

Report Number: WUCS-92-25

1992-07-10

# Can PAC Learning Algorithms Tolerate Random Attribute Noise?

Authors: Sally A. Goldman and Robert H. Sloane

This paper studies the robustness of pac learning algorithms when the instances space is  $\{0,1\}^n$ , and the examples are corrupted by purely random noise affecting only the instances (and not the labels). In the past, conflicting results on this subject have been obtained -- the "best agreement" rule can only tolerate small amounts of noise, yet in some cases large amounts of noise can be tolerated.

We show that the truth lies somewhere in between these two alternatives. For uniform attribute noise, in which each attribute is flipped independently at random with the same probability, we present an algorithm that pac learns monomials for any (unknown) noise rate less than  $1/2$ . Contrasting this positive result, we show that product random attribute noise, where each attribute  $i$  is flipped randomly and independently with its own probability  $p_i$ , is nearly as harmful as malicious noise-- no algorithm can tolerate more than a very small amount of such noise.

... **Read complete abstract on page 2.**

Follow this and additional works at: [http://openscholarship.wustl.edu/cse\\_research](http://openscholarship.wustl.edu/cse_research)



Part of the [Computer Engineering Commons](#), and the [Computer Sciences Commons](#)

---

## Recommended Citation

Goldman, Sally A. and Sloane, Robert H., "Can PAC Learning Algorithms Tolerate Random Attribute Noise?" Report Number: WUCS-92-25 (1992). *All Computer Science and Engineering Research*.  
[http://openscholarship.wustl.edu/cse\\_research/588](http://openscholarship.wustl.edu/cse_research/588)

---

# Can PAC Learning Algorithms Tolerate Random Attribute Noise?

## **Complete Abstract:**

This paper studies the robustness of pac learning algorithms when the instances space is  $\{0,1\}^n$ , and the examples are corrupted by purely random noise affecting only the instances (and not the labels). In the past, conflicting results on this subject have been obtained -- the "best agreement" rule can only tolerate small amounts of noise, yet in some cases large amounts of noise can be tolerated.

We show that the truth lies somewhere in between these two alternatives. For uniform attribute noise, in which each attribute is flipped independently at random with the same probability, we present an algorithm that pac learns monomials for any (unknown) noise rate less than  $1/2$ . Contrasting this positive result, we show that product random attribute noise, where each attribute  $i$  is flipped randomly and independently with its own probability  $p_i$ , is nearly as harmful as malicious noise-- no algorithm can tolerate more than a very small amount of such noise.

**Can PAC Learning Algorithms Tolerate  
Random Attribute Noise?**

**Sally A. Goldman and Robert H. Sloan**

**WUCS-92-25**

**July 1992**

**Department of Computer Science  
Washington University  
Campus Box 1045  
One Brookings Drive  
Saint Louis, MO 63130-4899**



# Can PAC Learning Algorithms Tolerate Random Attribute Noise?

Sally A. Goldman\*  
Department of Computer Science  
Washington University  
St. Louis, Missouri 63130

Robert H. Sloan†  
Dept. of Electrical Engineering and Computer Science  
University of Illinois at Chicago  
Chicago, IL 60680

WUCS-92-25‡

July 10, 1992

## Abstract

This paper studies the robustness of pac learning algorithms when the instance space is  $\{0, 1\}^n$ , and the examples are corrupted by purely random noise affecting only the instances (and not the labels). In the past, conflicting results on this subject have been obtained—the “best agreement” rule can only tolerate small amounts of noise, yet in some cases large amounts of noise can be tolerated.

We show that the truth lies somewhere between these two alternatives. For *uniform attribute noise*, in which each attribute is flipped independently at random with the same probability, we present an algorithm that pac learns monomials for any (unknown) noise rate less than  $1/2$ . Contrasting this positive result, we show that *product random attribute noise*, where each attribute  $i$  is flipped randomly and independently with its own probability  $p_i$ , is nearly as harmful as malicious noise—no algorithm can tolerate more than a very small amount of such noise.

