\eta = 2\epsilon/(1 + 2\epsilon) - \Delta$, for every $\Delta > 0$. We also show (Theorem 4.5) that this sample size is optimal in the sense that *every* learner using randomized hypotheses needs at least $d\epsilon/\Delta^2$ training examples for learning any target class of VC dimension $d$.

## 2. *Definitions and Notation*

We recall the definitions of PAC learning and malicious PAC learning of a given *target class* $\mathscr{C}$, where $\mathscr{C}$ is a set of $\{0, 1\}$-valued functions $C$, here called *concepts*, defined on some common domain $X$. We call *instance* every $x \in X$ and *labeled instance* or *example* every pair $(x, y) \in X \times \{0, 1\}$. In Valiant's PAC learning

model [Valiant 1984], the learning algorithm (or learner) gets as input a *sample*, that is, a multiset $((x_1, C(x_1)), \ldots, (x_m, C(x_m)))$ of desired size $m < \infty$. Each instance $x_t$ in the sample given to the learner must be independently drawn from the same distribution $D$ on $X$ and labeled according to the same target function $C \in \mathscr{C}$. Both $C$ and $D$ are fixed in advance and unknown to the learner. In the malicious PAC model, the input sample is corrupted by an adversary using noise rate $\eta > 0$ according to the following protocol. First, a sample $((x_1, C(x_1)), \ldots, (x_m, C(x_m)))$ of the desired size is generated exactly as in the noise-free PAC model. Second, before showing the sample to the learner, each example $(x_t, C(x_t))$ is independently marked with fixed probability $\eta$. Finally, the adversary replaces each marked example $(x_t, C(x_t))$ in the sample by a pair $(\hat{x}_t, \hat{y}_t)$ arbitrarily chosen from $X \times \{0, 1\}$ and then feeds the corrupted sample to the learner. We call the collection of marked examples the *noisy part* of the sample and the collection of unmarked examples the *clean part* of the sample. Note that learning in this model is harder than with the definition of malicious PAC learning given by Kearns and Li [1993]. There, the examples were sequentially ordered in the sample and the adversary's choice for each marked example had to be based only on the (possibly marked) examples occurring *earlier* in the sample sequence.[2] We call *KL-adversary* this weaker type of adversary. We also consider a third type of malicious adversary, which we call "nonadaptive". Whereas the corruption strategy for our malicious adversary is a function from the set of samples to the set of corrupted samples, the corruption strategy for a nonadaptive adversary is a function from the set of examples to the set of corrupted examples. In other words, a nonadaptive adversary must decide, before seeing the sample, on a fixed rule for replacing each pair $(x, C(x))$, for $x \in X$, with a corrupted pair $(x', y')$. Note that learning with a KL-adversary is easier than learning with our malicious adversary, but harder than learning with a nonadaptive adversary.

To meet the PAC learning criterion, the learner, on the basis of a polynomially-sized sample (which is corrupted in case of malicious PAC learning), must output an hypothesis $H$ that with high probability is a close approximation of the target $C$. Formally, an algorithm $A$ is said to *PAC learn* a target class $\mathscr{C}$ using hypothesis class $\mathscr{H}$ if, for all distributions $D$ on $X$, for all targets $C \in \mathscr{C}$, and for all $1 \geq \epsilon, \delta > 0$, given as input a sample of size $m$, $A$ outputs an hypothesis $H \in \mathscr{H}$ such that its *error probability*, $D(H \neq C)$, is strictly smaller than $\epsilon$ with probability at least $1 - \delta$ with respect to the sample random draw, where $m = m(\epsilon, \delta)$ is some polynomial in $1/\epsilon$ and $\ln(1/\delta)$. We call $\epsilon$ the *accuracy* parameter and $\delta$ the *confidence* parameter. We use $H \neq C$ to denote $\{x: H(x) \neq C(x)\}$. A hypothesis $H$ is called $\epsilon$-*good* (with respect to a distribution $D$) if it satisfies the condition $D(H \neq C) < \epsilon$; otherwise, it is called $\epsilon$-*bad*. Similarly, an algorithm $A$ is said to *learn* a target class $\mathscr{C}$ using hypothesis class $\mathscr{H}$ in the malicious PAC model with noise rate $\eta$ if $A$ learns $\mathscr{C}$ in the PAC model when the input sample is corrupted by any adversary using noise rate $\eta$. Motivated by the fact (shown in Kearns and Li [1993] and mentioned in the introduction) that a noise rate $\eta \geq \epsilon/(1 + \epsilon)$ forbids PAC learning with accuracy $\epsilon$, we allow the sample size $m$ to depend polynomially also on $1/\Delta$, where $\Delta = \epsilon/(1 + \epsilon) - \eta$.

---

[2] All the results from Kearns and Li [1993] we mention here hold in our harder noise model as well.

We will occasionally use *randomized* learning algorithms that have a sequence of tosses of a fair coin as additional input source. In this case the definition of PAC learning given above is modified so that $D(C \neq H) < \epsilon$ must hold with probability at least $1 - \delta$ also with respect to the algorithm's randomization. Finally, we will also use *randomized hypotheses* or *coin rules*. A coin rule is any function $F: X \rightarrow [0, 1]$ where $F(x)$ is interpreted as the probability that the Boolean hypothesis defined by the coin rule takes value 1 on $x$. Coin rules are formally equivalent to *p-concepts*, whose learnability has been investigated by Kearns and Schapire [1994]. However, here we focus on a completely different problem, that is, the malicious PAC learning of Boolean functions using *p*-concepts as hypotheses. If a learner uses coin rules as hypotheses, then the PAC learning criterion $D(C \neq H) < \epsilon$, where $H$ is the learner's hypothesis, is replaced by $\mathbf{E}_{x \sim D}|F(x) - C(x)| < \epsilon$, where $F$ is the coin rule output by a learner and $\mathbf{E}_{x \sim D}$ denotes expectation with respect to the distribution $D$ on $X$. Note that $|F(x) - C(x)|$ is the probability of misclassifying $x$ using coin rule $F$. Thus, $\mathbf{E}_{x \sim D}|F(x) - C(x)|$, which we call the *error probability* of $F$, is the probability of misclassifying a randomly drawn instance using coin rule $F$. Furthermore, since Proposition 4.1 in Section 4 shows that every noise rate larger or equal to $2\epsilon/(1 + 2\epsilon)$ prevents PAC learning with accuracy $\epsilon$ using randomized hypotheses, we allow the sample size of algorithms outputting randomized hypotheses to depend polynomially on $1/\Delta$, where $\Delta = 2\epsilon/(1 + 2\epsilon) - \eta$.

In addition to the usual asymptotical notations, let $\tilde{O}(f)$ be the order obtained from $O(f)$ by dropping polylogarithmic factors.

## 3. *Malicious Noise and Deterministic Hypotheses*

This section presents three basic results concerning the sample size needed for PAC learning in the presence of malicious noise. Theorems 3.4 and 3.7 establish the general lower bound $\Omega(\epsilon/\Delta^2 + d/\Delta)$ that holds for any learning algorithm. In Subsection 3.3, we show that this bound cannot be significantly improved. Finally, Theorem 3.9 presents the stronger lower bound $\Omega(d\epsilon/\Delta^2)$ that holds for the minimum disagreement strategy.

We make use of the following definitions and facts from probability theory. Let $S_{N,p}$ be the random variable that counts the number of successes in $N$ independent trials, each trial with probability $p$ of success. A real number $s$ is called *median* of a random variable $S$ if $\Pr\{S \leq s\} \geq 1/2$ and $\Pr\{S \geq s\} \geq 1/2$.

FACT 3.1. [JOGDEO AND SAMUELS 1968].    *For all $0 \leq p \leq 1$ and all $N \geq 1$, the median of $S_{N,p}$ is either $\lfloor Np \rfloor$ or $\lceil Np \rceil$. Thus,*

$$Pr\left\{ S_{N,p} \leq \lceil Np \rceil \right\} \geq \frac{1}{2} \ and \ Pr\left\{ S_{N,p} \geq \lfloor Np \rfloor \right\} \geq \frac{1}{2}.$$

FACT 3.2.    *If $0 < p < 1$ and $q = 1 - p$, then for all $N \geq 37/(pq)$,*

$$Pr\left\{ S_{N,p} \geq \lfloor Np \rfloor + \lfloor \sqrt{Npq - 1} \rfloor \right\} > \frac{1}{19} \tag{1}$$

$$Pr\left\{ S_{N,p} \le \lceil Np \rceil - \lfloor \sqrt{Npq - 1} \rfloor \right\} > \frac{1}{19}. \tag{2}$$

PROOF. The proof is in the Appendix. □

There are, in probability theory, various results that can be used to prove propositions of the form of Fact 3.2. A general but powerful tool is the Berry–Esseen theorem [Chow and Teicher 1988, Theorem 1, page 304] on the rate of convergence of the central limit theorem. This gives a lower bound on the left-hand side of (1) and (2) that is worse by approximately a factor of 2 than the bound 1/19 proven in Fact 3.2. More specialized results on the tails of the Binomial distribution were proven by Bahadur [1960], Bahadur and Ranga-Rao [1960], and Littlewood [1969]. We derive (1) and (2) by direct manipulation of the bound in Bahadur [1960], which is in a form suitable for our purposes.

FACT 3.3. *For every $0 < \beta < \alpha \le 1$, for every random variable $S \in [0, N]$ with* $\mathbf{E}S = \alpha N$, *it holds that $Pr\{S \ge \beta N\} > (\alpha - \beta)/(1 - \beta)$.*

PROOF. It follows by setting $z = Pr\{S \ge \beta N\}$ and solving the following for $z$:

$$\alpha N = \mathbf{E}S = \mathbf{E}[S|S < \beta N](1 - z) + \mathbf{E}[S|S \ge \beta N]z < \beta N(1 - z) + Nz. \quad □$$

3.1. A GENERAL LOWER BOUND ON THE SAMPLE SIZE. Two $\{0, 1\}$-valued functions $C_0$ and $C_1$ are called *disjoint* if there exists no $x \in X$ such that $C_0(x) = C_1(x) = 1$. A target class $\mathscr{C}$ is called *trivial* if any two targets $C_0, C_1 \in \mathscr{C}$ are either identical or disjoint. Kearns and Li [1993] have shown that nontrivial target classes cannot be PAC learned with accuracy $\epsilon$ if the malicious noise rate $\eta$ is larger or equal than $\epsilon/(1 + \epsilon)$. The proof is based on the statistical indistinguishability of two targets $C_0$ and $C_1$ that differ on some domain point $x$ which has probability $\epsilon$, but coincide on all other points with nonzero probability. The malicious nonadaptive adversary will present $x$ with the false label with probability $\eta_{\text{det}} := \epsilon/(1 + \epsilon)$. Hence, $x$ appears with the true label with probability $(1 - \eta_{\text{det}})\epsilon$. As $(1 - \eta_{\text{det}})\epsilon = \eta_{\text{det}}$, there is no chance to distinguish between $C_0$ and $C_1$.

Our first lower bound is based on a similar reasoning: For $\eta < \eta_{\text{det}}$, the targets $C_0$ and $C_1$ can be distinguished, but as $\eta$ approaches $\eta_{\text{det}}$, the discrimination task becomes arbitrarily hard. This idea is made precise in the proof of the following result.

THEOREM 3.4. *For every nontrivial target class $\mathscr{C}$, for every $0 < \epsilon < 1, 0 < \delta \le 1/38$, and $0 < \Delta = o(\epsilon)$, the sample size needed for PAC learning $\mathscr{C}$ with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, is greater than*

$$\frac{9\eta(1 - \eta)}{37\Delta^2} = \Omega\left(\frac{\eta}{\Delta^2}\right).$$

PROOF. For each nontrivial target class $\mathscr{C}$, there exist two points $x_0, x_1 \in X$ and two targets $C_0, C_1$ such that $C_0(x_0) = C_1(x_0) = 1$, $C_0(x_1) = 0$, and $C_1(x_1) = 1$. Let the distribution $D$ be such that $D(x_0) = 1 - \epsilon$ and $D(x_1) = \epsilon$. We will use a malicious adversary that corrupts examples with probability $\eta$. Let

$A$ be a (possibly randomized) learning algorithm for $\mathscr{C}$ that uses sample size $m = m(\epsilon, \delta, \Delta)$. For the purpose of contradiction, assume that $A$ PAC learns $\mathscr{C}$ against the (nonadaptive) adversary that flips a fair coin to select target $C \in \{C_0, C_1\}$ and then returns a corrupted sample $S'$ of size $m$ whose examples in the noisy part are all equal to $(x_1, 1 - C(x_1))$.

Let $p_A(m)$ be $\Pr\{H \neq C\}$, where $H$ is the hypothesis generated by $A$ on input $S'$. Since $H(x_1) \neq C(x_1)$ implies that $H$ is not an $\epsilon$-good hypothesis, we have that $p_A(m) \leq \delta \leq 1/38$ must hold. For the above malicious adversary, the probability that an example shows $x_1$ with the wrong label is $\eta$. The probability to see $x_1$ with the true label is a bit higher, namely $(1 - \eta)\epsilon = \eta + \Delta + \epsilon\Delta$. Let $B$ be the Bayes strategy that outputs $C_1$ if the example $(x_1, 1)$ occurs more often in the sample than $(x_1, 0)$, and $C_0$, otherwise. It is easy to show that $B$ minimizes the probability to output the wrong hypothesis. Thus, $p_B(m) \leq p_A(m)$ for all $m$. We now show that $m \leq 9\eta(1 - \eta)/(37\Delta^2)$ implies $p_B(m) > 1/38$. For this purpose, define events $BAD_1(m)$ and $BAD_2(m)$ over runs of $B$ that use sample size $m$ as follows. $BAD_1(m)$ is the event that at least $\lceil (\eta + \Delta)m \rceil + 1$ examples are corrupted, $BAD_2(m)$ is the event that the true label of $x_1$ is shown at most $\lceil (\eta + \Delta)m \rceil$ times. Clearly, $BAD_1(m)$ implies that the false label of $x_1$ is shown at least $\lceil (\eta + \Delta)m \rceil + 1$ times. Thus, $BAD_1(m)$ and $BAD_2(m)$ together imply that the hypothesis returned by $B$ is wrong. Based on the following two claims, we will show that whenever $m$ is too small, $\Pr\{BAD_1(m) \wedge BAD_2(m)\} > 1/38$.

