

A Study of Maximal-Coverage Learning Algorithms

Hussein Almuallim

Thomas G. Dietterich

Department of Computer Science

Oregon State University

Corvallis, OR 97331

almualh@cs.orst.edu

tgd@cs.orst.edu

March 15, 1994

Abstract

The *coverage* of a learning algorithm is the number of concepts that can be learned by that algorithm from samples of a given size. This paper asks whether good learning algorithms can be designed by maximizing their coverage. The paper extends a previous upper bound on the coverage of any Boolean concept learning algorithm and describes two algorithms—Multi-Balls and Large-Ball—whose coverage approaches this upper bound. Experimental measurement of the coverage of the ID3 and FRINGE algorithms shows that their coverage is far below this bound. Further analysis of Large-Ball shows that although it learns many concepts, these do not seem to be very interesting concepts. Hence, coverage maximization alone does not appear to yield practically-useful learning algorithms. The paper concludes with a definition of coverage *within a bias*, which suggests a way that coverage maximization could be applied to strengthen weak preference biases.

1 Introduction

Research in computational learning theory (e.g., [Valiant 84], [Natarajan 87], [COLT 88]–[COLT 91]) has provided many insights into the capabilities and limitations of inductive learning from examples. However, an important shortcoming of most work in this area is that it focuses on learning concepts drawn from *prespecified classes* of concepts (e.g., linearly separable functions, k -DNF formulae). This style of research begins by choosing a restricted class of concepts and then finding a polynomial bound—called the *sample complexity*—such

that if a sample of size larger than the sample complexity is available, any concept from the concept class that is consistent with the sample will be approximately correct with high probability.

Work of the above type usually leads to a learning algorithm that is *specialized* in learning the prescribed class of concepts, and an upper bound on the number of training examples required by the algorithm to guarantee successful learning. For real-world applications, such findings can be viewed as follows: A learning algorithm L designed to learn a class of concepts \mathcal{C} is guaranteed to succeed¹ in application domains in which the target concept belongs to \mathcal{C} (i.e., the restrictions used to define \mathcal{C} are satisfied by the target concept), provided that a sufficient number of training examples is given to the algorithm. Of course, no such guarantees are given if the target concept is not in \mathcal{C} .

This naturally means that one should seek algorithms that learn concept classes that are as large (i.e., less restricted) as possible. Rivest [Rivest 87], for instance, mentions this goal most explicitly by saying:

“One goal of research in machine learning is to identify the largest possible class of concepts that are learnable from examples.”

This goal is also declared (although less explicitly) in many papers in the related literature (e.g., [Valiant 84], [Natarajan 87], [COLT 88]–[COLT 91]).

Nevertheless, it is a well-known fact that learning larger classes of concepts necessarily requires a larger number of training examples [Blumer et.al. 87]. Such trade-off between the size of the class of concepts being learned and the required number of training examples dictates how far one can go in attempting to learn larger and larger classes of concepts.

Traditionally, this issue has been addressed by identifying new classes of concepts that are as large as possible but still require a training sample of size bounded by some polynomial. Such an approach, however, does not enjoy great practical merit. In fact, the idea of learning prescribed classes of concepts in general suffers from two important problems:

- Training examples are usually hard to obtain. In a typical inductive learning task, one has only a limited number of training examples, much less than the polynomial bounds provided by learning theory.
- The concept class is usually unknown. In most application settings, there is often considerable flexibility (and concomitant lack of prior knowledge) concerning the choice of which concept class to explore. In fact, many of the concept classes studied in computational learning theory have never been supported by any practical justification.

Due to these difficulties, the learning algorithms and sample complexity bounds developed in computational learning theory have rarely been of practical value.

Recently, an alternative theoretical framework was introduced [Dietterich 89]. Instead of fixing a class of concepts and then deriving the sample complexity, this framework turns the problem around by asking: Given a *fixed number* of training examples, what is the largest

¹The guarantees are on being approximately correct with high confidence.

collection of concepts that some algorithm can learn? The intuition behind this framework is that, in the absence of additional information, one should prefer the learning algorithm that has the highest chance of learning the unknown concept—that is, the algorithm that can learn the largest number of concepts. In short, this framework could provide an approach to discovering an “optimal” bias for inductive learning in the absence of prior knowledge.

The goal of this paper is to explore this approach. We define the *coverage* of a learning algorithm to be the number of concepts learnable by the algorithm from a given sample size (and other relevant parameters). There are three questions raised by this approach:

1. For given sample size m , accuracy parameter ϵ and confidence parameter δ , what is the largest possible coverage that any algorithm can achieve?
2. Can we design a learning algorithm that attains this optimal coverage?
3. What is the coverage of existing learning algorithms?