---

\*Supported in part by a GE Foundation Junior Faculty Grant and NSF grant CCR-9110108. Part of this research was conducted while the author was at the M.I.T. Laboratory for Computer Science and supported by NSF grant DCR-8607494 and a grant from the Siemens Corporation. Net address: sg@cs.wustl.edu.

†Supported in part by NSF grant CCR-9108753. Part of this research was conducted while the author was at Harvard and supported by ARO grant DAAL 03-86-K-0171. Net address: sloan@uicbert.eecs.uic.edu.

‡An earlier version of this report appeared as technical report WUCS-91-29.



# 1 Introduction

This paper studies the robustness of *pac* learning algorithms for learning boolean functions (i.e. the instance space is  $\{0, 1\}^n$ ). In particular, we examine random noise processes that corrupt examples by independently inverting each attribute bit with a probability given by the noise rate. We assume throughout that the classification of each instance is always correctly reported.

In the past, conflicting results on handling random attribute noise have been obtained. On the one hand, Sloan [12] has show that the “best-agreement” rule can only tolerate very small amounts of random attribute noise—suggesting that random attribute noise may be difficult to overcome. However, by using a different strategy, Shackelford and Volper [11] have obtained an algorithm that tolerates a large amount of random attribute noise (at a *known* noise rate) for learning  $k$ -DNF formulas. Thus their result suggests that random attribute noise may be like random classification noise, where large amounts of noise can sometimes be tolerated.

We show that the truth lies somewhere between these two alternatives. In order to fully understand the difficulty of overcoming random attribute noise, we first carefully examine the way in which the noise is formally modeled. Most previous work has considered *uniform random attribute noise*, in which each attribute is flipped independently at random with the same probability. The algorithm of Shackelford and Volper demonstrates that when the exact noise rate is known, at least for learning  $k$ -DNF formulas, a large amount of uniform random attribute noise can be tolerated. In this paper we extend that positive result by showing that even if the noise rate is *unknown*, the class of monomials is still efficiently learnable under a large amount of uniform random attribute noise.

Contrasting this positive result, we show that *product random attribute noise*, where each attribute  $i$  is flipped randomly and independently with its own probability  $p_i$  (all  $p_i$  are less than some given upper bound for the noise rate), is nearly as harmful as malicious noise. That is, no algorithm (regardless of sample complexity or computation time) can tolerate more than a very small amount of product random attribute noise. On the whole, these results are surprising. Intuitively, one would think that random labeling noise destroys much more information than random attribute noise, but, in fact, *pac* learning is possible with large amounts of random *labeling* noise [1].

## 2 Review of *pac* learning from noisy data

The method most commonly used for *pac* learning in the presence of noise is to pick a concept that has the best (or at least very good) agreement with a sample of data

corrupted by noise. It has been shown that for both discrete [1] and continuous [8] instance spaces, the hypothesis that minimizes disagreements meets the pac criterion when the examples are modified by random labeling noise. Sloan [12] has extended those results to the case of malicious labeling noise. Similarly, Kearns and Li [6] have shown that this method of minimizing disagreements can tolerate a small amount of malicious noise in discrete instance spaces. So these results specify the amount of noise that can be tolerated ignoring the issue of computation time. (Of course, if a hypothesis minimizing disagreements can be found in polynomial time then the above techniques produce efficient learning algorithms.) Finally, as Blumer et al. mention [2], their VC dimension methods can be used to prove that this minimal disagreement method also works for handling small amounts of malicious noise in continuous instance spaces. In the case of uniform random attribute noise, if one uses the minimal disagreement method, then the minimum error rate obtainable (i.e. the minimum “epsilon”) is bounded below by the noise rate [12]. We note that for arbitrary adversarial malicious noise, that is the maximum noise rate that *any* algorithm can tolerate [5].