CLAIM 3.5.   *For all $m \geq 1$, $Pr\{BAD_2(m)|BAD_1(m)\} \geq 1/2$.*

PROOF (OF THE CLAIM).   Given $BAD_1(m)$, there are less than $(1 - \eta - \Delta)m$ uncorrupted examples. Each uncorrupted example shows the true label of $x_1$ with probability $\epsilon$. In the average, the true label is shown less than $(1 - \eta - \Delta)\epsilon m = (1 - \eta_{\text{det}})\epsilon m = \eta_{\text{det}}m = (\eta + \Delta)m$ times. The claim now follows from Fact 3.1.   □

CLAIM 3.6.   *If $37/(\eta(1 - \eta)) \leq m \leq 9\eta(1 - \eta)/(37\Delta^2)$, then $Pr\{BAD_1(m)\} > 1/19$.*

PROOF (OF THE CLAIM).   Let $S_{m,\eta}$ denote the number of corrupted examples. Fact 3.2 implies that for all $m \geq 37/(\eta(1 - \eta))$,

$$\Pr\left\{ S_{m,\eta} \geq \lfloor m\eta \rfloor + \lfloor \sqrt{m\eta(1 - \eta) - 1} \rfloor \right\} > \frac{1}{19}.$$

The claim follows if

$$\lfloor m\eta \rfloor + \lfloor \sqrt{m\eta(1 - \eta) - 1} \rfloor > \lceil \eta m + \Delta m \rceil + 1.$$

This condition is implied by

$$m\eta + \sqrt{m\eta(1 - \eta) - 1} \geq \eta m + \Delta m + 3,$$

which, in turn, is implied by

$$\frac{1}{2}\sqrt{m\eta(1 - \eta) - 1} \geq 3 \quad \text{and} \quad \frac{1}{2}\sqrt{m\eta(1 - \eta) - 1} \geq \Delta m.$$

The latter two conditions easily follow from the lower and the upper bound on $m$ specified in the statement of the claim. $\square$

From these two claims, we obtain that, for $37/(\eta(1 - \eta)) \leq m \leq 9\eta(1 - \eta)/37\Delta^2$, it holds that $p_B(m) > 1/38$. Note that $\Delta \leq \epsilon/K$ (for a sufficiently large constant $K$) implies that the specified range for $m$ contains at least one integer, that is, the implication is not vacuous. As the Bayes strategy $B$ is optimal, it cannot be worse than a strategy which ignores sample points, thus the error probability $p_B(m)$ does not increase with $m$. We may therefore drop the condition $m \geq 37/(\eta(1 - \eta))$. This completes the proof.

The proof of our next lower bound combines the technique from Ehrenfeucht et al. [1989] for showing the lower bound on the sample size in the noise-free PAC learning model with the argument of statistical indistinguishability. Here, the indistinguishability is used to force with probability 1/2 a mistake on a point $x$, for which $D(x) = \eta/(1 - \eta)$. To ensure that, with probability greater than $\delta$, the learner outputs an hypothesis with error at least $\epsilon$, we use $t$ other points that are observed very rarely in the sample. This entails that the learning algorithm must essentially perform a random guess on half of them.

THEOREM 3.7. *For any target class $\mathscr{C}$ with VC dimension $d \geq 3$, and for every $0 < \epsilon \leq 1/8$, $0 < \delta \leq 1/12$, and every $0 < \Delta < \epsilon/(1 + \epsilon)$, the sample size needed for PAC learning $\mathscr{C}$ with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, is greater than*

$$\frac{d - 2}{32\Delta(1 + \epsilon)} = \Omega\left(\frac{d}{\Delta}\right).$$

Note that for $\Delta = \epsilon/(1 + \epsilon)$, that is, $\eta = 0$, this reduces to the known lower bound on the sample size for noise-free PAC learning.

PROOF. Let $t = d - 2$ and let $X_0 = \{x_0, x_1, \ldots, x_t, x_{t+1}\}$ be a set of points shattered by $\mathscr{C}$. We may assume without loss of generality, that $\mathscr{C}$ is the powerset of $X_0$. We define distribution $D$ as follows:

$$D(x_0) = 1 - \frac{\eta}{1 - \eta} - 8\left(\epsilon - \frac{\eta}{1 - \eta}\right),$$

$$D(x_1) = \cdots = D(x_t) = \frac{8(\epsilon - (\eta/(1 - \eta)))}{t}$$

$$D(x_{t+1}) = \frac{\eta}{1 - \eta}.$$

(Note that $\epsilon \leq 1/8$ implies that $D(x_0) \geq 0$.) We will use a nonadaptive adversary that, with probability $\eta$, corrupts each example by replacing it with $x_{t+1}$ labeled incorrectly. Therefore, $x_{t+1}$ is shown incorrectly labeled with probability $\eta$ and correctly labeled with probability $(1 - \eta)D(x_{t+1}) = \eta$. Thus, true and false labels for $x_{t+1}$ are statistically indistinguishable. We will call $x_1, \ldots, x_t$ *rare points* in the sequel. Note that when $\eta$ approaches $\eta_{\text{det}}$ the probability to select a rare point approaches 0. Let $A$ be a (possibly randomized) learning algorithm for

$\mathscr{C}$ using sample size $m = m(\epsilon, \delta, \Delta)$. Consider the adversary that flips a fair coin to select target $C \in \mathscr{C}$ and then returns a corrupted sample (of size $m$) whose examples in the noisy part are all equal to $(x_{t+1}, 1 - C(x_{t+1}))$.

Let $e_A$ be the random variable denoting the error $\Pr\{H \neq C\}$ of $A$'s hypothesis $H$. Then, by pigeonhole principle,

$$\Pr\{e_A \geq \epsilon\} > \frac{1}{12} \tag{3}$$

implies the existence of a concept $C_0 \in \mathscr{C}$ such that the probability that $A$ does not output an $\epsilon$-good hypothesis for $C_0$ is greater than $1/12 \geq \delta$. Let us assume for the purpose of contradiction that $m \leq t/(32\Delta(1 + \epsilon))$. It then suffices to show that (3) holds.

Towards this end, we will define events $BAD_1$, $BAD_2$, and $BAD_3$, over runs of $A$ that use sample size $m$, whose conjunction has probability greater than $1/12$ and implies (3). $BAD_1$ is the event that at least $t/2$ rare points are not returned as examples by the adversary. In what follows, we call *unseen* the rare points that are not returned by the adversary. Given $BAD_1$, let $BAD_2$ be the event that the hypothesis $H$ classifies incorrectly at least $t/8$ points among the set of $t/2$ unseen points with lowest indices. Finally, let $BAD_3$ be the event that hypothesis $H$ classifies $x_{t+1}$ incorrectly. It is easy to see that $BAD_1 \wedge BAD_2 \wedge BAD_3$ implies (3), because the total probability of misclassification adds up to

$$\frac{t}{8} \cdot \frac{8(\epsilon - (\eta/(1 - \eta)))}{t} + \frac{\eta}{1 - \eta} = \epsilon.$$

We finally have to discuss the probabilities of the 3 events. Only noise-free examples potentially show one of the rare points. The probability that this happens is

$$8\left(\epsilon - \frac{\eta}{1 - \eta}\right)(1 - \eta) = 8(\epsilon(1 - \eta) - \eta) = 8\Delta(1 + \epsilon).$$

Since $m \leq t/(32\Delta(1 + \epsilon))$, the examples returned by the adversary contain at most $t/4$ rare point in the average. It follows from Markov inequality that the probability that these examples contain more than $t/2$ rare points is smaller than $1/2$. Thus, $\Pr\{BAD_1\} > 1/2$. For each unseen point, there is a chance of $1/2$ of misclassification. Thus, $\Pr\{BAD_2 \mid BAD_1\}$ is the probability that a fair coin flipped $t/2$ times shows heads at least $t/8$ times. Fact 3.3 applied with $\alpha = 1/2$, and $\beta = 1/4$, implies that this probability is greater than $1/3$. Note that events $BAD_1$ and $BAD_2$ do not break the symmetry between the two possible labels for $x_{t+1}$ (although they change the probability that sample point $x_{t+1}$ is drawn at all). As the Boolean labels of $x_{t+1}$ are statistically indistinguishable, we get $\Pr\{BAD_3 \mid BAD_1, BAD_2\} = \Pr\{BAD_3\} = 1/2$. Thus,

$$\Pr\{BAD_1 \wedge BAD_2 \wedge BAD_3\} = \Pr\{BAD_1\} \cdot \Pr\{BAD_2 | BAD_1\} \cdot \Pr\{BAD_3\}$$

$$> \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{12},$$

which completes the proof.  □

Combining Theorems 3.4 and 3.7, we have:

COROLLARY 3.8.  *For any nontrivial target class $\mathscr{C}$ with, and for every $0 < \epsilon \leq 1/8$, $0 < \eta \leq 1/38$, and for every $0 < \Delta < \epsilon/(1 + \epsilon)$, the sample size needed for PAC learning $\mathscr{C}$ with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, is greater than*

$$\Omega\left(\frac{d}{\Delta} + \frac{\eta}{\Delta^2}\right) = \Omega\left(\frac{d}{\Delta} + \frac{\epsilon}{\Delta^2}\right).$$

Even though the lower bound $\Omega(\eta/\Delta^2 + d/\Delta)$ holds for the nonadaptive adversary, in Subsection 3.3 we show a matching upper bound (ignoring logarithmic factors) that holds for the (stronger) malicious adversary.

The lower bound proven in Theorem 3.4 has been recently improved to

$$\Omega\left(\frac{\eta}{\Delta^2} \ln \frac{1}{\delta}\right)$$

by Gentile and Helmbold [1998], who introduced an elegant information-theoretic approach avoiding the analysis of the Bayes risk associated with the learning problem. It is not clear whether their approach can be applied to obtain also the other lower bound term, $\Omega(d/\Delta)$, that we prove in Theorem 3.7.

The next result is a lower bound on the sample size needed by the minimum disagreement strategy (called MDS henceafter.)

3.2. A LOWER BOUND ON THE SAMPLE SIZE FOR MDS.  Given any sample $S'$ of corrupted training examples, MDS outputs a hypothesis $H \in \mathscr{C}$ with the fewest disagreements on $S'$.

THEOREM 3.9.  *For any target class $\mathscr{C}$ with VC dimension $d \geq 3$, every $0 < \epsilon \leq 1/38$, $0 < \delta \leq 1/74$, and every $0 < \Delta = o(\epsilon)$, the sample size needed by the Minimum Disagreement strategy for learning $\mathscr{C}$ with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, is greater than*

$$\frac{4(1 - \eta)(1 - \epsilon)\lceil (d - 1)/38 \rceil \epsilon}{37(1 + \epsilon)^2 \Delta^2} = \Omega\left(\frac{d\epsilon}{\Delta^2}\right).$$

PROOF.  The proof uses $d$ shattered points, where $d - 1$ of them (the rare points) have a relatively small probability. The total probability of all rare points is $c\epsilon$ for some constant $c$. Let $\mu$ be the mean and $\sigma$ the standard deviation for the number of true labels of a rare point within the (corrupted) training sample. If the rare points were shown $\mu$ times, the malicious adversary would have no chance to fool MDS. However, if the sample size $m$ is too small, the standard deviation $\sigma$ for the number of true labels of a rare point $x$ gets very big. Hence, with constant probability, the number of true labels of a rare point is smaller than (roughly) $\mu - \sigma$. If this happens, we call $x$ a hidden point. It follows that there is also a constant probability that a constant fraction of the rare points are hidden. This gives the malicious adversary the chance to present more false than

true labels for each hidden point. We now make these ideas precise. Our proof needs the following technical assumption:

$$m \geq \frac{37\lceil (d-1)/38 \rceil}{\epsilon(1-\epsilon)(1-\eta)}. \tag{4}$$

This condition can be forced by invoking the general lower bound $\Omega(d/\Delta)$ from Theorem 3.7 for $\Delta \leq \epsilon/K$ and a sufficiently large constant $K$.

For the purpose of contradiction, we now assume that

$$m \leq \frac{4(1-\eta)(1-\epsilon)\lceil (d-1)/38 \rceil \epsilon}{37(1+\epsilon)^2\Delta^2}. \tag{5}$$

Let $BAD_1$ be the event that at least $\lfloor \eta m \rfloor$ examples are corrupted by the malicious adversary. According to Fact 3.1, $BAD_1$ has probability at least 1/2. Let $t = d - 1$ and let $X_0 = \{x_0, \ldots, x_t\}$ be a set of points shattered by $\mathscr{C}$. Distribution $D$ is defined by

$$D(x_0) = 1 - t\left\lceil \frac{t}{38} \right\rceil^{-1}\epsilon, \quad D(x_1) = \cdots = D(x_t) = \left\lceil \frac{t}{38} \right\rceil^{-1}\epsilon.$$

Points $x_1, \ldots, x_t$ are called *rare*. Consider a fixed rare point $x_i$. Each example shows $x_i$ with its true label with probability

$$p = \left\lceil \frac{t}{38} \right\rceil^{-1}\epsilon(1-\eta) = \left\lceil \frac{t}{38} \right\rceil^{-1}(\eta + \Delta(1+\epsilon)).$$

Let $T_i$ denote the number of examples that present the true label of $x_i$. Call $x_i$ *hidden* if

$$T_i \leq \lceil pm \rceil - \lfloor \sqrt{mp(1-p) - 1} \rfloor.$$

Inequality (4) implies that $m \geq 37/(p(1-p))$. Thus, according to Fact 3.2, $x_i$ is hidden with probability greater than 1/19. Using the fact that

$$\Pr\{x_i \text{ is hidden}\} = \Pr\{x_i \text{ is hidden } |BAD_1\}\Pr\{BAD_1\}$$
$$+ \Pr\{x_i \text{ is hidden } |\neg BAD_1\}(1 - \Pr\{BAD_1\})$$

and

$$\Pr\{x_i \text{ is hidden } |BAD_1\} \geq \Pr\{x_i \text{ is hidden } |\neg BAD_1\},$$

it follows that

$$\Pr\{x_i \text{ is hidden } |BAD_1\} \geq \Pr\{x_i \text{ is hidden}\} > \frac{1}{19}.$$

Given $BAD_1$, let $T$ be the (conditional) random variable which counts the number of hidden points. The expectation of $T$ is greater than $t/19$. According to Fact 3.3 (with $\alpha = 1/19$ and $\beta = 1/38$), the probability that at least $t/38$ rare

points are hidden is greater than 1/37. Thus with probability greater than $\delta = 1/74$, there are (at least) $\lfloor \eta m \rfloor$ corrupted examples and (at least) $\lceil t/38 \rceil$ hidden points. This is assumed in the sequel.