This paper contributes to answering each of these questions. First, we generalize the upper bound on coverage given in [Dietterich 89]. Next, we present two learning algorithms and determine their coverage analytically. The coverage of the first algorithm, Multi-Balls, is shown to be quite close to the upper bound. The coverage of the second algorithm, Large-Ball, turns out to be even better than Multi-Balls in many situations. Third, we considerably improve upon Dietterich’s limited experiments for estimating the coverage of existing learning algorithms. We find that the coverage of Large-Ball exceeds the coverage of ID3 [Quinlan 86] and FRINGE [Pagallo and Haussler 90] by more than an order of magnitude in most cases.

These results are very thought-provoking, because, upon careful analysis, it becomes clear that the Large-Ball algorithm is rather trivial and uninteresting. In the final part of the paper, we conclude that coverage analysis does not—by itself—provide a framework for deriving an optimal inductive bias. It does, however, provide a framework for designing optimal-coverage algorithms within a given bias.

2 Definitions and Notation

We consider the space of Boolean concepts defined on n Boolean features. Let U_n be the set of all the 2^n truth assignments to the n features. A *concept* is an arbitrary set $c \subseteq U_n$. An *example* of a concept c is a pair $\langle X, c(X) \rangle$ where $c(X) = 1$ if $X \in c$ and 0 otherwise. The example is called *positive* in the first case, and *negative* in the second.

We assume the uniform distribution over U_n . However, all our results can be easily extended to the distributions where the probability is 0 on a subset of U_n and uniform on the rest. This is done by substituting the number of instances in U_n having non-zero probability in place of every occurrence of 2^n in the results.

A *training sample* of a concept c is a collection of examples drawn randomly from U_n and labeled according to c . The number of examples in this collection is called the *sample size*, denoted by m . Except in our experimental work, we assume that examples in a sample

are drawn independently (i.e., with replacement), and thus, a sample of size m does not necessarily contain m distinct examples. It should be noted that, assuming $m \ll 2^n$, the difference between sampling with and without replacement is not significant.

The *disagreement* between a training sample and a concept is the number of examples in the sample that are incorrectly classified by the concept.

The *distance* between two concepts c and h is the number of assignments $X \in U_n$ such that $c(X) \neq h(X)$. The *error* between c and h is the distance divided by 2^n , which is equivalent to the probability that a randomly chosen X will be classified differently by the two concepts. For any $0 < \epsilon < 1$, we say that h is ϵ -close to c if the error between the two concepts is at most ϵ . We let $Ball(c, \epsilon)$ denote the set of concepts that are ϵ -close to c . Note that for any concept $c' \in Ball(c, \epsilon)$, the distance between c' and c is at most $\epsilon 2^n$. Therefore, the number of concepts in $Ball(c, \epsilon)$ is given by $|Ball(c, \epsilon)| = \sum_{i=0}^{\lfloor \epsilon 2^n \rfloor} \binom{2^n}{i}$. We call c and $\lfloor \epsilon 2^n \rfloor$ the *center* and *radius* of the ball, respectively.

A *learning algorithm* is a mapping from the space of samples to the space of concepts. The output of the algorithm is called an *hypothesis*. An hypothesis is *consistent* if it has no disagreement with the training sample.

We adopt PAC learning [Blumer et.al. 87] as the criterion for successful learning, but we restrict this to learning under the uniform distribution only. We say that an algorithm L *learns* a concept c for given m, ϵ and δ , if with probability at least $1 - \delta$, L returns some hypothesis h that is ϵ -close to c when given a randomly drawn sample of c of size m . Formally, let S_c^m denote a random sample of c of size m . We say that L learns c with respect to m, ϵ and δ if

$$Pr[error(h, c) \leq \epsilon] \geq 1 - \delta$$

where h is the hypothesis returned by L given S_c^m , and where the probability is computed over all the samples of c of size m . ϵ and δ are called the *accuracy* and *confidence* parameters, respectively.

In general, ϵ and δ are in the range $0 < \epsilon, \delta < 1$. In practice, however, only values that are close to 0 are interesting. For this reason, we will sometimes explicitly assume for instance that $0 < \epsilon < \frac{1}{4}$ and/or $0 < \delta < \frac{1}{2}$, with the understanding that these are reasonable assumptions in practice. Further, to simplify our results, we will only consider the values of ϵ such that $\epsilon 2^n$ is a positive integer. Clearly, this is not a serious assumption when n is sufficiently large.