Although the method of minimizing disagreements is not effective against random attribute noise, there are techniques for coping with uniform random attribute noise. In particular, Shackelford and Volper [11] have an algorithm that tolerates large amounts of random attribute noise for learning  $k$ -DNF formulas. That algorithm, however, has one very unpleasant requirement: it must be given the exact noise rate (or at least a very good estimate of the noise rate) as an input. Recently, Littlestone [9] has looked at how Winnow can tolerate several different models of attribute noise. He first considers an adversarial model of attribute noise. Assuming that the uncorrupted instances are linearly separable, under this noise model, the mistake bound for Winnow just increases by an additive term proportional to the weighted sum of the number of errors occurring in each relevant attribute. Next he considers the situation in which there is a large amount of redundant information provided by attributes that are separately indicative of the correct classification of each instance. In this case he shows that a large amount of noise can be handled where the irrelevant attributes are affected by arbitrary adversarial noise and the relevant attributes are affected by random noise independently of one another.

In this paper we present two results about handling attribute noise. In Section 4, we present an algorithm for learning monomials that tolerates large uniform random attribute noise, for any noise rate less than  $1/2$ , without any a priori knowledge of the noise rate. Contrasting this positive result, in Section 5 we show that with product random attribute noise, the minimum error rate obtainable is bounded below by one-half of the noise rate, regardless of the technique (or computation time) of the learning algorithm.



### 3 Notation

We assume that the reader is familiar with the model of pac learning introduced by Valiant [13]. Good discussions of the details of the model are given by Kearns et al. and by Haussler et al. [7, 3]. Briefly, a *concept* is a subset of some *instance space*  $X$ , and a concept class is some subset of  $2^X$ . An *example* of a concept  $c$  is a pair  $(x, s)$ , where  $x \in X$ , and  $s$  is 1 if  $x \in c$  and 0 otherwise. We call a sequence of examples a *sample*.

We assume a fixed (but unknown) probability distribution  $\mathbf{D}$  on  $X$ . Furthermore, the learner who is trying to learn concept  $c$  has available to it a black box or oracle called EX such that each call to EX returns a labeled instance,  $(x, s)$  where  $x$  is drawn at random from  $\mathbf{D}$  and labeled according to  $c$ . The learner's goal is the following: Given parameters  $0 < \epsilon, \delta \leq 1$ , draw a sample from EX, and output some representation of a concept  $\hat{c}$  such that

$$\Pr[\mathbf{D}(c \Delta \hat{c}) > \epsilon] \leq \delta,$$

where  $\Delta$  denotes symmetric difference, and the probability is over the calls to EX and any coin flips used by the learning algorithm. Such a  $\hat{c}$  is called  $\epsilon$ -good.

The ordinary definition of pac learning (from noiseless data) assumes that EX returns correct data. In this paper we are concerned with the case in which our instances come from some noise oracle, instead of the usual noise-free oracle, EX. Each noise oracle represents some noise process being applied to the examples from EX. The output from the noise process is all the learner can observe. The "desired," noiseless output of each oracle would thus be a correctly labeled example  $(x, s)$ , where  $x$  is drawn according to  $\mathbf{D}$ . We now describe the actual outputs from the following noise oracles:  $\text{MAL}_\nu$  [14],  $\text{URA}_\nu$  [12], and  $\text{PRA}_\nu$ .

- When  $\text{MAL}_\nu$  is called, with probability  $1 - \nu$ , it does indeed return a correctly labeled  $(x, s)$  where  $x$  is drawn according to  $\mathbf{D}$ . With probability  $\nu$  it returns an example  $(x, s)$  about which no assumptions whatsoever may be made. In particular, this example may be maliciously selected by an adversary who has infinite computing power, and has knowledge of the target concept,  $\mathbf{D}$ ,  $\nu$ , and the internal state of the algorithm calling this oracle. This *malicious noise* oracle models the situation where the learner usually gets a correct example, but some small fraction  $\nu$  of the time the learner gets noisy examples and the nature of the noise is unknown or unpredictable.
- The oracle  $\text{URA}_\nu$  makes sense only when the instance space is  $\{0, 1\}^n$  (i.e., we are learning boolean functions). The oracle  $\text{URA}_\nu$  calls EX and obtains some  $(x_1 \cdots x_n, s)$ .  $\text{URA}_\nu$  then adds noise to this example by independently flipping

each bit  $x_i$  to  $\bar{x}_i$  with probability  $\nu$  for  $1 \leq i \leq n$ . Note that the label of the “true” example is never altered by  $\text{URA}_\nu$ . This *uniform random attribute noise* oracle models a situation where the attributes of the examples are subject to noise, but that noise is as benign as possible. For example, the attributes might be sent over a noisy channel.