The total probability of $\lceil t/38 \rceil$ hidden points (measured by $D$) is exactly $\lceil t/38 \rceil$ $\lceil t/38 \rceil^{-1} \epsilon = \epsilon$. It suffices therefore to show that there are enough corrupted examples to present each of the $\lceil t/38 \rceil$ hidden points with more false than true labels. The total number of true labels for $\lceil t/38 \rceil$ hidden points can be bounded from above as follows:

$$\left\lceil \frac{t}{38} \right\rceil \cdot \left( \lceil pm \rceil - \lfloor \sqrt{mp(1-p)} - 1 \rfloor \right) \leq \eta m + \Delta(1+\epsilon)m + 2\left\lceil \frac{t}{38} \right\rceil$$
$$- \left\lceil \frac{t}{38} \right\rceil \sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1}.$$

The number of false labels that the adversary can use is greater than $\eta m - 1$ and should exceed the number of true labels by at least $\lceil t/38 \rceil$. The adversary can therefore force an $\epsilon$-bad hypothesis of MDS if

$$\eta m - 1 \geq \eta m + \Delta(1+\epsilon)m + 3\left\lceil \frac{t}{38} \right\rceil - \left\lceil \frac{t}{38} \right\rceil \sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1}$$

or equivalently if

$$\left\lceil \frac{t}{38} \right\rceil \sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1} \geq \Delta(1+\epsilon)m + 3\left\lceil \frac{t}{38} \right\rceil + 1. \qquad (6)$$

We will develop a sufficient condition which is easier to handle. The right-hand side of (6) contains the three terms $z_1 = 3\lceil t/38 \rceil$, $z_2 = 1$, $z_3 = \Delta(1 + \epsilon)m$. Splitting the left-hand side $Z$ of (6) in three parts, we obtain the sufficient condition $Z/2 \geq z_1$, $Z/6 \geq z_2$, $Z/3 \geq z_3$, which reads (after some algebraic simplifications) in expanded form as follows:

$$\sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1} \geq 6$$

$$\left\lceil \frac{t}{38} \right\rceil \sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1} \geq 6$$

$$\left\lceil \frac{t}{38} \right\rceil \sqrt{\frac{m\epsilon(1-\eta)(1-\epsilon)}{\lceil t/38 \rceil} - 1} \geq 3\Delta(1+\epsilon)m$$

An easy computation shows that these three conditions are implied by (4) and (5). This completes the proof. $\square$

It is an open question whether a similar lower bound can be proven for the KL-adversary, or even the weaker nonadaptive adversary.

3.3. LEARNING THE POWERSET AND $k$-INTERVALS WITH DETERMINISTIC HYPOTHESES. Let $\mathscr{C}_k$ be the class of all subsets over $k$ points and let $\mathscr{I}_k$ be the class of unions of at most $k$ intervals on the unit interval $[0, 1]$. The Vapnik–Chervonenkis dimension of $\mathscr{I}_k$ is $2k$. In this subsection, we show that the lower bound proven in Subsection 3.1 cannot be improved in general. That is, we show that for each $d \geq 1$, the class $\mathscr{C}_d$ of all subsets over $d$ points and the class $\mathscr{I}_k$, $k = d/2$, can be PAC learned with accuracy $\epsilon > 0$ and malicious noise rate $\eta < \epsilon/(1 + \epsilon)$ using a sample of size $\bar{O}(\eta/\Delta^2 + d/\Delta)$, where $\Delta = \epsilon/(1 + \epsilon) - \eta$.

According to Subsection 3.2, strategy MDS is not optimal for PAC learning powersets in the presence of malicious noise. Instead, we use a modification of MDS, called RMD-POW, which uses a majority vote on the sample (like MDS) to decide the labels of some of the domain points and tosses a fair coin to decide the labels of the remaining ones. The heart of RMD-POW is the subroutine that splits the domain into two appropriate subgroups.

It turns out that the result for $k$-intervals can be derived as a corollary to the result for powersets if the latter result is proven in a slightly generalized form. Instead of considering $\mathscr{C}_d$ and an arbitrary distribution $D$ on $\{1, \ldots, d\}$, we consider $\mathscr{C}_N$, $N \geq d$, and a restricted subclass of distributions, where the restriction will become vacuous in the special case that $N = d$. Thus, although the result is distribution-specific in general, it implies a distribution-independent result for $\mathscr{C}_d$.

The restriction which proves useful later can be formalized as follows. Given parameters $d$, $\Delta$ and a distribution $D$ on $X_N = \{1, \ldots, N\}$, we call point $i \in X_N$ *light* if $D(i) < \Delta/(3d)$, and *heavy* otherwise. We say that $D$ is a *legal distribution* on $X_N$ if it induces at most $d$ light points. The main use of this definition is that even a hypothesis that is incorrect on all light points is only charged by less than $\Delta/3$ (plus its total error on heavy points, of course). Note furthermore that the existence of a legal distribution on $X_N$ implies that $N \leq d + 3d/\Delta$, and every distribution is legal if $N = d$.

The rest of this subsection is structured as follows: We first prove a general result for powerset $\mathscr{C}_N$, $N \geq d$, and legal distributions with the obvious corollary for powerset $\mathscr{C}_d$ and arbitrary distributions. Afterwards, we present a general framework of partitioning the domain $X$ of a concept class into so-called *bins* in such a way that labels can be assigned binwise without much damage. A bin $B \subseteq X$ is called *homogeneous* with respect to target concept $C$ if $C$ assigns the same label to all points in $B$. Otherwise, $B$ is called *heterogeneous*. If there are at most $d$ heterogeneous bins of total probability at most $\Delta/3$, we may apply the subroutine for $\mathscr{C}_N$ (working well for legal domain distributions), where heterogeneous bins play the role of light points and homogeneous bins the role of heavy (or another kind of unproblematic) points. From these general considerations, we derive the learnability result for $k$-intervals by showing that (for this specific concept class) an appropriate partition of the domain can be efficiently computed from a given sample with high probability of success.

THEOREM 3.10. *There exists a randomized algorithm, RMD-POW, which achieves the following. For every $N \geq d \geq 1$ and every $1 \geq \epsilon$, $\delta$, $\Delta > 0$, RMD-POW PAC learns the class $\mathscr{C}_N$ under legal distributions with accuracy $\epsilon$, confidence $\delta$, tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, and using a sample of size $\bar{O}(\epsilon/\Delta^2 + d/\Delta)$.*

**Input:** Parameters $\alpha, L, n$. Domain size $d$, accuracy $\varepsilon$, confidence $\delta$.
—Get sample $(i_1, y_1), \ldots, (i_m, y_m)$ of size $m = \widetilde{O}(\varepsilon/\Delta^2 + d/\Delta)$.
—For each point $i \in X_N = \{1, \ldots, N\}$.
  (1)  If $i$ is in strong majority or belongs to a sparse band, then let $H(i)$ be the most frequent label with which $i$ appears in the sample;
  (2)  else, let $H(i)$ be a random label.
—Output the hypothesis $H$.

FIG. 1.   Pseudo-code for the randomized algorithm RMD-POW (see Theorem 3.10).

PROOF.    Given parameters $\Delta$, $d$ and $N \leq d + 3d/\Delta$, a learning algorithm has to deliver a good approximation to an unknown subset $C$ of domain $X_N = \{1, \ldots, N\}$ with high probability of success. Let $D$ denote the unknown domain distribution, and $p_i = D(i)$ for $i = 1, \ldots, N$. It is assumed that $D$ induces at most $d$ light points, that is;

$$I_{\text{light}} = \left\{ i \in X_N : p_i < \frac{\Delta}{3d} \right\}$$

contains at most $d$ elements. Clearly, the total error (of any hypothesis) on light points is bounded by $D(I_{\text{light}}) < \Delta/3$. It is therefore sufficient to deliver, with probability at least $1 - \delta$ of success, a hypothesis whose total error on

$$I_{\text{heavy}} = \left\{ i \in X_N : p_i \geq \frac{\Delta}{3d} \right\}$$

is bounded by $\epsilon - \Delta/3$. In the sequel, we describe the algorithm RMD-POW and show that it achieves this goal. RMD-POW is parameterized by three parameters $\alpha$, $L$, $n$. We denote the logarithm to base $\beta$ by $\log_\beta$ and assume that $\alpha$, $L$, $n$ are initialized as follows:

$$\alpha = \sqrt[3]{5/3} - 1 \tag{7}$$

$$L = \left\lceil \log_{1+\alpha} \frac{6d(1 + \alpha)^2 \epsilon}{\Delta} \right\rceil \tag{8}$$

$$n = \left\lceil 50 \ln\left(\frac{4L}{\delta}\right) \right\rceil \tag{9}$$

An informal description of RMD-POW is given below. Its pseudo-code may be found in Figure 1. Algorithm RMD-POW can deliver any subset of $X_N$ as hypothesis, and can therefore decide the label of each point in $X_N$ independently. This is done as follows. Based on the sample, the domain is divided into two main groups. The label of each point $i$ in the first group is decided by taking a majority vote on the occurrences of $(i, 0)$ and $(i, 1)$ in the sample. The labels of the points in the second group are instead chosen in a random way.

To bound the total error of the hypothesis chosen by the algorithm, we divide each of the two above groups into subgroups, and then separately bound the contributions of each subgroup to the total error. Within each such subgroup, we approximately bound the total probability of the domain points for which the

algorithm chooses the wrong label by the total frequency of corrupted examples of points in the subgroup. Since, for a large enough sample, the sample frequency of corrupted examples is very close to the actual noise rate $\eta$, and since the noise rate is bounded away from the desired accuracy $\epsilon$, we can show that the total probability of the points labeled incorrectly, summed over all subgroups, is at most the desired accuracy $\epsilon$.

Given a sample $(i_1, y_1), \ldots, (i_m, y_m)$ drawn from the set $\{1, \ldots, N\} \times \{0, 1\}$, let $\nu_{0,i}$ and $\nu_{1,i}$ be the frequencies with which each point $i \in \{1, \ldots, N\}$ appears in the sample with label respectively 0 and 1. For each $i$, we define $\ell_i = \min\{\nu_{0,i}, \nu_{1,i}\}$ and $u_i = \max\{\nu_{0,i}, \nu_{1,i}\}$. We say that a domain point $i$ is in *strong majority* (with respect to the sample and to parameter $\alpha$ defined in (7)) if $u_i > (1 + \alpha)\ell_i$, and is in *weak majority* otherwise. We divide some of the points into $L$ *bands*, where $L$ is the parameter defined in (8). A point $i$ is in band $k$, for $k = 1, \ldots, L$, if $i$ is in weak majority and $(1 + \alpha)^{-k} \epsilon < \ell_i \leq (1 + \alpha)^{1-k} \epsilon$. We further divide the bands into *sparse bands*, containing less than $n$ elements, and *dense bands*, containing at least $n$ elements, where $n$ is the parameter defined in (9). Let $I_{\mathrm{maj}}$, $I_{\mathrm{sparse}}$, and $I_{\mathrm{dense}}$ be the sets of all domain points respectively in strong majority, sparse bands and dense bands. For fixed choice of input parameters, we denote RMD-POW's hypothesis by $H$.

For simplicity, for each point $i$ we will write $t_i$ and $f_i$ to denote, respectively, $\nu_{C(i),\, i}$ and $\nu_{1-C(i),\, i}$. That is, $t_i$ and $f_i$ are the sample frequencies of, respectively, clean and corrupted examples associated with each point $i$. We define

$$f_{\mathrm{maj}} = \sum_{i \in I_{\mathrm{maj}}} f_i , \; f_{\mathrm{sparse}} = \sum_{i \in I_{\mathrm{sparse}}} f_i , \; f_{\mathrm{dense}} = \sum_{i \in I_{\mathrm{dense}}} f_i .$$

First, we upper bound in probability the sum $f_{\mathrm{maj}} + f_{\mathrm{sparse}} + f_{\mathrm{dense}}$. Let $\hat{\eta}$ be the frequency of corrupted examples in the sample. By using (25) from Lemma A.1 with $p = \eta$ and $\lambda = \Delta/(3\eta)$, we find that

$$f_{\mathrm{maj}} + f_{\mathrm{sparse}} + f_{\mathrm{dense}} \leq \hat{\eta} \leq \left( 1 + \frac{\Delta}{3\eta} \right) \eta = \eta + \frac{\Delta}{3} \qquad (10)$$

holds with probability at least $1 - \delta/4$ whenever the sample size is at least

$$\left( \frac{27\eta}{\Delta^2} \right) \ln \left( \frac{4}{\delta} \right) = \tilde{O} \left( \frac{\eta}{\Delta^2} \right) .$$

Second, we lower bound in probability the sample frequency $t_i$ of uncorrupted examples for each $i \in I_{\mathrm{heavy}}$. Note that the probability that a point $i$ appears uncorrupted in the sample is at least $(1 - \eta)p_i$. Also, $|I_{\mathrm{heavy}}| \leq N$, as there are at most $N$ points. By using (23) from Lemma A.1 with $p = \Delta/(3d)$ and $\lambda = \alpha/(1 + \alpha)$, we find that

$$t_i \geq \frac{1 - \eta}{1 + \alpha} p_i = \left( 1 - \frac{\alpha}{1 + \alpha} \right)(1 - \eta)p_i \qquad \text{for all } i \in I_{\mathrm{heavy}} \qquad (11)$$

holds with probability at least $1 - \delta/4$ whenever the sample size is at least

$$\frac{6(1 + \alpha)^2 d}{(1 - \eta)\alpha^2 \Delta} \ln \frac{4N}{\delta} = \tilde{O} \left( \frac{d}{\Delta} \right) .$$

Thus, there is a sample size with order $\bar{O}(\eta/\Delta^2 + d/\Delta)$ of magnitude, such that (10) and (11) simultaneously hold with probability at least $1 - \delta/2$. At this point recall that $N$ linearly depends on $d$.