For given n, m, ϵ and δ , the *coverage* of a learning algorithm is the number of concepts the algorithm learns with respect to these parameters.

3 Upper Bound on Coverage

We begin by proving an upper bound on the best coverage that any algorithm can attain. An upper bound of this type has been proven for the case where the training sample is drawn randomly *without* replacement [Dietterich 89]. In the following, we generalize Dietterich's

result and show that the same upper bound also holds for the case where sampling is done with replacement. In addition, we provide a closed-form expression for this bound.

Theorem 1 *Assuming that $m \leq (1 - 2\epsilon)2^n$, the coverage of any learning algorithm under the uniform distribution can not exceed*

$$\frac{1}{1 - \delta} 2^m \sum_{i=0}^{\epsilon 2^n} \binom{2^n - m}{i}$$

concepts, for sample size m , accuracy parameter ϵ and confidence parameter δ .

Proof: The proof of this theorem uses Lemmas 6 and 7 which are given in the appendix. Let L be any learning algorithm and let $C \subseteq 2^{U_n}$ be the set of concepts learned by L for the given learning parameters. Let us denote by p_c the probability that a sample of size m for a concept c be mapped by L to some hypothesis within ϵ of c . By definition, for any $c \in C$, p_c must be at least $1 - \delta$. Thus, it must be true that

$$\sum_{c \in C} p_c \geq |C|(1 - \delta)$$

and therefore

$$|C| \leq \frac{1}{1 - \delta} \sum_{c \in C} p_c \leq \frac{1}{1 - \delta} \sum_{c \in 2^{U_n}} p_c$$

since $C \subseteq 2^{U_n}$. In the following, we let

$$W = \sum_{c \in 2^{U_n}} p_c .$$

The theorem holds by proving an upper bound on W as follows.

Let \mathcal{S} be the set of all possible outcomes of randomly drawing m objects from U_n . For any $s \in \mathcal{S}$, let $\sharp s$ denote the number of distinct objects in s . Note that $\sharp s = m$ for every $s \in \mathcal{S}$ in the case of sampling without replacement, and that $1 \leq \sharp s \leq m$ in the case of sampling with replacement. Also, let $Pr[s]$ denote the probability of the outcome s .

Now, for some $s \in \mathcal{S}$, suppose t is a training sample obtained by arbitrarily labeling the objects in s as positive or negative, and let h be the hypothesis returned by L when given t . Mapping t to h contributes the amount of $Pr[s]$ to p_c for every concept c that is (i) consistent with t and (ii) within ϵ of h . Therefore, we can compute W by summing up this contribution for all possible outcomes $s \in \mathcal{S}$ and all possible ways of labeling these. That is,

$$W = \sum_{s \in \mathcal{S}} Pr[s] \sum_{t \in Labeling(s)} N(t)$$

where $Labeling(s)$ is the set of all training samples obtained by labeling s , and $N(t)$ is the number of concepts that are consistent with t and within ϵ from the hypothesis returned by L when given t . By Lemma 6, for any $t \in Labeling(s)$

$$N(t) \leq \sum_{i=0}^{\epsilon 2^n} \binom{2^n - \sharp s}{i}$$

which means that

$$\begin{aligned}
W &\leq \sum_{s \in \mathcal{S}} \left\{ Pr[s] \times |Labeling(s)| \times \sum_{i=0}^{\epsilon 2^n} \binom{2^n - \#s}{i} \right\} \\
&= \sum_{s \in \mathcal{S}} \left\{ Pr[s] \times 2^{\#s} \times \sum_{i=0}^{\epsilon 2^n} \binom{2^n - \#s}{i} \right\} \\
&= \sum_{s \in \mathcal{S}} \left\{ Pr[s] \times \sum_{i=0}^{\epsilon 2^n} 2^{\#s} \times \binom{2^n - \#s}{i} \right\}.
\end{aligned}$$

For both cases of sampling with and without replacement, $\#s$ is at most m . Therefore, using Lemma 7 we can write

$$\begin{aligned}
W &\leq \sum_{s \in \mathcal{S}} Pr[s] \sum_{i=0}^{\epsilon 2^n} 2^m \binom{2^n - m}{i} \\
&= 2^m \sum_{i=0}^{\epsilon 2^n} \binom{2^n - m}{i} \sum_{s \in \mathcal{S}} Pr[s] \\
&= 2^m \sum_{i=0}^{\epsilon 2^n} \binom{2^n - m}{i}
\end{aligned}$$

which proves the theorem. \square

The following theorem puts the upper bound of Theorem 1 in closed-form.