- The oracle  $\text{PRA}_\nu$  also only applies when we are learning boolean functions. This oracle calls EX and obtains some  $(x_1 \cdots x_n, s)$ . The oracle  $\text{PRA}_\nu$  then adds noise by independently flipping each bit  $x_i$  to  $\bar{x}_i$  with some fixed probability  $\nu_i \leq \nu$  for each  $1 \leq i \leq n$ . This *product random attribute noise* oracle provides a more realistic model of random attribute noise than  $\text{URA}_\nu$ .<sup>1</sup>

The noise oracles we will focus on here are  $\text{URA}_\nu$  and  $\text{PRA}_\nu$ .

## 4 Learning monomials from noisy data

In this section we present an algorithm for learning monomials from data corrupted with uniform random attribute noise with any noise rate less than  $1/2$ . The examples will come from  $\text{URA}_\nu$ , and the learner will receive *no* prior information about  $\nu$ .

The key idea we exploit is the following: Imagine that the literal  $x_1$  is included in the target concept. Then whenever the learner receives a positive instance with the first bit off, it must be that the bit was on in the “noise free” instance, and flipped by the noise oracle. Hence, the ratio

$$\frac{\text{number of times bit 1 is off in positive instance}}{\text{total number of positive instances}} \tag{1}$$

provides a good estimate of the noise rate. Notice also that if  $x_1$  is *not* in the formula, then the ratio specified in (1) still is bounded below by the noise rate. Thus we will estimate the noise rate to be the minimum, over all literals, of the ratio specified in (1) for the literal  $x_1$ .

Since we are able to obtain good estimates for the noise rate, we can apply Angluin and Laird’s [1] technique of successive approximation to obtain an upper bound for the noise rate that is sufficiently close to the actual noise rate. Finally, using this estimate of the noise rate, we apply Shackelford and Volper’s [11] algorithm when specialized to the case of monomials. However, the correctness proof provided by Shackelford and Volper assumes the *exact* noise rate is provided. While one could extend their proof to apply when only given a sufficiently close estimate of the noise

---

<sup>1</sup>Technically,  $\text{PRA}_\nu$  specifies a family of oracles, each member of which is specified by  $n$  variables,  $\nu_1, \dots, \nu_n$ , where  $0 \leq \nu_i \leq \nu$ .

rate, we give a much more direct proof of correctness that is easy to follow and yields a better bound on the time and sample complexity.

Let the literals be numbered from 1 to  $2n$ , and for each literal  $i$ , let

$$\begin{aligned} q_i &= \Pr[\text{Literal } i \text{ is off in a positive instance from EX (noiseless data)}] \\ p_i &= \Pr[\text{Literal } i \text{ is off in a positive instance from } \text{URA}_\nu]. \end{aligned}$$

Our goal is to output a conjunction that contains every literal that is in the target monomial, and no literal with a high value of  $q_i$ . Of course, we cannot directly estimate the  $q_i$ 's since we only see examples from  $\text{URA}_\nu$ . Our method is to accurately estimate all the  $p_i$ 's using examples from  $\text{URA}_\nu$  and then use these estimates to determine which literals have a high value for  $q_i$  and should thus be excluded. Observe that for all  $i$

$$\begin{aligned} p_i &= q_i(1 - \nu) + (1 - q_i)\nu \\ &= \nu + q_i(1 - 2\nu) \end{aligned} \tag{2}$$

Thus for any literal  $i$  that is in the target monomial,  $p_i = \nu$ . Furthermore, since  $\nu < 1/2$ , for any literal  $i$  that is not in the target monomial  $p_i \geq \nu$ . We will show that by accurately estimating all the  $p_i$ 's we can obtain a good estimate for the noise rate by simply taking the minimum estimated value over all the  $p_i$ 's. We then output the conjunction of all literals  $i$  having values of  $p_i$  close to that minimum. The algorithm is specified in full in Figure 1.