Let $I_{\text{wrong}} = \{i \in I_{\text{heavy}}: C(i) \neq H(i)\}$. Remember that $D(I_{\text{light}}) \leq \Delta/3$ and that it suffices to bound in probability $D(I_{\text{wrong}}) + \Delta/3$ from above by $\epsilon$. Claim 3.14 shows that, if (10) and (11) hold, then all heavy points are in the set $I_{\text{maj}} \cup I_{\text{sparse}} \cup I_{\text{dense}}$. Thus,

$$D(I_{\text{wrong}}) \leq D(I_{\text{wrong}} \cap I_{\text{maj}}) + D(I_{\text{wrong}} \cap I_{\text{sparse}}) + D(I_{\text{wrong}} \cap I_{\text{dense}}). \quad (12)$$

Now, Claims 3.11–3.13 show how the three terms in the right-hand side of (12) can be simultaneously bounded. In the rest of this subsection, we prove Claims 3.11–3.14. We start by bounding the error made by $H$ on heavy points $i \in I_{\text{maj}}$.

CLAIM 3.11. STRONG MAJORITY.    *If (11) holds, then*

$$D(I_{\text{wrong}} \cap I_{\text{maj}}) \leq \frac{f_{\text{maj}}}{1 - \eta}.$$

PROOF.    Recall that, for each $i \in I_{\text{maj}}$, $H(i) \neq C(i)$ if and only if $t_i = \ell_i$. Hence, if (11) holds, we find that for every $i \in I_{\text{wrong}} \cap I_{\text{maj}}$, $(1 - \eta)p_i \leq (1 + \alpha)t_i = (1 + \alpha)\ell_i \leq (1 + \alpha)/(1 + \alpha)u_i = f_i$. As $\sum_{i \in I_{\text{wrong}} \cap I_{\text{maj}}} f_i \leq f_{\text{maj}}$, the proof is concluded.    □

We now bound the error occurring in the sparse bands by proving the following:

CLAIM 3.12. SPARSE BANDS.    *Let the sample size be at least*

$$\frac{6(1 + \alpha)\epsilon L^2 n^2}{\Delta^2} \ln \frac{4d}{\delta} = \bar{O}\left(\frac{\epsilon}{\Delta^2}\right).$$

*Then (10) and (11) together imply that $D(I_{\text{wrong}} \cap I_{\text{sparse}}) \leq f_{\text{sparse}}/(1 - \eta) + \Delta/(3 \cdot (1 - \eta))$ holds with probability at least $1 - \delta/4$ with respect to the sample random draw.*

PROOF.    Recall that there are $L$ bands and each sparse band contains at most $n$ elements. We first prove that

$$t_i \geq (1 - \eta)p_i - \frac{\Delta}{3Ln} \qquad \text{for all } i \in I_{\text{wrong}} \cap I_{\text{sparse}} \quad (13)$$

holds in probability. To show this, we use (23) from Lemma A.1 to write the following

$$\Pr\{S_m \leq (p - \lambda)m\} = \Pr\left\{S_m \leq \left(1 - \frac{\lambda}{p}\right)mp\right\}$$

$$\leq \exp\left(-\frac{\lambda^2 m}{2p}\right) \leq \exp\left(-\frac{\lambda^2 m}{2p'}\right), \quad (14)$$

where the last inequality holds for all $p' \geq p$ by monotonicity. Now assume (10) and (11) both hold and choose $i$ such that $p_i > (1 + \alpha)\epsilon/(1 - \eta)$. Then $i \in I_{\text{heavy}}$ and $t_i \geq \epsilon$. As $\hat{\eta} < \epsilon$ by (10), $i \notin I_{\text{wrong}}$. Hence, (10) and (11) imply that $p_i \leq (1 + \alpha)\epsilon/(1 - \eta)$ holds for all $i \in I_{\text{wrong}}$. We then apply (14) to each $i \in I_{\text{wrong}} \cap I_{\text{sparse}}$. Setting $p = (1 - \eta)p_i$, $p' = (1 + \alpha)\epsilon \geq p$, and $\lambda = \Delta/(3Ln)$ we find that (13) holds with probability at least $1 - \delta/4$ whenever the sample size is at least

$$\frac{18(1 + \alpha)\epsilon L^2 n^2}{\Delta^2} \ln \frac{4d}{\delta} = \tilde{O}\left(\frac{\epsilon}{\Delta^2}\right).$$

Finally, from (13) we get that

$$D(I_{\text{wrong}} \cap I_{\text{sparse}}) \leq \sum_{I_{\text{wrong}} \cap I_{\text{sparse}}} \left(\frac{t_i}{1 - \eta} + \frac{\Delta}{(1 - \eta)3Ln}\right)$$

$$\leq \sum_{I_{\text{wrong}} \cap I_{\text{sparse}}} \frac{f_i}{1 - \eta} + \frac{\Delta}{3(1 - \eta)} = \frac{f_{\text{sparse}}}{1 - \eta} + \frac{\Delta}{3(1 - \eta)}.$$

This concludes the proof. □

We move on to bounding the error made on points in dense bands.

CLAIM 3.13. DENSE BANDS. *If* (11) *holds, then*

$$D(I_{\text{wrong}} \cap I_{\text{dense}}) \leq \frac{f_{\text{dense}}}{1 - \eta}$$

*holds with probability at least* $1 - \delta/4$ *with respect to the algorithm randomization.*

PROOF. For each $k = 1, \ldots, L$, let $I_k$ be the set of all heavy points in the $k$th band. Furthermore, let $t_{\max}^k = \max\{t_i : i \in I_k \cap I_{\text{wrong}}\}$ and $f_{\min}^k = \min\{f_i : i \in I_k \cap I_{\text{wrong}}\}$. Since all points in $I_k$ are in weak majority and by definition of bands, we have that $t_{\max}^k \leq (1 + \alpha)^2 f_{\min}^k$ holds for each $k = 1, \ldots, L$. Furthermore, using (11), $p_j \leq (1 + \alpha)t_j/(1 - \eta)$, for each $j \in I_k$. As for each dense band $|I_k| \geq n \geq 50\ln(4L/\delta)$, using (24) from Lemma A.1 we can guarantee that $|I_k \cap I_{\text{wrong}}| \leq (3/5)|I_k|$ holds simultaneously for all bands $k = 1, \ldots, L$ with probability at least $1 - \delta/4$. Combining everything we get

$$\sum_{I_k \cap I_{\text{wrong}}} p_i \leq \frac{1 + \alpha}{1 - \eta} \sum_{I_k \cap I_{\text{wrong}}} t_i$$

$$\leq \frac{1 + \alpha}{1 - \eta} \cdot \frac{3}{5}|I_k|t_{\max}^k$$

$$\leq \frac{(1 + \alpha)^3}{1 - \eta} \cdot \frac{3}{5}|I_k|f_{\min}^k$$

$$\leq \frac{(1 + \alpha)^3}{1 - \eta} \cdot \frac{3}{5}\sum_{I_k} f_i.$$

By choosing $\alpha = (5/3)^{1/3} - 1$ so that $(3/5)(1 + \alpha)^3 = 1$, we get

$$D(I_{\text{dense}} \cap I_{\text{wrong}}) = \sum_k \sum_{I_k \cap I_{\text{wrong}}} p_i \leq \sum_k \sum_{I_k} \frac{f_i}{1 - \eta} = \frac{f_{\text{dense}}}{1 - \eta}$$

concluding the proof.  □

CLAIM 3.14.   *If* (10) *and* (11) *hold, then* $I_{heavy} \subseteq I_{maj} \cup I_{sparse} \cup I_{dense}$.

PROOF.   We have to show that each heavy point $i$ in weak majority belongs to a band, or equivalently,

$$(1 + \alpha)^{-L}\epsilon < \ell_i \leq \epsilon \quad \text{for all } i \in I_{\text{heavy}}\backslash I_{\text{maj}}.$$

If (10) holds, then $\ell_i \leq f_i \leq \eta + (\Delta/3) < \epsilon$ for each point $i$. Also, if (11) holds, then, for each $i \in I_{\text{heavy}}\backslash I_{\text{maj}}$,

$$\ell_i \geq \frac{u_i}{1 + \alpha} \geq \frac{t_i}{1 + \alpha} \geq \frac{1 - \eta}{(1 + \alpha)^2}p_i \geq \frac{(1 - \eta)\Delta}{3d(1 + \alpha)^2}.$$

Using the definition of $L$ and the fact that $\eta < 1/2$, we can conclude the proof of the claim as follows:

$$1 + \alpha^{-L}\epsilon < (1 + \alpha)^{\log_1 + \alpha(6d(1 + \alpha)^2\epsilon/\Delta)}\epsilon \leq \frac{\Delta}{6d(1 + \alpha)^2} < \frac{(1 - \eta)\Delta}{3d(1 + \alpha)^2}. \qquad \square$$

Putting (10), (11) and the preceding claims together, we find that, with probability at least $1 - \delta$, a hypothesis $H$ is delivered such that

$$D\{i: H(i) \neq C(i)\} < D(I_{\text{wrong}}) + \frac{\Delta}{3}$$

$$\leq \frac{f_{\text{maj}} + f_{\text{sparse}} + f_{\text{dense}} + \Delta/3}{1 - \eta} + \frac{\Delta}{3}$$

$$\leq \frac{\eta + \Delta}{1 - \eta}$$

$$\leq \frac{\eta_0}{1 - \eta_0}$$

$$= \epsilon.$$

This concludes the proof.  □

We want to use the algorithm for powerset $\mathscr{C}_N$ and legal distributions as a subroutine to learn other concept classes. For this purpose, we need the following observations concerning the proof of Theorem 3.10:

—The total error on heavy points is bounded above (in probability) by $(f_{heavy} + \Delta/3)/(1 - \eta)$, where $f_{heavy}$ denotes the sample frequency of corrupted examples associated with heavy points.

—The required sample size does not change its order of magnitude if heavy points $i$ were only "heavy up to a constant", that is $D(i) \geq \Delta/(kd)$ for some fixed but arbitrary constant $k \geq 3$.

Imagine that the learning algorithm gets, as auxiliary information in addition to the corrupted sample, a partition of the domain $X$ into bins $B_1, \ldots, B_N$ such that the following holds:

—$D(B_i) \leq \Delta/(3d)$ for each heterogeneous bin $B_i$. Here $d$ denotes the VC dimension of the target class and $D$ the domain distribution.

—There are at most $d$ heterogeneous bins.

—Each homogeneous bin $B_i$ satisfies at least one of the following properties:

  —*Almost Heavy Bin.* $D(B_i) \geq \Delta/(144d)$.

  —*Unprofitable Bin.* $f_i \geq D(B_i)$, where $f_i$ denotes the sample frequency of corrupted examples associated with sample points hitting $B_i$.

  —*Nice Bin.* The corrupted sample shows the true label of $B_i$ with a strong majority over the wrong label.

It should be clear from the proof of Theroem 3.10 that, given a sufficiently large corrupted sample and the auxiliary information, RMD-POW can be successfully applied with heterogeneous bins in the role of light points and almost heavy homogeneous bins in the role of heavy points. For nice bins, RMD-POW will find the correct label. On unprofitable bins, RMD-POW might fail, but this does not hurt because the adversary's "investment" was too large. These considerations lead to the following:

COROLLARY 3.15. *There exists a randomized algorithm, RMD-kINV, which achieves the following. For every $k \geq 1$ and every $1 \geq \epsilon, \delta, \Delta > 0$, RMD-kINV PAC learns the class $\mathcal{I}_k$ with accuracy $\epsilon$, confidence $\delta$, tolerating malicious noise rate $\eta = \epsilon/(1 + \epsilon) - \Delta$, and using a sample of size $\tilde{O}(\epsilon/\Delta^2 + d/\Delta)$ with $d = 2k$.*

PROOF. Note that concept class $\mathcal{I}_k$ has VC dimension $d = 2k$. According to the preceding discussion, it suffices to show that an appropriate partition of $[0, 1]$ into bins can be computed with probability at least $1 - \delta/2$ of success, using a (corrupted) sample of size $2m = \tilde{O}(\epsilon/\Delta^2 + d/\Delta)$.

The bin partition is computed from the given corrupted sample $S$ in two stages. Stage 1 uses the first $m$ sample points (subsample $S_1$) and Stage 2 the last $m$ sample points (subsample $S_2$).

*Stage* 1. Each point $i \in [0, 1]$ occuring in $S_1$ with a relative frequency of at least $\Delta/(24d)$ is put into a separate *point-bin*. Let $B_1, \ldots, B_{N_1}$ denote the resulting sequence of point-bins. The points in subdomain $X_2 := [0, 1]\backslash B_1 \cup \cdots \cup B_{N_1})$ are called *empirically light*.