Theorem 2 For $0 < \epsilon < \frac{1}{4}$ and $m < \frac{1}{4}2^n$, the quantity of Theorem 1 is further bounded above by

$$\frac{2^{(1-\epsilon \log_2 e)m+1} \binom{2^n}{\epsilon 2^n}}{1-\delta} \leq \frac{2^{(1-1.44\epsilon)m+1+H(\epsilon)2^n}}{1-\delta}$$

where e is the base of the natural logarithm, and $H(\epsilon) = \epsilon \log_2 \frac{1}{\epsilon} + (1-\epsilon) \log_2 \frac{1}{1-\epsilon}$.

Proof: The first inequality follows directly from Lemma 9. The second inequality follows from Lemma 5. \square

This result shows that given a training sample of a reasonable size, any learning algorithm can learn only a small proportion of the concept space. As a numerical example, consider the case where $n = 20, m = 100,000$ and $\delta = \epsilon = 0.05$. Note that while having 20 features in a practical domain is not unusual, 100,000 examples is a rather large sample size. In this case, $H(\epsilon) \approx 0.286$. The above result states that no learning algorithm can learn more than

$$\frac{2^{92,788+0.286 \times 2^{20}}}{0.95}$$

concepts. This is less than $\frac{1}{600,000}$ of the $2^{2^{20}}$ possible concepts definable over 20 features—a strikingly small fraction.

For the extreme case where $\delta = 0$, we can derive a much tighter bound:

Theorem 3 *If $\delta = 0$ and $\epsilon < \frac{1}{4}$, then the coverage of any learning algorithm is at most*

$$\sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$$

concepts, for accuracy parameter ϵ and confidence parameter δ .

Proof: Let L be any learning algorithm. Any two concepts c_1, c_2 , such that $c_1 \neq \neg c_2$, must share some samples in common. Since $\delta = 0$, for L to learn both concepts, L must map *all* of these samples to some hypothesis h that is within ϵ of both concepts. Such an h exists only if $\text{distance}(c_1, c_2) \leq 2\epsilon 2^n$.

Now, suppose that C is the set of concepts learned by L for $\delta = 0$, $\epsilon < \frac{1}{4}$ and some particular m . One of the following two cases must hold:

- There exists no concept c such that both c and $\neg c$ are in C . This implies that the distance between every pair of concepts in C is at most $2\epsilon 2^n$, and hence, $|C| \leq \sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$ as desired.
- There exists some concept c such that $c, \neg c \in C$. Let c' be any concept such that $c' \neq c$ and $c' \neq \neg c$. c' must, therefore, share some samples with c and some samples with $\neg c$. Thus, for c' to be learned by L , it must be within $2\epsilon 2^n$ of both c and $\neg c$. However, this is impossible since $\epsilon < \frac{1}{4}$. This implies that C contains no concepts other than c and $\neg c$ which means a coverage of only 2.

This proves the theorem. \square

This suggests that Theorem 1 is not tight when δ is very small. It also implies that the degree of freedom provided by the confidence parameter, δ , in the PAC definition is very important. Any algorithm that does not exploit this freedom to output (with probability within δ) a totally incorrect hypothesis can have only very limited coverage.

4 The Multi-Balls Learning Algorithm

Given these upper bounds on coverage, can we design algorithms that achieve these coverages?

Let c_1 and c_2 be two concepts with distance d , and suppose that we desire to construct a learning algorithm L that learns both c_1 and c_2 . To do this, we must consider how L should treat every possible training sample consistent with c_1 , c_2 , or both.

Obviously, any sample that is consistent with only one of c_1 or c_2 can be mapped to some hypothesis that is within ϵ of the consistent concept. The key question is how to map a sample that is consistent with both c_1 and c_2 . Now, if $\frac{d}{2^n} \leq 2\epsilon$, then there exists some concept h that is within ϵ of both c_1 and c_2 . In this case, all we need to do is to map the sample to h . However, if $\frac{d}{2^n} > 2\epsilon$, then there exists no concept that is within ϵ of both c_1 and c_2 , and therefore, we must map the sample either in favor of c_1 or in favor of c_2 , but not

Algorithm Two-Balls (*Sample*)

1. If $\text{disagreement}(\text{Sample}, c_1) < \text{disagreement}(\text{Sample}, \neg c_1)$ then return c_1 .
2. Else return $\neg c_1$. Break ties arbitrarily.