**Theorem 1** *The algorithm specified in Figure 1 pac learns any monomial given data from  $\text{URA}_\nu$ . The sample complexity is  $O\left(\frac{n^2}{(1-2\nu)^2\epsilon^3} \ln \frac{n}{\delta} + \frac{1}{(1-2\nu)^2} \ln \frac{n}{\delta(1-2\nu)}\right)$ , and the time complexity is  $O\left(\frac{n^3}{(1-2\nu)^2\epsilon^3} \ln \frac{n}{\delta} + \frac{n}{(1-2\nu)^2} \ln \frac{n}{\delta(1-2\nu)}\right)$ .*

**Proof:** It follows from Theorem 3 of Angluin and Laird [1] that after step 2 of the algorithm with probability at least  $1 - \delta/2$ , the algorithm halts with  $r \leq 1 + \lceil \lg \frac{1}{1-2\nu} \rceil$ ,  $\nu \leq \nu_b < 1/2$ , and  $\frac{1}{1-2\nu_b} \leq \frac{2}{1-2\nu}$ . Thus the total sample complexity needed for step 2 is at most

$$\sum_{r=1}^{1 + \lceil \lg \frac{1}{1-2\nu} \rceil} 2^{2r+3} \ln \left( \frac{n 2^{r+3}}{\delta} \right) = O \left( \left( \frac{1}{1-2\nu} \right)^2 \ln \frac{n}{\delta(1-2\nu)} \right).$$

Let  $p_+$  denote the probability of drawing a positive example from  $\text{URA}_\nu$ . (Note that since the noise process does not affect the labels,  $p_+$  is also the probability of drawing a positive example directly from EX.) Applying Hoeffding's Inequality [4] it is easily shown that if  $p_+ \geq \epsilon$  then using a sample of size  $\max \left\{ \frac{2m}{\epsilon}, \frac{2}{\epsilon^2} \ln \frac{4}{\delta} \right\}$  ensures with probability at least  $1 - \delta/4$  that the algorithm will obtain at least  $m$  positive examples in step 4. Of course, if  $p_+ < \epsilon$  then the hypothesis FALSE is  $\epsilon$ -good.

Inputs:  $n, \epsilon, \delta$ , access to  $\text{URA}_\nu$  ( $\nu$  unknown).

Output: Some monomial (possibly “FALSE”).

1. Let  $\hat{\nu}_b = 1/4$ ,  $r = 1$ ,  $\gamma = (\frac{1}{2})^{r+2}$   
 (Comment:  $\hat{\nu}_b$  is current guess at an upper bound on the unknown  $\nu$ .)
2. Repeat until done
  - (a) Draw  $m_r = 2^{2r+3} \ln \frac{n2^{r+3}}{\delta}$  examples from  $\text{URA}_\nu$
  - (b) For each literal  $i$   
 $\hat{p}_i = (\text{number of times literal } i \text{ is off in sample})/m_r$
  - (c)  $\hat{p}_m = \min_i \hat{p}_i$
  - (d) If  $\hat{p}_m < \hat{\nu}_b - \gamma$ , let  $\nu_b = \hat{\nu}_b$  and exit repeat loop
  - (e) Else  $r = r + 1$ ,  $\hat{\nu}_b = \frac{1}{2} - (\frac{1}{2})^{r+1}$ ,  $\gamma = \gamma/2$
3. Let  $m = \frac{32n^2}{(1 - 2\nu_b)^2 \epsilon^2} \ln \frac{16n}{\delta}$
4. Draw  $m' = \max \left\{ \frac{2m}{\epsilon}, \frac{2}{\epsilon^2} \ln \frac{4}{\delta} \right\}$  examples from  $\text{URA}_\nu$ 
  - (a) If  $m$  positive instances are not obtained, halt and output “FALSE”
  - (b) Else let  $S$  be a sample of  $m$  positive examples
5. For each literal  $i$   
 $\hat{p}_i = (\text{number of times literal } i \text{ is off in } S)/m$ .
6.  $\hat{\nu} = \min_i \hat{p}_i$ .
7. Output conjunction of all literals  $i$  such that  $\hat{p}_i \leq \hat{\nu} + \frac{\epsilon(1 - 2\nu_b)}{4n}$