*Stage* 2. In order to put the points from $X_2$ into bins, a sorted list $S_2'$ (in increasing order), containing all points from $X_2$ occuring at least once in $S_2$, is processed from left to right. We keep track of a variable $f$ which is initialized to zero and sums up the subsample frequencies (with respect to $S_2$) of the points

from $S_2'$ that are already processed. As soon as $f$ reaches threshold $\Delta/(24d)$, we complete the next bin and reset $f$ to zero. For instance, if $f$ reaches this threshold for the first time when $z_1 \in S_2'$ is currently processed, we create bin $[0, z_1]$, reset $f$ to zero, and proceed with the successor of $z_1$ in $S_2'$. When $f$ reaches the threshold for the second time when $z_2' \in S_2'$ is currently processed, bin $(z_1, z_2]$ is created, and so on. The last bin $B_{N_2}'$ has endpoint $z_{N_2} = 1$ (belonging to $B_{N_2}'$ iff 1 is empirically light).[3]

Note that the bin partition does only depend on the instances within $S$ and not on their labels. (Clearly, the labels are used when we apply procedure RMD-POW as a subroutine in order to determine the bin labels.)

If sample size $2m$ is appropriately chosen from $\tilde{O}(\epsilon/\Delta^2 + d/\Delta)$, the following holds with probability at least $1 - \delta/2$ of success:

*Condition* 1.    $\hat{\eta} \le \eta + \Delta/3 < 1/2$, where $\hat{\eta} = \max\{\hat{\eta}_1, \hat{\eta}_2\}$ and $\hat{\eta}_i$ denotes the fraction of corrupted examples within subsample $S_i$.

*Condition* 2.    $D(i) < \Delta/(6d)$ for all empirically light points $i$.

*Condition* 3.    For every bin $B$ created in Stage 2 we have $D(B) < \Delta/(3d)$.

For Condition 1, we can argue as in the proof of Theorem 3.10.

Loosely speaking, Condition 2 follows because the adversary cannot make points appearing substantially lighter than they actually are. A more formal argument is as follows. The class of singletons over domain $[0, 1]$ has VC dimension 1. Assume $D(i) \ge \Delta/(6d)$. An easy application of Lemma A.2 shows (in probability) that $i$ occurs in the uncorrupted subsample for Stage 1 (whose corruption yields $S_1$) more than $\Delta/(12d)$ times. Since $\hat{\eta}_1 < 1/2$, it occurs in the corrupted subsample $S_1$ more than $\Delta/(24d)$ times. This implies that each empirically light point has probability less than $\Delta/(6d)$.

Condition 3 can be seen as follows: The class of intervals over subdomain $X_2$ has VC dimension 2. Consider a bin $B_j'$ created in Stage 2. Let $B_j'' := B_j'\backslash\{z_j\}$.

Note that the relative frequency of hitting $B_j''$ with points from $S_2$ is smaller than $\Delta/(24d)$. Applying the same kind of reasoning as for empirically light points in Stage 1, it follows (in probability) that $D(B_j'') < \Delta/(6d)$. If $z_j$ belongs to $B_j'$ (which is always the case unless perhaps $j = N_2$ and $z_j = 1$), then $z_j$ is empirically light. Applying Condition 2, we get $D(B_j') < \Delta/(3d)$.

We conclude the proof by showing that these three conditions imply that the bin partition has the required properties. Clearly, all point-bins are homogeneous. Since the at most $k$ intervals of the target concept have at most $d = 2k$ endpoints, at most $d$ of the bins created in Stage 2 are heterogeneous. Thus, there are at most $d$ heterogeneous bins altogether, and each of them has probability at most $\Delta/(3d)$. Each point-bin is hit by $S_1$ with a relative frequency of at least $\Delta/(24d)$. Similarly, each bin created in Stage 2 is hit by $S_2$ with a relative frequency of at least $\Delta/(24d)$. Since $|S_1| = |S_2| = m$, each homogeneous bin is hit by $S$ with a relative frequency of at least $\Delta/(48d)$. Let $B$ be a homogeneous bin. Remember that we consider $B$ as almost heavy if $D(B) \ge \Delta/(144d)$. Assume $B$ is not almost heavy. If the sample frequency of corrupted examples hitting $B$ is at least $\Delta/(144d)$, then $B$ is unprofitable. Assume this is

---

[3] If variable $f$ is below its threshold when endpoint 1 is reached, the last bin must formally be treated like a heterogeneous bin. For the sake of simplicity, we ignore this complication in the sequel.

not the case. But then $S$ presents the true label of $B$ with a relative frequency of at least $\Delta/(48d) - \Delta/(144d) = \Delta/(24d)$, twice as often as the wrong label. It follows that $B$ is a nice bin. This concludes the proof of the corollary. $\square$

For class $\mathcal{I}_k$, we used a bin class of constant VC dimension with at most $d$ heterogeneous bins. We state without proof that a partition using a bin class of VC dimension $d_1$ and at most $d_2$ heterogeneous bins leads to an application of algorithm RMD-POW requiring an additional term of order $\tilde{O}(d_1 d_2/\Delta)$ in the sample size. As long as $d_1 d_2 = O(d)$, the sample size has the same order of magnitude as in Corollary 3.15.

## 4. *Malicious Noise and Randomized Hypotheses*

In this section, we investigate the power of randomized hypotheses for malicious PAC learning. We start by observing that an easy modification of Kearns and Li [1993, Theorem 1] yields the following result. (Recall that a target class is nontrivial if it contains two functions $C$ and $C'$ and there exist two distinct points $x, x' \in X$ such that $C(x) = C'(x) = 1$ and $C(x') \neq C'(x')$.)

PROPOSITION 4.1. *For all nontrivial target classes $\mathcal{C}$ and all $\epsilon < 1/2$, no algorithm can learn $\mathcal{C}$ with accuracy $\epsilon$, even using randomized hypotheses, and tolerating malicious noise rate $\eta \geq 2\epsilon/(1 + 2\epsilon)$.*

Let $\eta_{\mathrm{rand}} = \eta_{\mathrm{rand}}(\epsilon) = 2\epsilon/(1 + 2\epsilon)$ (we omit the dependence on $\epsilon$ when it is clear from the context.) As the corresponding information-theoretic bound $\eta_{\mathrm{det}} = \epsilon/(1 + \epsilon)$ for learners using deterministic hypotheses is strictly smaller than $\eta_{\mathrm{rand}}$, one might ask whether this gap is real, that is, whether randomized hypotheses really help in this setting. In Subsection 4.1, we give a positive answer to this question by showing a general strategy that, using randomized hypotheses, learns any target class $\mathcal{C}$ tolerating any malicious noise rate $\eta$ bounded by a constant fraction of $(7/6)\epsilon/(1 + (7/6)\epsilon)$ and using sample size $\tilde{O}(d/\epsilon)$, where $d$ is the VC dimension of $\mathcal{C}$. Note that $(7/6)\epsilon/(1 + (7/6)\epsilon) > \eta_{\mathrm{det}}$, whereas no learner using deterministic hypotheses can tolerate a malicious noise rate $\eta \geq \eta_{\mathrm{det}}$, even allowing an *infinite* sample. Furthermore, the sample size used by our strategy is actually *independent* of $\eta$ and is of the same order as the one needed in the noise-free case. Finally, in Subsection 4.2 we show an algorithm for learning the powerset of $d$ points, for every $d \geq 1$, with malicious noise rates arbitrarily close to $\eta_{\mathrm{rand}}$.

The problem of finding a general strategy for learning an arbitrary concept class with randomized hypotheses and malicious noise rate arbitrarily close to $2\epsilon/(1 + 2\epsilon)$ remains open.

4.1. A GENERAL UPPER BOUND FOR LOW NOISE RATES.    We show the following result.

THEOREM 4.2. *For every target class $\mathcal{C}$ with VC dimension $d$, every $0 < \epsilon, \delta \leq 1$, and every fixed constant $0 \leq c < 7/6$, a sample size of order $d/\epsilon$ (ignoring logarithmic factors) is necessary and sufficient for PAC learning $\mathcal{C}$ using randomized hypotheses, with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = c\epsilon/(1 + c\epsilon)$.*

Input: Target class $\mathcal{C}$, sample $S$. Accuracy and confidence
parameters $\varepsilon, \delta$, upper bound $\eta < \dfrac{(7/6)\varepsilon}{1 + (7/6)\varepsilon}$ on true malicious noise.

Initialization: Let $\gamma = \dfrac{(7/6) - c}{(7/6) + c}$, where $c = \dfrac{\eta}{(1 - \eta)\varepsilon}$.

**Phase 1.**

   (1) Remove from $\mathcal{C}$ all concepts $H$ such that $\widehat{\mathrm{err}}(H) > (1 + \gamma)\eta$.
       Let $\mathcal{K}$ be the set of remaining concepts.
   (2) If $\mathcal{K}$ is empty, then output a default hypothesis.

**Phase 2.**

   (1) If there exists an independent set $\{F, G, H\} \subseteq \mathcal{K}$, then output the majority vote
       of $F$, $G$, and $H$.
   (2) Otherwise, pick a maximal independent set $\mathcal{U}$ of size 1 or 2.
       If $\mathcal{U} = \{H\}$, then output $H$. If $\mathcal{U} = \{G, H\}$, then output the coin rule $\frac{1}{2}(G + H)$.

FIG. 2. A description of the randomized algorithm SIH used in the proof of Theorem 4.2.

This, combined with Theorems 3.4 and 3.7 shows that, if the noise rate is about $\eta_{\mathrm{det}}$, then a sample size of order $d/\Delta + \epsilon/\Delta^2$ is only needed if the final hypothesis has to be deterministic.

The idea behind the proof of Theorem 4.2 is the following: A sample size $\tilde{O}(d/\epsilon)$ is too small to reliably discriminate the target (or other $\epsilon$-good concepts) from $\epsilon$-bad concepts. (Actually, the adversary can make $\epsilon$-bad hypotheses perform better on the sample than the target $C$.) It is, however, possible to work in two phases as follows. Phase 1 removes some concepts from $\mathcal{C}$. It is guaranteed that all concept with an error rate "significantly larger" than $\epsilon$ are removed, and that the target $C$ is not removed. Phase 2 reliably checks whether two concepts are independent in the sense that they have a "small" joint error probability (the probability to produce a wrong prediction on the same randomly drawn example). Obviously, the majority vote of three pairwise independent hypotheses has error probability at most 3 times the "small" joint error probability of two independent concepts. If there is no independent set of size 3, there must be a maximal independent set $\mathcal{U}$ of size 2 (or 1). We will show that each maximal independent set contains a concept with error probability significantly smaller than $\epsilon$. It turns out that, if either $\mathcal{U} = \{G\}$ or $\mathcal{U} = \{G, H\}$, then either $G$ or the coin rule $1/2(G + H)$, respectively, is $\epsilon$-good. The resulting algorithm, SIH (which stands for Search Independent Hypotheses), is described in Figure 2.

PROOF (OF THEOREM 4.2). The lower bound is obvious because it holds for the noise-free case [Ehrenfeucht et al. 1989]. For the upper bound, we begin with the following preliminary considerations. Let $X$ be the domain of target class $\mathcal{C}$, $D$ any distribution on $X$, and $C \in \mathcal{C}$ be the target. Given a hypothesis $H \in \mathcal{C}$, $E(H) = \{x: H(x) \neq C(x)\}$ denotes its *error set*, and $\mathrm{err}(H) = D(E(H))$ its error probability. The *joint error probability* of two hypotheses $G, H \in \mathcal{C}$ is given by $\mathrm{err}(G, H) = D(E(G) \cap E(H))$. Our proof will be based on the fact that (joint) error probabilities can be accurately empirically estimated. Let $S$ be the sample. We denote the relative frequency of mistakes of $H$ on the whole sample $S$ by $\widehat{\mathrm{err}}(H)$. The partition of $S$ into a clean and a noisy part leads to the decomposition $\widehat{\mathrm{err}}(H) = \widehat{\mathrm{err}}^c(H) + \widehat{\mathrm{err}}^n(H)$, where upper indices $c$ and $n$ refer to

the clean and the noisy part of the sample, respectively. Note that term $\widehat{err}^c(H)$ is an empirical estimation of $\widehat{err}(H)$, whereas $\widehat{err}^n(H)$ is under control of the adversary. The terms $\widehat{err}(G, H)$, $\widehat{err}^c(G, H)$, and $\widehat{err}^n(G, H)$ are defined analogously.

Let $\gamma, \lambda > 0$ denote two fixed constants to be determined by the analysis and let $\hat{\eta}$ denote the empirical noise rate. A standard application of the Chernoff–Hoeffding bound (23) and Lemma A.2 shows that, for a suitable choice of $m = \tilde{O}(d/\epsilon)$, the following conditions are simultaneously satisfied with probability $1 - \delta$:

*Condition* 1.    $\hat{\eta} \leq (1 + \gamma)\eta$.

*Condition* 2.    If $H \in \mathscr{C}$ satisfies $err(H) \geq \lambda\epsilon$, then

$$\widehat{err}^c(H) \geq (1 - \gamma)(1 - \eta)err(H).$$

*Condition* 3.    If $G, H \in \mathscr{C}$ satisfy $err(G, H) \geq \lambda\epsilon$, then

$$\widehat{err}^c(G, H) \geq (1 - \gamma)(1 - \eta)err(G, H).$$

We just mention that for proving Conditions 2 and 3 we use the fact that the VC dimensions of the classes of the error sets and the joint error sets are both $O(d)$.

We now describe the learning algorithm SIH illustrated in Figure 2.