Figure 1: The Two-Balls Algorithm. c_1 is a Built-in Constant Concept.

both. In these cases, if the correct concept is c_2 and we map the sample in favor of c_1 , we will commit a *mistake*, and we can only afford to do this with probability δ . The probability of getting a sample that is consistent with both c_1 and c_2 is $(\frac{2^n-d}{2^n})^m$, where m is the sample size. This quantity is decreasing as d increases, so if we choose d sufficiently large, we can keep the probability of a mistake below δ .

In short, if c_1 and c_2 are close together, then there is no problem, because we can choose an h ϵ -close to both. Conversely, if they are far apart, there is also no problem, because the probability of a mistake can be bounded by δ . This suggests that a good strategy for designing learning algorithms with high coverage is to choose a collection of concepts that is as large as possible, such that the distance between each pair of concepts in the collection is either

- sufficiently large to suppress the probability of getting a sample consistent with both concepts, or
- within $2\epsilon 2^n$, so that we can find concept(s) within ϵ of both concepts.

This means that the concepts to be learned must be clustered as one or more balls in the space of concepts. What we need to do in order to construct an appropriate algorithm is to keep the radius of each ball small enough, and at the same time, make the distance between the centers of the balls large enough.

As a trivial case, suppose that we want to learn the set of all concepts that are within ϵ of a fixed concept c —the set $Ball(c, \epsilon)$. This is achieved simply by returning c as the hypothesis regardless of the training sample. This leads to a coverage of $\sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$.

A less trivial case is to learn 2 ϵ -balls of concepts. This is accomplished by the “Two-Balls” algorithm given in Figure 1. This algorithm returns, as the hypothesis, some fixed concept c_1 or its complement, whichever is *closer* to the training sample. For any concept c in $Ball(c_1, \epsilon)$, the probability of drawing an example of c that disagrees with c_1 (and thus, agrees with $\neg c_1$) is at most ϵ . The same argument applies to the concepts in $Ball(\neg c_1, \epsilon)$. Therefore, if $\epsilon \leq \delta$, then a sample of size 1 (that is, $m = 1$) is sufficient to learn these two balls. The following theorem gives the sample size sufficient to achieve this goal in the general case.

Theorem 4 For accuracy and confidence parameters ϵ and δ , the coverage of the Two-Balls algorithm is

$$2 \sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$$

when given a sample of size

$$m > \frac{12\epsilon}{(1-2\epsilon)^2} \ln \frac{1}{\delta}$$

assuming that sampling is with replacement under the uniform distribution and that $\epsilon < \frac{1}{2}$.

Proof: Without loss of generality, assume that c_1 in the algorithm is just the *nil* concept. Let c be a concept with distance at most $\epsilon 2^n$ from the *nil* concept. It is enough to show that the given sample size is sufficient to make the algorithm learn c .

Let Z denote the number of positive examples in a sample of c of size m . Since c has at most $\epsilon 2^n$ positive examples, Z can be viewed as a binomial random variable of m trials and at most ϵ as the ratio of success. A sample of c is mapped to the *true* (instead of the *nil*) concept only if it has more positive examples than negative examples. The probability of this is bounded above by

$$Pr[Z \geq \frac{m}{2}] \leq \sum_{i=\lceil \frac{m}{2} \rceil}^m \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i}.$$

Using Chernoff's bound (Lemma 4), this is at most $e^{-\frac{(1-2\epsilon)^2 m}{12\epsilon}}$. Therefore, c is learned if

$$e^{-\frac{(1-2\epsilon)^2 m}{12\epsilon}} < \delta,$$

which is satisfied if

$$m > \frac{12\epsilon}{(1-2\epsilon)^2} \ln \frac{1}{\delta}$$

as desired. \square

According to this theorem, even when δ is as small as 0.001 (i.e. 99.9% confidence), for any ϵ in the range $0 < \epsilon \leq 0.1$, the required sample size is only 13 examples, *independent* of n .

Since the sample size is usually much larger than this, a direct generalization of the above trivial cases is to attempt to learn as many balls of concepts as permitted by the sample size. The idea is to choose a collection of well-separated concepts (the centers of the balls) and attempt to learn all the concepts clustered around each of these centers. More specifically, we start by fixing two positive integers d and k , and then construct a set $\mathcal{H} = \{h_1, h_2, h_3, \dots, h_k\}$ of k concepts such that the distance between each pair of concepts in \mathcal{H} is at least d . Then given a training sample, the concept in \mathcal{H} that has the minimum disagreement with the sample is returned as the hypothesis.