**Figure 1:** Algorithm for pac learning monomials under random attribute noise with an unknown noise rate.

We now show that if at least  $m$  positive examples are obtained in step 4, then the hypothesis output in step 7 is  $\epsilon$ -good with probability at least  $1 - \delta$ . Again, by applying Hoeffding's Inequality [4], it is easily shown that by using a sample of size

$$m = \frac{32n^2}{(1 - 2\nu_b)^2 \epsilon^2} \ln \frac{16n}{\delta} = O\left(\frac{n^2}{(1 - 2\nu)^2 \epsilon^2} \ln \frac{n}{\delta}\right),$$

the probability that *all* of the estimates  $\hat{p}_i$  are within  $\epsilon(1 - 2\nu_b)/8n$  of their true value  $p_i$  is at least  $1 - \delta/4$ . That is,

$$\Pr\left[\bigwedge_{(1 \leq i \leq 2n)} p_i - \frac{\epsilon(1 - 2\nu_b)}{8n} < \hat{p}_i < p_i + \frac{\epsilon(1 - 2\nu_b)}{8n}\right] \geq 1 - \delta/4.$$

We now assume that  $\nu \leq \nu_b$ , at least  $m$  positive examples are obtained in step 4, and all the  $\hat{p}_i$ 's are within the above tolerance (these conditions are all satisfied with probability at least  $1 - \delta$ ). Then:

1. The estimate  $\hat{\nu}$  of  $\nu$  is accurate. Namely,

$$\nu - \frac{\epsilon(1 - 2\nu_b)}{8n} < \hat{\nu} < \nu + \frac{\epsilon(1 - 2\nu_b)}{8n}.$$

2. Any literal that is in the target monomial will be placed in the algorithm's hypothesis.
3. For any literal  $i$  that is *not* in the target monomial but is placed in the algorithm's hypothesis  $q_i < \epsilon/2n$ .

Since for all  $i$ ,  $p_i \geq \nu$  and for any literal in the target monomial  $p_i = \nu$ , item 1 above easily follows from the fact that for all  $i$ ,  $\hat{p}_i$  is within  $\epsilon(1 - 2\nu_b)/8n$  of  $p_i$ . For any literal  $i$  that is in the target monomial,  $p_i = \nu$ , and thus

$$\hat{p}_i < \hat{\nu} + \frac{\epsilon(1 - 2\nu_b)}{4n}.$$

Hence for every literal  $i$  that is in the target monomial,  $\hat{p}_i$  will satisfy the inequality in step 7 of the algorithm and thus be placed in the hypothesis.

Finally, to prove that item 3 above holds, we show that if  $q_i \geq \epsilon/2n$  then  $\hat{p}_i$  will not satisfy the inequality in step 7. Applying equation (2) we get

$$p_i \geq \nu + \frac{\epsilon(1 - 2\nu_b)}{2n}.$$

Finally since  $\hat{p}_i$  and  $\hat{\nu}$  are within the given tolerance of their true values, it follows that:

$$\hat{p}_i > \hat{\nu} + \frac{\epsilon(1 - 2\nu_b)}{4n}.$$

Thus the choice of literals made by the algorithm in step 7 ensures that every literal  $i$  in the output formula has  $q_i < \epsilon/2n$ .

To complete the proof, observe that literals in the output monomial are a superset of the literals in the target monomial. Therefore, the algorithm's hypothesis is false whenever the target concept is false. Since every literal  $i$  in the output has  $q_i < \epsilon/2n$  and there are at most  $2n$  literals in the output formula, the probability that the output is false when the target concept is true is at most  $\epsilon$ .

□

*Remark:* We can obtain a very similar result for a noise model where with probability  $1 - \nu$  the example is noise free, and with probability  $\nu$  a single one of the  $n$  bits is picked at random and flipped.