In Phase 1, all concepts $H \in \mathscr{C}$ satisfying $\widehat{err}(H) > (1 + \gamma)\hat{\eta}$ are removed. Let $\mathscr{H} = \{H \in \mathscr{C} : \widehat{err}(H) \leq (1 + \gamma)\hat{\eta}\}$ denote the set of remaining concepts. Note that $\widehat{err}(C) \leq \hat{\eta}$. Applying Condition 1, it follows that target $C$ belongs to $\mathscr{H}$. Applying Condition 2 with constant $\lambda \leq c(1 + \gamma)/(1 - \gamma)$, it follows that all concepts $H \in \mathscr{H}$ satisfy:

$$err(H) \leq \frac{(1 + \gamma)\eta}{(1 - \gamma)(1 - \eta)}. \tag{15}$$

We are now in position to formalize the notion of independence, which is central for Phase 2 of SIH. Let us introduce another parameter $\alpha$ whose value will also be determined by the analysis. We say that $G, H \in \mathscr{H}$ are *independent* if $\widehat{err}(G, H) \leq (\alpha + \gamma)\hat{\eta}$. A subset $\mathscr{U} \subseteq \mathscr{H}$ is called independent if its hypotheses are pairwise independent.

CLAIM.    *If $\widehat{err}^c(H) \geq (1 - \alpha)\hat{\eta}$, then $H$ and $C$ are independent.*

To prove the claim note that, since $C$ is the target, $\widehat{err}(H, C) \leq \widehat{err}^n(H)$. The definition of $\mathscr{H}$ and the decomposition of $\widehat{err}$ into $\widehat{err}^c$ and $\widehat{err}^n$ imply that

$$\widehat{err}^n(H) = \widehat{err}(H) - \widehat{err}^c(H) \leq (1 + \gamma)\hat{\eta} - (1 - \alpha)\hat{\eta} = (\alpha + \gamma)\hat{\eta},$$

proving the claim.

From Claim and Condition 2 applied with $\lambda \leq c(1 - \alpha)/(1 - \gamma)$ we obtain the following facts:

*Fact* 1.   If

$$\text{err}(H) \geq \frac{(1 - \alpha)\eta}{(1 - \gamma)(1 - \eta)},$$

then $H$ and target $C$ are independent.

*Fact* 2.   Each maximal independent set $\mathcal{U} \subseteq \mathcal{H}$ contains at least one hypothesis whose error is smaller than $(1 - \alpha)\eta/((1 - \gamma)(1 - \eta))$. In particular, if $\mathcal{U} = \{H\}$, then

$$\text{err}(H) \leq \frac{(1 - \alpha)\eta}{(1 - \gamma)(1 - \eta)}.$$

If $\mathcal{U} = \{G, H\}$, then one among the two quantities $\text{err}(G)$ and $\text{err}(H)$ is smaller than or equal to $(1 - \alpha)\eta/((1 - \gamma)(1 - \eta))$.

We now move on to the description of Phase 2 (see Figure 2). Note that Phase 2 of SIH either terminates with a deterministic hypothesis, or terminates with the coin rule $1/2(G + H)$. The following case analysis will show that the final hypothesis output by SIH is $\epsilon$-good unless it is the default hypothesis.

Let us first consider the case that the final hypothesis is the majority vote $\text{MAJ}_{F,G,H}$ of three independent hypothesis $F$, $G$, and $H$ (Step 1 in Phase 2). Then an error occurs exactly on those instances $x$ that are wrongly predicted by at least two hypotheses of $F$, $G$, $H$, that is

$$\text{err}(\text{MAJ}_{F,G,H}) \leq \text{err}(F, G) + \text{err}(F, H) + \text{err}(G, H).$$

By definition of independence, we know that

$$\widehat{\text{err}}^c(X, Y) \leq \widehat{\text{err}}(X, Y) \leq (\alpha + \gamma)\eta$$

for each pair $(X, Y)$ of distinct hypothesis in $\{F, G, H\}$. Then, observing that $c\epsilon = \eta/(1 - \eta)$ and applying Condition 3 to each such pair with $\lambda \leq c(\alpha + \gamma)/(1 - \gamma)$, we get that

$$\text{err}(\text{MAJ}_{F,G,H}) \leq \frac{3(\alpha + \gamma)\eta}{(1 - \gamma)(1 - \eta)}. \tag{16}$$

If the final hypothesis is the coin rule $1/2(G + H)$ (from Step 2 in Phase 2), we may apply (15) and Fact 2 to bound the error probability as follows:

$$\text{err}\left(\frac{1}{2}(G + H)\right) = \frac{1}{2}(\text{err}(G) + \text{err}(H))$$

$$< \frac{1}{2}\left(\frac{(1 - \alpha)\eta}{(1 - \gamma)(1 - \eta)} + \frac{(1 + \gamma)\eta}{(1 - \gamma)(1 - \eta)}\right). \tag{17}$$

**Initialization**
Input: A sample $(x_1, l_1), \ldots, (x_m, l_m)$, $x_i \in \{1, \ldots, d\}$, $l \in \{0, 1\}$.
For $i := 1, \ldots, d$ Compute

$$p_0(i) := \frac{|\{j \mid x_j = i,\ l_j = 0\}|}{m}$$

$$p_1(i) := \frac{|\{j \mid x_j = i,\ l_j = 1\}|}{m}$$

$$h(i) := \frac{(p_1(i))^2}{(p_0(i))^2 + (p_1(i))^2}$$

**Prediction**
Input: Some $i \in \{1, \ldots, d\}$
Output: Label 1 with probability $h(i)$ and label 0 with probability $1 - h(i)$.

FIG. 3. Pseudo-code for Algorithm SQ-RULE. In the initialization phase some parameters are computed from the sample. These are then used in the working phase to make randomized predictions.

If the final hypothesis is $H$ (from Step 2), then

$$\text{err}(H) < \frac{(1 - \alpha)\eta}{(1 - \gamma)(1 - \eta)}$$

by Fact 2. This error bound is smaller than the bound (17). We now have to find choices for the parameters $\alpha$, $\gamma$, $\lambda$ such that (16) and (17) are both upper bounded by $\epsilon$ and the previously stated conditions

$$\lambda \leq \frac{c(1 + \gamma)}{(1 - \gamma)}, \ \lambda \leq \frac{c(1 - \alpha)}{(1 - \gamma)}, \ \lambda \leq \frac{c(\alpha + \gamma)}{(1 - \gamma)}$$

on $\lambda$ hold. Equating bounds (16) and (17) and solving for $\alpha$ gives $\alpha = 2/7 - 5/7\gamma$. Substituting this into (16), setting the resulting formula to $\epsilon$ and solving for $\gamma$ yields $\gamma = ((7/6) - c)/((7/6) + c)$. (Observe that the choice of $\eta = c\epsilon/(1 + (c\epsilon)$ implies the equation $\eta/(1 - \eta) = c\epsilon$.) This in turn leads to the choice $\alpha = 3 (c - (1/2))/(c + (7/6))$. According to these settings, one can finally choose $\lambda = 1/3$.  □

REMARK 4.3.    *The result of Theorem 4.2 gives rise to a challenging combinatorial problem: Given a target class, find 3 independent hypotheses, or alternatively, a maximal independent set of less than 3 hypotheses. This replaces the "consistent hypothesis" paradigm of noise-free PAC learning and the "minimizing disagreement" paradigm of agnostic learning. There are examples of concept classes such that, for certain samples, one can find three or more independent hypotheses.*

4.2  AN ALMOST OPTIMAL COIN RULE FOR THE POWERSET.    In this subsection, we introduce and analyze a simple algorithm called Square Rule (SQ-RULE) for learning with coin rules the powerset $\mathscr{C}_d$ of $d$ elements in presence of a malicious noise rate arbitrarily close to $\eta_{\text{rand}}$ and using almost optimal sample size. (See also Figure 3.) Algorithm SQ-RULE works as follows: Let $H(p, q) = q^2/(p^2 +$
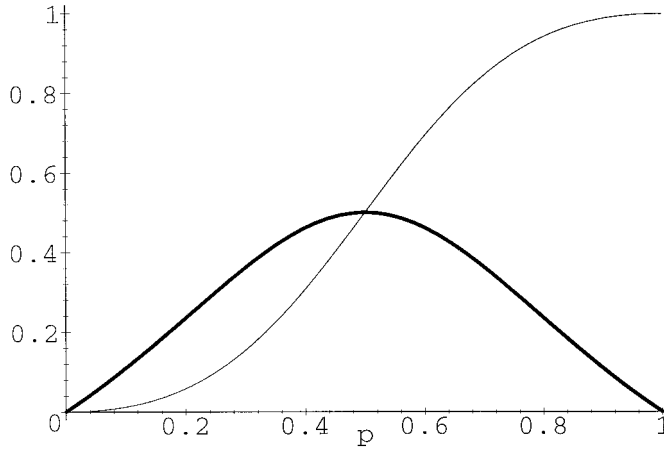
FIG. 4. Curve of the coin rule $H(p, q) = q^2/(p^2 + q^2)$ (thin) and of the return to investment ratio $H(p, q) \cdot p/q$ (thick); the curves are scaled to $p + q = 1$.

$q^2$). For a given sample $S$, let $\hat{p}_0(x)$ and $\hat{p}_1(x)$, respectively, denote the relative frequency of $(x, 0)$ and $(x, 1)$. On input $S$, SQ-RULE outputs the coin rule[4] $F(x) = H(\hat{p}_0(x), \hat{p}_1(x)) = \hat{p}_1(x)^2/(\hat{p}_0(x)^2 + \hat{p}_1(x)^2)$. We now show a derivation of the coin rule $F$. Consider a single point $x$ and let $\hat{p}_0 = \hat{p}_0(x)$ and $\hat{p}_1 = \hat{p}_1(x)$. Now, if the true label of $x$ is 0, we say that the "empirical return" of the adversary is $\hat{p}_0$. Likewise, the adversary's "investment" for the false label 1 is $\hat{p}_1$. Also, as $F$ incorrectly classifies $x$ with probability $H(\hat{p}_0, \hat{p}_1)$, the "empirical return to investment ratio", which we denote by $\rho$, is $H(\hat{p}_0, \hat{p}_1) \cdot \hat{p}_0/\hat{p}_1$. Similarly, if the true label of $x$ is 1, then $\rho$ is $(1 - H(\hat{p}_0, \hat{p}_1)) \cdot \hat{p}_1/\hat{p}_0$. The function $F$ that minimizes $\rho$ over all choices of $\hat{p}_0$ and $\hat{p}_1$ is found by letting $H(\hat{p}_0, \hat{p}_1) \cdot \hat{p}_0/\hat{p}_1 = (1 - H(\hat{p}_0, \hat{p}_1)) \cdot \hat{p}_1) \cdot \hat{p}_1/\hat{p}_0$ and solving for $H$ to obtain $H(\hat{p}_0, \hat{p}_1) = \hat{p}_1^2/(\hat{p}_0^2 + \hat{p}_1^2)$. (A plot of the functions $H$ and $\rho$ is shown in Figure 4. In the same figure we also plot the return to investment ratio, showing that the best strategy for the adversary is to balance the labels whence this ratio becomes 1/2.) Note that, as our final goal is to bound the quantity $\mathbf{E}_{x \sim D}|F(x) - C(x)|$, we should actually choose $F$ so to minimize the *expected* return to investment ratio, that is, the ratio $|H(\hat{p}_0, \hat{p}_1) - C(x)| \cdot D(x)/\hat{p}_{1-C(x)}$. However, as we will show in a moment, an estimate of the unknown quantity $D(x)$ will suffice for our purposes.

THEOREM 4.4. *For every $d \geq 1$ and every $0 < \epsilon, \delta, \Delta \leq 1$, algorithm SQ-RULE learns the class $\mathscr{C}_d$ with accuracy $\epsilon$, confidence $\delta$, tolerating malicious noise rate $\eta = 2\epsilon/(1 + 2\epsilon) - \Delta$, and using a sample of size $O(d\epsilon/\Delta^2)$.*

PROOF. We start the proof with the following preliminary considerations. Let $X = \{1, \ldots, d\}$, let $D$ be any distribution on $X$, and let $C$ be the target function. Let $t(x) = (1 - \eta)D(x)$, that is, $t(x)$ denotes the probability that $x$ is presented by the adversary with the true label $C(x)$. Fix a sample $S$. The relative frequency of $(x, C(x))$ in $S$ is denoted by $\hat{t}(x)$. We assume without loss of generality, that the adversary does never present an instance with its true label in

---

[4] The function $H$ was also used by Kearns et al. [1994] in connection with agnostic learning.

noisy trials. (The performance of the coin rule $F$ gets better in this case.) We denote the relative frequency of $(x, 1 - C(x))$ in $S$ by $\hat{f}(x)$. The relative frequency of noisy trials in $S$ is denoted by $\hat{\eta}$. Clearly, $\hat{\eta} = \Sigma_{x \in X} \hat{f}(x)$. Applying Lemmas A.1 and A.2, it is not hard to show that there exists a sample size $m = \bar{O}(d\epsilon/\Delta^2)$ such that with probability $1 - \delta$ the sample $S$ satisfies the following conditions:

$$\hat{\eta} \leq \eta + \frac{\Delta}{2}, \tag{18}$$

$$\forall x \in X : t(x) \geq \frac{\Delta}{24d} \Rightarrow \hat{t}(x) \geq \frac{t(x)}{2}, \tag{19}$$

$$\forall M \subseteq X : \sum_{x \in M} t(x) \leq 16\epsilon \Rightarrow \sum_{x \in M} \hat{t}(x) \geq \sum_{x \in M} t(x) - \frac{\Delta}{8}. \tag{20}$$

To prove (18), we apply (25) with $p' = \eta$ and $\lambda = \Delta/(2\eta)$. To prove (19), we apply (23) with $p = \Delta/(24d)$ and $\lambda = 1/2$. Finally, to prove (20), we use (23) to find that

$$\Pr\left\{ S_m \leq \left( 1 - \frac{\lambda}{p} \right) mp \right\} \leq \exp\left( -\frac{\lambda^2 m}{2p} \right) \leq \exp\left( -\frac{\lambda^2 m}{2p'} \right),$$

where the last inequality holds for all $p' \geq p$ by monotonicity. Setting $\lambda = \Delta/8$ and $p' = 16\epsilon$ concludes the proof of (20). These three conditions are assumed to hold in the sequel. An instance $x$ is called *light* if $t(x) < \Delta/(24d)$, and *heavy* otherwise. Note that $\eta < \eta_{\text{rand}} \leq 2/3$ (recall that $\eta_{\text{rand}} = 2\epsilon/(1 + 2\epsilon)$.) Thus, $D(x) < \Delta/(8d)$ for all light points. The total contribution of light points to the error probability of the coin rule $F$ is therefore less than $\Delta/8$. The following analysis focuses on heavy points; note that for these points the implication in (19) is valid. We will show that the total error probability on heavy points is bounded by $\epsilon - \Delta/8$.