The goal of the algorithm is to learn all the concepts in $\bigcup_{h \in \mathcal{H}} \text{Ball}(h, \epsilon)$ which will give a coverage of $k \sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$, provided that the intersection between the balls is empty (that is,

$d > 2\epsilon 2^n$). The question is, of course, how to determine the appropriate values of d and k such that this goal is accomplished. Particularly, we need to worry about the following:

1. As explained earlier, d must be large enough so that the interaction between concepts in different balls is kept within what is allowed by the confidence parameter δ .
2. The number of concepts that we can construct such that the minimum distance between each pair is d drops sharply as d increases. Therefore, making d too large causes k (and hence, the coverage of the algorithm) to be too small.

The following two lemmas show how to choose appropriate values for d and k .

Lemma 1 *For $0 < \epsilon < \frac{1}{4}$ and $2\epsilon < \alpha < \frac{1}{2}$, let $\mathcal{H} = \{h_1, h_2, \dots, h_k\}$ be a set of concepts such that the distance between each pair $h_i, h_j \in \mathcal{H}$ is at least $d = \lceil \alpha 2^n \rceil$, and let $c \in \text{Ball}(h, \epsilon)$ for some $h \in \mathcal{H}$. Assume that L is a learning algorithm that on any sample S outputs some hypothesis $h_i \in \mathcal{H}$ that has minimal disagreement with S . Then, under the uniform distribution, the probability that a sample of c of size m is mapped by L to an hypothesis other than h is at most*

$$\min_{0 < \beta < 1} k \cdot \{e^{-2(1-\beta)^2 \alpha^2 m} + e^{-\beta \frac{(\alpha-2\epsilon)^2}{2\alpha} m}\} . \quad (1)$$

Proof: Let $g \neq h$ be a *specific* concept in \mathcal{H} and let us bound the probability that a sample of c is mapped to g (instead of h). We consider the case in which this probability is maximized. First, we assume that g is as close as possible to h —that is, $\text{distance}(h, g) = d$. Second, we place c to be as far as possible from h and as close as possible to g —that is, $\text{distance}(c, h) = \epsilon 2^n$ and $\text{distance}(c, g) = d - \epsilon 2^n$.

Let A be the set of objects where c agrees with g but not with h and let B be the set of objects where c agrees with h but not with g . Thus, $|A| = \epsilon 2^n$ and $|B| = d - \epsilon 2^n$. A sample S of c is mapped to g only if $\text{disagreement}(S, g) \leq \text{disagreement}(S, h)$. That is, only if S has more examples drawn from A than from B . Now, let us define the following three random variables:

- X : the number of examples in S drawn from A .
- Y : the number of examples in S drawn from B .
- $Z = X + Y$: the number of examples in S drawn from $A \cup B$.

Then, the probability that S is mapped to g is just $\Pr[X \geq Y]$.

Note that Z is just a binomial random variable with m trials and $\frac{d}{2^n}$ as the ratio of success. Given that, out of the m examples in S , there are i examples drawn from $A \cup B$ (i.e. given that $Z = i$), the event $X \geq Y$ is equivalent to $Y \leq \frac{i}{2}$. In this case, Y can be

viewed as a binomial random variable of i trials and $\frac{d-\epsilon 2^n}{d}$ as the ratio of success. Therefore, we can write

$$\begin{aligned}
Pr[X \geq Y] &= \sum_{i=0}^m Pr[Z = i] Pr[X \geq Y | Z = i] \\
&= \sum_{i=0}^m Pr[Z = i] Pr[Y \leq \frac{i}{2} | Z = i] \\
&= \sum_{i=0}^m \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} \sum_{j=0}^{\frac{i}{2}} \binom{i}{j} \left(\frac{d - \epsilon 2^n}{d}\right)^j \left(\frac{\epsilon 2^n}{d}\right)^{i-j} \right].
\end{aligned}$$

Applying Hoeffding's bound (Lemma 3) on the inner summation, this is at most

$$\begin{aligned}
&\sum_{i=0}^m \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} e^{-2i\left(\frac{d-2\epsilon 2^n}{2d}\right)^2} \right] \\
&\leq \sum_{i=0}^{\lfloor \beta \frac{d}{2^n} m \rfloor} \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} e^{-2i\left(\frac{d-2\epsilon 2^n}{2d}\right)^2} \right] + \\
&\quad \sum_{i=\lceil \beta \frac{d}{2^n} m \rceil}^m \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} e^{-2i\left(\frac{d-2\epsilon 2^n}{2d}\right)^2} \right]
\end{aligned}$$

for any $0 < \beta < 1$. Letting $i = 0$ in the first exponential and $i = \beta \frac{d}{2^n} m$ in the second, the above is at most