## 5 Product random attribute noise

In this section we show that product random attribute noise makes pac learning almost as difficult as malicious noise. Kearns and Li [6] showed that for any nontrivial concept class, it is impossible to pac learn to accuracy  $\epsilon$  with examples from  $\text{MAL}_\nu$  unless

$$\nu < \epsilon/(1 + \epsilon).$$

Our result for product random attribute noise is similar, with a slightly weaker bound.

**Theorem 2** *Let  $\mathcal{C}$  be any concept class over the domain  $\{0,1\}^n$  that includes the concepts  $x_i$  and  $x_j$  for some  $i \neq j$ . It is possible to pac learn  $\mathcal{C}$  to accuracy  $\epsilon$  with examples from  $\text{PRA}_\nu$  only if  $\nu < 2\epsilon$ .*

**Proof:** We use the method of induced distributions [6].

Say  $\mathcal{C}$  contains the concepts  $x_1$  and  $x_2$ . In what follows, we will put zero probability weight on instances containing 1's in positions 3 through  $n$ , and thus may assume without loss of generality that  $\nu_k = 0$  for  $3 \leq k \leq n$ . In fact, we can assume that our entire instance space is  $\{00, 01, 10, 00\}$ , since all instances seen will have 0's in all other attributes.

Fix some value of  $\nu$  in the range  $0 \leq \nu < 1/2$ . Consider a distribution  $\mathbf{D}$  which assigns weight  $(1 - \nu)/2$  to 00 and 11 and weight  $\nu/2$  to 01 and 10. In Table 1 we show the two noise-free probability distributions on examples obtained by labeling the instances drawn from  $\mathbf{D}$  according to concept  $x_1$  or concept  $x_2$ .

Now consider what happens under the following two learning problems:

1. For the first learning problem let  $x_1$  be the target concept, let  $\mathbf{D}$  be the distribution on instances (so  $\mathbf{D}_1$  is the noise-free distribution on examples), and let

Labeled Instance	$\mathbf{D}_1$	$\mathbf{D}_2$	Observed dist.
(00, -)	$(1 - \nu)/2$	$(1 - \nu)/2$	$(1 - \nu)^2/2$
(00, +)	0	0	$\nu^2/2$
(01, -)	$\nu/2$	0	$\nu(1 - \nu)/2$
(01, +)	0	$\nu/2$	$\nu(1 - \nu)/2$
(10, -)	0	$\nu/2$	$\nu(1 - \nu)/2$
(10, +)	$\nu/2$	0	$\nu(1 - \nu)/2$
(11, -)	0	0	$\nu^2/2$
(11, +)	$(1 - \nu)/2$	$(1 - \nu)/2$	$(1 - \nu)^2/2$

**Table 1:** Two induced noise-free probability distributions obtained from distribution  $\mathbf{D}$  for concept class  $\mathcal{C}$ , and the distribution obtained after noise for the two learning problems used in the proof of Theorem 2.

the noise oracle be  $\text{PRA}_\nu$  with  $\nu_1 = \nu$  and  $\nu_2 = 0$ . The observed distribution on examples is shown in the last column of Table 1.

- For the second learning problem, let  $x_2$  be the target concept, let  $\mathbf{D}$  be the distribution (so  $\mathbf{D}_2$  is the noise-free distribution on examples), and let the noise oracle be  $\text{PRA}_\nu$  with  $\nu_1 = 0$  and  $\nu_2 = \nu$ . The observed distribution on examples is also the one shown in the last column of Table 1.

These two learning problems have an identical probability distribution on the observed (noisy) samples. Therefore, no pac learning algorithm has any basis for distinguishing between these two scenarios. Thus with probability at least  $1 - \delta$  the learning algorithm must output a concept  $c$  such that:

$$\mathbf{D}(c_1 \triangle c) < \epsilon \text{ and } \mathbf{D}(c_2 \triangle c) < \epsilon.$$

By the triangle inequality it follows that:

$$\mathbf{D}(c_1 \triangle c) + \mathbf{D}(c_2 \triangle c) \geq \mathbf{D}(c_1 \triangle c_2).$$

Thus with probability at least  $1 - \delta$  the learning algorithm must a concept  $c$  such that  $\mathbf{D}(c_1 \triangle c_2) < 2\epsilon$ . Finally, since  $\mathbf{D}(c_1 \triangle c_2) = \nu$ , no learning algorithm can succeed unless  $\nu < 2\epsilon$ .  $\square$

*Remark:* Notice that the lower bound of Theorem 2 is an information-theoretic representation-independent hardness result. No algorithm, regardless of either its sample size or computation time, can escape this bound. Furthermore, we have made no assumptions on the representation class from which the hypothesis may be

selected—this bound on the tolerable noise rate holds for *any* hypothesis the learner may output.