It will be instructive to consider the error probability on $x$ of our coin rule $F$ as the return of the adversary at $x$ (denoted by RETURN$(x)$ henceforth) and the quantity $\hat{f}(x)$, defined above, as its investment at $x$. Our goal is to show that the total return of the adversary is smaller than $\epsilon - \Delta/8$, given that its total investment is $\hat{\eta}$. The function $R(p, q) = pq/(p^2 + q^2)$ plays a central role in the analysis of the relation between return and investment. (A plot of this function is shown in Figure 4.) Function $R$ attains its maximal value $1/2$ for $p = q$. For $q \leq p/4$ or $p \leq q/4$, the maximal value is $4/17$, see also Figure 5. Before bounding the total return, we will analyze the term RETURN$(x)$.

If $C(x) = 0$, then

$$\hat{f}(x) = \hat{p}_1(x), \hat{t}(x) = \hat{p}_0(x), \text{RETURN}(x) = F(x) \cdot D(x) = \frac{\hat{p}_1(x)^2 \cdot D(x)}{\hat{p}_0(x)^2 + \hat{p}_1(x)^2}.$$
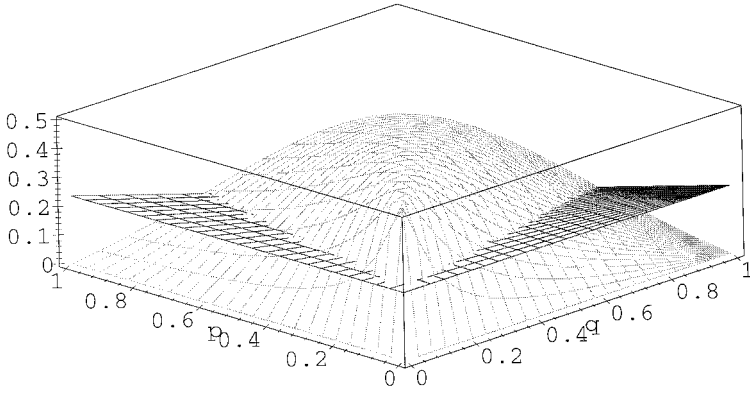
FIG. 5. Plot of the return curve $R(p, q) = (pq)/(p^2 + q^2)$ and of the constant function $c(p, q) = 4/17$, for $p, q \in [0, 1]$, $p + q = 1$.

If $C(x) = 1$, then

$$\hat{f}(x) = \hat{p}_0(x), \ \hat{t}(x) = \hat{p}_1(x), \ \text{RETURN}(x) = (1 - F(x)) \cdot D(x) = \frac{\hat{p}_0(x)^2 \cdot D(x)}{\hat{p}_0(x)^2 + \hat{p}_1(x)^2}.$$

Note that in both cases $\hat{f}(x) \cdot \hat{t}(x) = \hat{p}_0(x) \cdot \hat{p}_1(x)$ and $\text{RETURN}(x) = \hat{f}(x)^2 \cdot D(x)/(\hat{p}_0(x)^2 + \hat{p}_1(x)^2)$. Setting $\hat{\alpha}(x) = t(x) - \hat{t}(x)$, we obtain $D(x) = t(x)/(1 - \eta) = (\hat{t}(x) + \hat{\alpha}(x))/(1 - \eta)$. For the sake of simplicity, we will use the abbreviations

$$\hat{p}_0 = \hat{p}_0(x), \quad \hat{p}_1 = \hat{p}_1(x), \quad \hat{f} = \hat{f}(x), \quad \hat{t} = \hat{t}(x), \quad \hat{\alpha} = \hat{\alpha}(x).$$

With these abbreviations, $\hat{t} \cdot \hat{f} = \hat{p}_0 \cdot \hat{p}_1$ is valid and the term $\text{RETURN}(x)$ can be written as follows:

$$\text{RETURN}(x) = \frac{\hat{f}^2 \cdot (\hat{t} + \hat{\alpha})}{(\hat{p}_0^2 + \hat{p}_1^2) \cdot (1 - \eta)} = \frac{1}{1 - \eta} \left( \frac{\hat{t} \cdot \hat{f}}{\hat{p}_0^2 + \hat{p}_1^2} \hat{f} + \frac{\hat{\alpha} \cdot \hat{f}^2}{\hat{p}_0^2 + \hat{p}_1^2} \right) \qquad (21)$$

$$= \frac{1}{1 - \eta} \left( R(\hat{p}_0, \hat{p}_1) \cdot \hat{f} + \frac{\hat{\alpha} \cdot \hat{f}^2}{\hat{p}_0^2 + \hat{p}_1^2} \right)$$

We now bound separately each one of the last two terms in (21). If $\hat{f} \geq \hat{t}/4$, then

$$\frac{R(\hat{p}_0, \hat{p}_1) \cdot \hat{f}}{(1 - \eta)} \leq \frac{\hat{f}}{2(1 - \eta)}.$$

Furthermore, as either $\hat{f} = \hat{p}_0$ or $\hat{f} = \hat{p}_1$,

$$\frac{\hat{\alpha} \cdot \hat{f}^2}{(\hat{p}_0^2 + \hat{p}_1^2)(1 - \eta)} \leq \frac{\hat{\alpha}}{1 - \eta}.$$

If $\hat{f} < \hat{t}/4$, we bound the last term in (21) using (19) and get

$$\frac{\hat{\alpha} \cdot \hat{f}^2}{(\hat{p}_0^2 + \hat{p}_1^2)} = \frac{(t - \hat{t}) \cdot \hat{f}^2}{(\hat{p}_0^2 + \hat{p}_1^2)} \leq \frac{\hat{t} \cdot \hat{f}^2}{(\hat{p}_0^2 + \hat{p}_1^2)} = \frac{(\hat{p}_0 \cdot \hat{p}_1) \cdot \hat{f}}{(\hat{p}_0^2 + \hat{p}_1^2)} = R(\hat{p}_0, \hat{p}_1) \cdot \hat{f}.$$

Hence, using (21) and $R(\hat{p}_0, \hat{p}_1) \leq 4/17 < 1/4$ for $\hat{f} < \hat{t}/4$, we finally get

$$\text{RETURN}(x) \leq \frac{2R(\hat{p}_0, \hat{p}_1) \cdot \hat{f}}{1 - \eta} < \frac{\hat{f}}{2(1 - \eta)}.$$

Piecing the above together, we obtain

$$\text{RETURN}(x) \leq \frac{\hat{f}}{2(1 - \eta)} + \begin{cases} \dfrac{\hat{\alpha}}{(1 - \eta)} & \text{if} \quad \hat{f} \geq \hat{t}/4, \\ 0 & \text{otherwise.} \end{cases} \qquad (22)$$

We are now in the position to bound the total return on all heavy instances $x$. For the first term in the left-hand-side of (22) we obtain the bound

$$\frac{1}{2(1 - \eta)} \sum \hat{f}(x) \leq \frac{\hat{\eta}}{2(1 - \eta)}$$

where the sum is over all heavy $x$. The treatment of the second term in the left-hand-side of (22) is more subtle. Let $M$ denote the set of heavy instances $x$ where $\hat{f}(x) \geq \hat{t}(x)/4$. $D(M)$ is therefore bounded as follows:

$$\sum_{x \in M} t(x) \leq 2 \sum_{x \in M} \hat{t}(x) \leq 8 \sum_{x \in M} \hat{f}(x) \leq 8\hat{\eta} < 16\epsilon.$$

From (20), we conclude that:

$$\frac{1}{1 - \eta} \sum_{x \in M} \hat{\alpha}(x) \leq \frac{1}{1 - \eta} \sum_{x \in M} (t(x) - \hat{t}(x)) \leq \frac{\Delta}{8(1 - \eta)}.$$

A straightforward computation shows that

$$\frac{\hat{\eta}}{2(1 - \eta)} + \frac{\Delta}{8(1 - \eta)} \leq \epsilon - \frac{\Delta}{8}.$$

As the probability of all light points is at most $\Delta/8$, the expected error of SQ-RULE is at most $\epsilon$. This completes the proof of Theorem 4.4. $\square$

The upper bound of Theorem 4.4 has a matching lower bound (up to logarithmic factors). The proof, which is a somewhat involved modification of the proof of Theorem 3.9, is only sketched.

THEOREM 4.5. *For every target class $\mathscr{C}$ with VC dimension $d \geq 3$, for every $0 < \epsilon \leq 1/38$, $0 < \delta \leq 1/74$, and for every $0 < \Delta = o(\epsilon)$, the sample size needed by any strategy (even using randomized hypotheses) for learning $\mathscr{C}$ with accuracy $\epsilon$, confidence $\delta$, and tolerating malicious noise rate $\eta = 2\epsilon/(1 + 2\epsilon) - \Delta$, is $\Omega(d\epsilon/\Delta^2)$.*

TABLE I.  Survey on sample size results of learning in the presence of noise, ignoring logarithmic terms

| Noise model | classification | | malicious with deterministic hyp. | | malicious with randomized hyp. | |
|---|---|---|---|---|---|---|
| Theoretical limit on noise rate | $\eta_0 := 1/2$ | | $\eta_0 := \epsilon/(1+\epsilon)$ | | $\eta_0 := 2\epsilon/(1+2\epsilon)$ | |
| Type of bound | upper | lower | upper | lower | upper | lower |
| $\eta = \eta_0 - \Delta$ $\Delta \to 0$ | $\dfrac{d}{\epsilon\Delta^2}$ | $\dfrac{d}{\epsilon\Delta^2}$ | $\dfrac{d}{\epsilon} + \dfrac{\epsilon}{\Delta^2}$ | $\dfrac{d}{\epsilon} + \dfrac{\epsilon}{\Delta^2}$ | $\dfrac{d\epsilon}{\Delta^2}$ | $\dfrac{d\epsilon}{\Delta^2}$ |
| $\eta = c \cdot \eta_0$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | — | — |
| $\eta < \dfrac{(7/6)\epsilon}{1+(7/6)\epsilon}$ | — | — | — | — | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ |
| Min. Dis. $\eta = c \cdot \eta_0$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | $\dfrac{d}{\epsilon}$ | — | — |
| Min. Dis. $\eta = \eta_0 - \Delta$ $\Delta \to 0$ | $\dfrac{d}{\epsilon\Delta^2}$ | $\dfrac{d}{\epsilon\Delta^2}$ | $\dfrac{d\epsilon}{\Delta^2}$ | $\dfrac{d\epsilon}{\Delta^2}$ | — | — |

NOTE: Empty entries are unknown or are combinations that make no sense. The lower part contains the results for the minimum disagreement strategy. For the sake of comparison, we also add the corresponding results for classification noise.

PROOF.     One uses $d$ shattered points with a suitable distribution $D$ and shows that, with constant probability, there exists a constant fraction of these points which occur with a frequency much lower than expected. The measure according to $D$ of these points is $2\epsilon$, but the adversary can balance the occurrences of both labels on these points while using noise rate less than $\eta_{\mathrm{rand}}$. Hence, even a randomized hypothesis cannot achieve an error smaller than $\epsilon$.     □

REMARK 4.6.     *Algorithm* SQ-RULE *can be modified to learn the class* $\mathfrak{I}_k$ *of unions of at most k intervals. Similarly to Corollary 3.15, one first computes a suitable partition of the domain into bins and applies the algorithm for the powerset afterwards as a subroutine.*

## 5. *Summary*

Table I shows the known results on learning in the malicious and classification noise models. The latter is a noise model where independently for every example the label is inverted with probability $\eta < 1/2$.[5]

There are still a few problems open. One is the question whether the strong adversary in the lower bound proofs of Theorems 3.9 and 4.5 can be replaced by the weaker KL-adversary. Also it would be interesting to see whether the constant 7/6 in Theorem 4.2 can improved to arbitrary constants $0 \le c < 2$. It seem that both questions are not easy to answer.

---

[5] See, for example, Laird [1988] for a survey on this noise model and for the upper bound on the sample size in the case of finite concept classes. See Simon [1996a] for the lower bound on the sample size, and Simon [1996b] for a generalization of Laird's upper bound to arbitrary concept classes.

*Appendix A. Some Statistical and Combinatorial Relations*

Let $S_{m,p}$ and $S'_{m,p'}$ be the sums of successes in a sequence of $m$ Bernouilli trials each succeeding with probability respectively at least $p$ and at most $p'$.

LEMMA A.1.    *For all $0 < \lambda < 1$,*

$$Pr\{S_{m,p} \leq (1 - \lambda)mp\} \leq \exp\left(-\lambda^2 \frac{mp}{2}\right) \tag{23}$$

$$Pr\{S_{m,p} \leq m(p - \lambda)\} \leq \exp(-2\lambda^2 m), \tag{24}$$

$$Pr\{S'_{m,p'} \geq (1 + \lambda)mp'\} \leq \exp\left(-\lambda^2 \frac{mp'}{3}\right). \tag{25}$$

Let $\mathscr{C}$ be a target class of VC dimension $d$ over some domain $X$. Let $D$ be a distribution over $X$. Let $S$ be an unlabeled sample of size $m$ drawn from $X$ under $D$. For $C \in \mathscr{C}$ let $D_S(C) = |\{x \in S: x \in C\}|/m$, the empirical probability of $C$.