$$\begin{aligned}
&\sum_{i=0}^{\lfloor \beta \frac{d}{2^n} m \rfloor} \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} \right] + \\
&\quad e^{-2\beta \frac{d}{2^n} \left(\frac{d-2\epsilon 2^n}{2d}\right)^2 m} \sum_{i=\lceil \beta \frac{d}{2^n} m \rceil}^m \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} \right] \\
&\leq \sum_{i=0}^{\lfloor \beta \frac{d}{2^n} m \rfloor} \left[\binom{m}{i} \left(\frac{d}{2^n}\right)^i \left(\frac{2^n - d}{2^n}\right)^{m-i} \right] + e^{-2\beta \frac{d}{2^n} \left(\frac{d-2\epsilon 2^n}{2d}\right)^2 m}.
\end{aligned}$$

Applying Hoeffding's bound (Lemma 3) again on the first term, the above is at most

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

Since the above bound holds for any $0 < \beta < 1$, we can, of course, minimize over all β in that range and get

$$\min_{0 < \beta < 1} \left\{ e^{-2(1-\beta)^2 \left(\frac{d}{2^n}\right)^2 m} + e^{-\beta \frac{(\frac{d}{2^n} - 2\epsilon)^2}{2 \frac{d}{2^n}} m} \right\}.$$

This quantity is decreasing in $\frac{d}{2^n}$ when $\frac{d}{2^n} > 2\epsilon$. Thus, we can replace $\frac{d}{2^n}$ by α since $2\epsilon < \alpha \leq \frac{d}{2^n}$. Since g can be any of the k concepts in \mathcal{H} , multiplying the above probability by k gives the desired result. \square

Note that d in this lemma is expressed as a fraction α of 2^n . This result says that if the conditions of the lemma are met, and if c is a concept in one of the k ϵ -balls, then the samples of c are usually mapped by L to the center of that ball (which is ϵ -close to c) except with a probability that is bounded by the quantity of Equation (1). Thus, α and k must be chosen so that this probability is at most δ . It is important to note that this probability is diminishing in α and m and independent of n .

The problem of finding a collection of bit vectors that maintain a given minimum pairwise distance is well studied in the field of Error-Correcting Coding Theory. Specifically, the following theorem states how large k can (at least) be for a given separation distance d .

Lemma 2 (The Gilbert-Varshamov Bound) *There exist at least 2^{l-r} bit vectors of length l and minimum pairwise distance d , where r is any integer satisfying*

$$\binom{l-1}{0} + \binom{l-1}{1} + \binom{l-1}{2} + \cdots + \binom{l-1}{d-2} < 2^r.$$

A proof of this lemma, in addition to a method of actually constructing the l -bit vectors as given above, can be found in many references in the literature of Error-Correcting Coding Theory (e.g., [Peterson and Weldon 72]).

For our purposes here, it is convenient to draw the following corollary from the Gilbert-Varshamov bound given above.

Corollary 1 *For any α , $0 < \alpha < \frac{1}{2}$, and any even positive integer l , there exist at least $2^{l-\lceil H(\alpha)l \rceil}$ bit vectors of length l and pairwise distance at least αl , where $H(\alpha) = \alpha \log_2 \frac{1}{\alpha} + (1-\alpha) \log_2 \frac{1}{1-\alpha}$.*

Proof: We just need to show that when $d = \lceil \alpha l \rceil$ and $r = \lceil H(\alpha)l \rceil$, the inequality of Lemma 2 is satisfied. The left-hand side of the inequality is just

$$\begin{aligned} \sum_{i=0}^{d-2} \binom{l-1}{i} &< \sum_{i=0}^{\alpha l} \binom{l}{i} \\ &\leq 2^{H(\alpha)l} \quad (\text{by Lemma 5}) \\ &\leq 2^{\lceil H(\alpha)l \rceil}, \end{aligned}$$

which equals the right-hand side. This proves the corollary. \square

Using Lemma 1 and Corollary 1, one can search for the appropriate value for d that leads to learning k different ϵ -balls, for k as large as possible. Let's now compute a lower bound on the coverage that can be achieved by this approach.