## 6 Final thoughts

We have studied the robustness of pac learning algorithms under several forms of random attribute noise. We presented a new algorithm for learning monomials that tolerates large amounts of uniform random attribute noise (any noise rate less than  $1/2$ ) without any prior knowledge of the noise rate. An intriguing open question is whether one can pac learn  $k$ -DNF formulas under uniform random attribute noise for an *unknown* noise rate.

However, we feel that our negative result for the more realistic product random attribute noise oracle makes it clear that, in general, under the pac learning model random attribute noise is quite harmful. This result was surprising to the authors. One expects to only be able to tolerate a small amount of truly malicious noise—it is obviously the worst sort of noise possible. Yet, one would expect that labeling noise would be worse than random attribute noise. Indeed, in one empirical test (of the ID-3 system), that is exactly what was found [10]. Yet, in spite of both these empirical results and our intuition, we have shown that in the pac model random attribute noise (when it is product) is significantly more harmful than random labeling noise. In fact, product random attribute noise is significantly more harmful than *malicious* labeling noise generated by a powerful adversary [12], and nearly as harmful as truly malicious noise.

## Acknowledgments

We would like to thank Les Valiant for suggesting this line of research. We thank Dana Angluin for pointing out an error (and providing a correction to this error) in our original statement of Theorem 2. We also thank the anonymous reviewers for their many useful comments. In particular, we thank them for pointing out that we could combine our techniques with those from Angluin and Laird to obtain a good estimate for the noise rate.

## References

- [1] Dana Angluin and Philip Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988.



- [2] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, 1989.
- [3] David Haussler, Michael Kearns, Nick Littlestone, and Manfred K. Warmuth. Equivalence of models for polynomial learnability. *Information and Computation*. To appear.
- [4] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, March 1963.
- [5] Michael Kearns. Thoughts on hypothesis boosting. (Unpublished), December 1988.
- [6] Michael Kearns and Ming Li. Learning in the presence of malicious errors. In *Proceedings of the Twentieth Annual ACM Symposium on Theory of Computing*, Chicago, Illinois, May 1988.
- [7] Michael Kearns, Ming Li, Leonard Pitt, and Leslie Valiant. On the learnability of boolean formulae. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, pages 285–295, New York, New York, May 1987.
- [8] Philip D. Laird. *Learning from Good and Bad Data*. Kluwer international series in engineering and computer science. Kluwer Academic Publishers, Boston, 1988.
- [9] Nicholas Littlestone. Redundant noisy attributes, attribute errors, and linearity-threshold learning using winnow. In *Fourth Workshop on Computational Learning Theory*, pages 147–156, 1991.
- [10] J. Ross Quinlan. The effect of noise on concept learning. In *Machine Learning, An Artificial Intelligence Approach (Volume II)*, chapter 6, pages 149–166. Morgan Kaufmann, 1986.
- [11] George Shackelford and Dennis Volper. Learning  $k$ -DNF with noise in the attributes. In *First Workshop on Computational Learning Theory*, pages 97–103, Cambridge, Mass. August 1988. Morgan Kaufmann.
- [12] Robert H. Sloan. Types of noise in data for concept learning. In *First Workshop on Computational Learning Theory*, pages 91–96. Morgan Kaufmann, 1988.
- [13] Leslie G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, November 1984.
- [14] Leslie G. Valiant. Learning disjunctions of conjunctions. In *Proceedings IJCAI-85*, pages 560–566. International Joint Committee for Artificial Intelligence, Morgan Kaufmann, August 1985.