LEMMA A.2. [VAPNIK 1982; BLUMER ET AL. 1989].    *For every $0 < \epsilon, \gamma \leq 1$ and every $0 < \delta < 1$, the probability that there exists a $C \in \mathscr{C}$ such that $D(C) > \epsilon$ and $D_S(C) \leq (1 - \gamma)D(C)$ is at most $8(2m)^d\exp(-\gamma^2\epsilon m/4)$, which in turn is at most $\delta$ if*

$$m \geq max\left\{\frac{8}{\gamma^2\epsilon} \ln\left(\frac{8}{\delta}\right), \frac{16d}{\gamma^2\epsilon} \ln\left(\frac{16}{\gamma^2\epsilon}\right)\right\}.$$

LEMMA A.3. [SAUER 1972; SHELAH 1972].    *Let $\mathscr{C}$ be a target class over $X$ of VC dimension $d$. For all $(x_1, \ldots, x_m) \in X^m$*

$$|\{(C(x_1), \ldots, C(x_m)): C \in \mathscr{C}\}| \leq \sum_{i=0}^{d} \binom{m}{i}.$$

PROOF OF FACT 3.2.    We prove inequality (1), the proof of (2) is similar. We proceed by establishing a series of inequalities. We shall also use Stirling's formula

$$\sqrt{2\pi N}\left(\frac{N}{e}\right)^N < N! < \sqrt{2\pi N}\left(\frac{N}{e}\right)^N \exp\left(\frac{1}{12N}\right). \tag{26}$$

Using (26) one can lower bound the binomial coefficient $\binom{N}{Np}$ as follows (assuming that $N$ is a multiple of $1/p$, which will be justified later in the proof)

$$\binom{N}{Np} = \frac{N!}{(Np)!(Nq)!}$$

$$> \frac{\sqrt{2\pi N}\left(\dfrac{N}{e}\right)^N}{\sqrt{2\pi Np}\left(\dfrac{Np}{e}\right)^{Np}\exp\left(\dfrac{1}{12Np}\right)\sqrt{2\pi Nq}\left(\dfrac{Nq}{e}\right)^{Nq}\exp\left(\dfrac{1}{12Nq}\right)}$$

$$= \frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{Npq}}\frac{1}{p^{Np}q^{Nq}}\exp\left(-\frac{1}{12Npq}\right).$$

This leads to

$$\sqrt{Npq} > \frac{1}{\sqrt{2\pi p^{Np}\,q^{Nq}}}\binom{N}{Np}^{-1}\exp\left(-\frac{1}{12Npq}\right). \tag{27}$$

Bahadur [1960] proved the following lower bound on the tail of the binomial distribution, where $0 \le k \le N$,

$$\Pr\{S_{N,p} \ge k\} \ge \binom{N}{k}p^k q^{(N-k)} \cdot \frac{q(k+1)}{k+1-p(N+1)} \cdot \left(1 + \frac{Npq}{(k-Np)^2}\right)^{-1}. \tag{28}$$

In order to be able to apply (27), we first remove the "floors" in (1). To this end, we replace $p$ by $p' = p - \gamma$ (and $q$ by $q' = q + \gamma$) such that $\lfloor Np \rfloor = Np'$. Then $Np'$ is integer and $p' > p - (1/N)$. We shall also need the following observation.

$$pq = (p'+\gamma)(q'-\gamma) = p'q' + \gamma(q'-p') - \gamma^2 = p'q' + \gamma(q'-p'-\gamma)$$

$$= p'q' + \gamma(q'-p') < p'q' + \gamma < p'q' + \frac{1}{N}. \tag{29}$$

Then (29) and $N \ge 37/(pq)$ imply that

$$Np'q' > Npq - 1 \ge 36. \tag{30}$$

Hence, (1) can be lower bounded as follows

$$\Pr\{S_{N,p} \ge \lfloor Np \rfloor + \lfloor \sqrt{Npq-1}\rfloor\}$$

$$\ge \Pr\{S_{N,p'} \ge Np' + \lfloor \sqrt{Npq-1}\rfloor\}$$

$$\ge \Pr\left\{S_{N,p'} \ge Np' + \left\lfloor \sqrt{N\left(p'q'+\frac{1}{N}\right)-1}\right\rfloor\right\}$$

$$\ge \Pr\{S_{N,p'} \ge Np' + \lfloor \sqrt{Np'q'}\rfloor\}. \tag{31}$$

In order to bound (31) we apply inequality (28) with $k = Np' + \lfloor \sqrt{Np'q'} \rfloor$ and $p$ and $q$ being replaced by $p'$ and $q'$, respectively. The three factors in the right-hand side of (28), denoted by $F_1$, $F_2$ and $F_3$, are separately bounded as follows:

$$F_1 = \binom{N}{Np' + \lfloor \sqrt{Np'q'} \rfloor} p'^{Np' + \lfloor \sqrt{Np'q'} \rfloor} q'^{N - Np' - \lfloor \sqrt{Np'q'} \rfloor} \tag{32}$$

$$= \binom{N}{Np' + \lfloor \sqrt{Np'q'} \rfloor} p'^{Np'} q'^{Nq'} \left(\frac{p'}{q'}\right)^{\lfloor \sqrt{Np'q'} \rfloor} \tag{33}$$

$$F_2 = \frac{q'(Np' + \lfloor \sqrt{Np'q'} \rfloor + 1)}{Np' + \lfloor \sqrt{Np'q'} \rfloor + 1 - Np' - p'}$$

$$= \frac{Np'q' + q'(\lfloor \sqrt{Np'q'} \rfloor + 1)}{\lfloor \sqrt{Np'q'} \rfloor + q'} > \sqrt{Np'q'} \tag{34}$$

$$> \frac{1}{\sqrt{2\pi} p'^{Np'} q'^{Nq'}} \binom{N}{Np'}^{-1} \exp\left(-\frac{1}{12Np'q'}\right). \tag{35}$$

(The inequality (35) follows from (27).)

$$F_3 = \left(1 + \frac{Np'q'}{(Np' + \lfloor \sqrt{Np'q'} \rfloor - Np')^2}\right)^{-1}$$

$$= \left(1 + \frac{Np'q'}{\lfloor \sqrt{Np'q'} \rfloor^2}\right)^{-1}$$

$$= \frac{\lfloor \sqrt{Np'q'} \rfloor^2}{\lfloor \sqrt{Np'q'} \rfloor^2 + Np'q'}$$

$$> \frac{(\sqrt{Np'q'} - 1)^2}{(\sqrt{Np'q'} + 1)^2 + Np'q'}$$

$$= \frac{Np'q' - 2\sqrt{Np'q'} + 1}{2Np'q' - 2\sqrt{Np'q'} + 1}$$

$$> \frac{Np'q' - 2\sqrt{Np'q'}}{2Np'q' - 2\sqrt{Np'q'}}$$

$$= \frac{1}{2}\left(1 - \frac{\sqrt{Np'q'}}{Np'q' - \sqrt{Np'q'}}\right)$$

$$= \frac{1}{2}\left(1 - \frac{1}{\sqrt{Np'q'} - 1}\right). \tag{36}$$

The following calculation shows how the product of (33) and (35) can be lower bounded. For notational convenience, let $T = (\exp(1/(12Np'q')) \sqrt{2\pi})^{-1}$ and let $K = \lfloor \sqrt{Np'q'} \rfloor$.

$$F_1 \cdot F_2 > \frac{\binom{N}{Np' + K}(p'/q')^K}{\sqrt{2\pi}\binom{N}{Np'}\exp(1/(12Np'q'))}$$

$$= T\left(\frac{p'}{q'}\right)^K \cdot \frac{N!(Np')!(N - Np')!}{N!(Np' + K)!(N - Np' - K)!}$$

$$= T\left(\frac{p'}{q'}\right)^K \frac{(Nq' - K + 1)\cdots(Nq')}{(Np' + 1)\cdots(Np + K)}$$

$$= T\left(\frac{p'}{q'}\right)^K \left(\prod_{i=1}^{K}\frac{Nq' - K + i}{Np' + i}\right) \tag{37}$$

$$> T\left(\frac{p'}{q'}\right)^K \left(\frac{Nq'}{Np' + K}\right)^K \tag{38}$$

$$= T\left(\frac{Np'}{Np' + K}\right)^K$$

$$= T\left(\frac{Np' + K}{Np'}\right)^{-K}$$

$$= T\left(1 + \frac{K}{Np'}\right)^{-K}$$

$$\geq T\left(1 + \frac{\sqrt{Np'q'}}{Np'}\right)^{-\sqrt{Np'q'}}$$

$$= T\left(1 + \frac{q'}{\sqrt{Np'q'}}\right)^{-\sqrt{Np'q'}} \tag{39}$$

$$\geq T \exp(-q') \tag{40}$$

$$= \frac{1}{\sqrt{2\pi}\exp(1/(12Np'q'))}\exp(-q') \tag{41}$$

$$\geq \frac{1}{\sqrt{2\pi} \exp(1/(12 \cdot 36))} \exp(-1) \geq 0.14642\ldots \qquad (42)$$

In (41) and (42), we used that $Np'q' > 36$, by (30). For the step from (37) to (38), we assume that $Nq' - k \geq Np'$. If $Nq' - k < Np'$, the steps from (38) in the above calculation are replaced by the following:

$$T\left(\frac{p'}{q'}\right)^K \left(\prod_{i=1}^{K} \frac{Nq' - K + i}{Np' + i}\right) \qquad (43)$$

$$> T\left(\frac{p'}{q'}\right)^K \left(\frac{Nq' - K}{Np'}\right)^K \qquad (44)$$

$$= T\left(\frac{Nq' - K}{Nq'}\right)^K$$

$$= T\left(\frac{Nq' - \lfloor \sqrt{Np'q'} \rfloor}{Nq'}\right)^{\lfloor \sqrt{Np'q'} \rfloor}$$

$$\geq T\left(\frac{Nq' - \sqrt{Np'q'}}{Nq'}\right)^{\sqrt{Np'q'}}$$

$$= T\left(1 - \frac{p'}{\sqrt{Np'q'}}\right)^{\sqrt{Np'q'}} \qquad (45)$$

$$\geq \frac{1}{\sqrt{2\pi} \exp(1/(12Np'q'))} \frac{10}{11} \exp(-p') \qquad (46)$$

$$\geq \frac{1}{\sqrt{2\pi} \exp(1/(12 \cdot 36))} \frac{10}{11} \exp(-1) \geq 0.133112\ldots \qquad (47)$$

The step from (45) to (46) follows from an elementary analysis of the function $(1 - a/b)^b - 10/11\exp(-a)$. Using (47) (which is less than the bound in (42)) and (36), we can lower bound the product $F_1 F_2 F_3$ as follows:

$$F_1 \cdot F_2 \cdot F_3 > 0.133112 \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{\lfloor \sqrt{Np'q'} \rfloor - 1}\right)$$

$$\geq 0.066556 \cdot \left(1 - \frac{1}{\lfloor \sqrt{36} \rfloor - 1}\right) \qquad (48)$$

$$> 0.05324\cdots > \frac{1}{19}. \qquad (49)$$

For the step from (48) to (49), we again used inequality (30). □

REFERENCES

BAHADUR, R. 1960. Some approximations to the binomial distribution function. *Ann. Math. Stat. 31*, 43–54.

BAHADUR, R., AND RANGA-RAO, R. 1960. On deviations of the sample mean. *Ann. Math. Stat. 31*, 1015–1027.

BLUMER, A., EHRENFEUCHT, A., HAUSSLER, D., AND WARMUTH, M. K. 1989. Learnability and the Vapnik-Chervonenkis dimension. *J. ACM 36*, 4, (Oct.), 929–965.

CHOW, Y. S., AND TEICHER, H. 1988. *Probability Theory*. Springer-Verlag, New York.

EHRENFEUCHT, A., HAUSSLER, D., KEARNS, M., AND VALIANT, L. 1989. A general lower bound on the number of examples needed for learning. *Inf. Comput. 82*, 3, 247–261.

GENTILE, C., AND HELMBOLD, D. 1998. Improved lower bounds for learning from noisy examples: An information-theoretic approach. In *Proceedings of the 11th Workshop on Computational Learning Theory (CoLT '98)* (Madison, Wisc., July 24–26). ACM, New York, pp. 104–115.

JOGDEO, K., AND SAMUELS, S. M. 1968. Monotone convergence of binomial probabilities and a generalization of ramanujan's equation. *Ann. Math. Sta. 39*, 4, 1191–1195.

KEARNS, M., AND LI, M. 1993. Learning in the presence of malicious errors. *SIAM J. Comput. 22*, 807–837.

KEARNS, M. J., AND SCHAPIRE, R. E. 1994. Efficient distribution-free learning of probabilistic concepts. *J. Comput. Syst. Sci. 48*, 3, 464–497.

KEARNS, M. J., SCHAPIRE, R. E., AND SELLIE, L. M. 1994. Toward efficient agnostic learning. *Mach. Learn. 17*, 2/3, 115–142.

LAIRD, P. D. 1988. Learning from good and bad data. In *Kluwer International Series in Engineering and Computer Science*. Kluwer Academic Publishers, Boston, Mass.

LITTLEWOOD, J. 1969. On the probability in the tail of a binomial distribution. *Adv. Appl. Prob. 1*, 43–72.

SAUER, N. 1972. On the density of families of sets. *J. Combin. Th. A 13*, 145–147.

SHELAH, S. 1972. A combinatorial problem; Stability and order for models and theories in infinitary languages. *Pacific J. Math. 41*, 247–261.

SIMON, H. U. 1996a. General bounds on the number of examples needed for learning probabilistic concepts. *J. Comput. Syst. Sci. 52*, 2, 239–255.

SIMON, H. U. 1996b. A general upper bound on the number of examples sufficient to learn in the presence of classification noise. Unpublished Manuscript.

VALIANT, L. G. 1984. A theory of the learnable. *Commun. ACM 27*, 11 (Nov.), 1134–1142.

VAPNIK, V. 1982. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, New York.