Figure 2 shows the “Multi-Balls” algorithm in which we give a specific way of choosing the values of α (and hence, d) and k . To be able to give a lower bound on the coverage of

Algorithm Multi-Balls (*Sample*, ϵ , δ)

1. Find α in the range $2\epsilon < \alpha < \frac{1}{2}$ such that

$$1 - H(\alpha) = 2[1 - \beta(\alpha)]^2 \alpha^{\frac{2}{2^n}} \log_2 e$$
 where

$$H(\alpha) = \alpha \log_2 \frac{1}{\alpha} + (1 - \alpha) \log_2 \frac{1}{1 - \alpha}, \text{ and}$$

$$\beta(\alpha) = 1 + \frac{(\alpha - 2\epsilon)^2}{8\alpha^3} - \sqrt{\left[1 + \frac{(\alpha - 2\epsilon)^2}{8\alpha^3}\right]^2 - 1}.$$
2. Let $k = \left\lfloor 2^{2^n - H(\alpha)2^n} \times \frac{\delta}{2} \right\rfloor$.
3. Construct $\mathcal{H} = \{h_1, h_2, h_3, \dots, h_k\}$ such that

$$\forall_{h_i, h_j \in \mathcal{H}} \text{distance}(h_i, h_j) \geq \lceil \alpha 2^n \rceil.$$
4. Return some hypothesis in \mathcal{H} that has minimal disagreement with *Sample* (break ties arbitrarily).

Figure 2: The Multi-Balls Algorithm.

this algorithm, we need the following definition.

Definition: For $0 \leq \theta \leq 1$ and $0 < \epsilon < \frac{1}{4}$, define $\rho(\theta, \epsilon)$ as

$$\rho(\theta, \epsilon) = \frac{1 - H(\hat{\alpha})}{\theta}$$

for $\hat{\alpha}$ being the solution of the equation

$$1 - H(\alpha) = 2[1 - \beta(\alpha)]^2 \alpha^2 \theta \log e \quad (2)$$

in the range $2\epsilon < \alpha < \frac{1}{2}$, where

$$\beta(\alpha) = 1 + \frac{(\alpha - 2\epsilon)^2}{8\alpha^3} - \sqrt{\left[1 + \frac{(\alpha - 2\epsilon)^2}{8\alpha^3}\right]^2 - 1}$$

and

$$H(\alpha) = \alpha \log_2 \frac{1}{\alpha} + (1 - \alpha) \log_2 \frac{1}{1 - \alpha}.$$

□

Lemma 10 shows that the above $\rho(\theta, \epsilon)$ function is well-defined. Although this function is not provided in closed form, it is easily computed for any given values of θ and ϵ using standard numerical methods. For illustration, Figure 3 plots this function over the range $0 \leq \theta \leq \frac{1}{2}$ for $\epsilon = 0.01, 0.05$ and 0.10 .

Using the above definition of ρ , a lower bound on the coverage of Multi-Balls can be stated as follows:

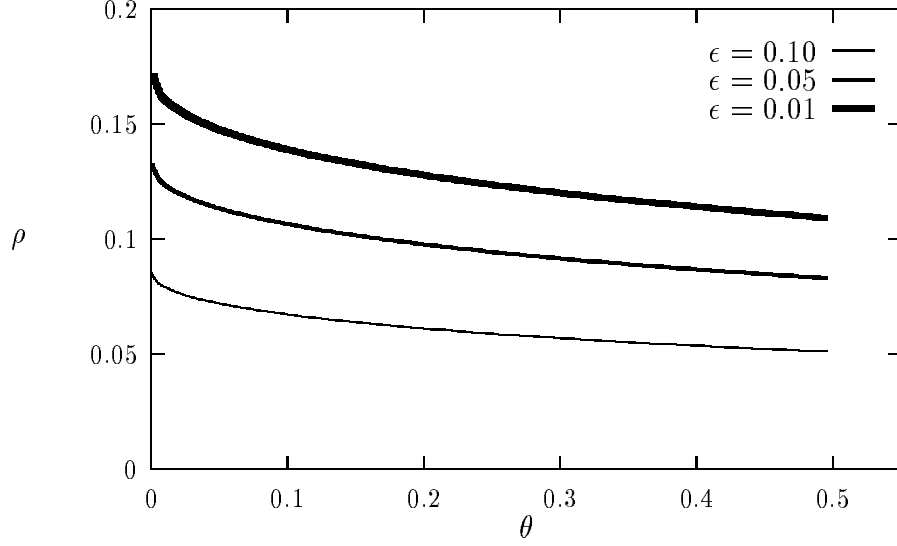


Figure 3: The Function $\rho(\theta, \epsilon)$.

Theorem 5 For $0 < \epsilon < \frac{1}{4}$ and $0 < \delta < 1$, the coverage of the Multi-Balls algorithm under the uniform distribution is at least

$$\left\lfloor \frac{\delta}{2} 2^{\rho(\frac{m}{2^n}, \epsilon) m} \right\rfloor \sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$$

for sample size m , accuracy parameter ϵ and confidence parameter δ .

Proof: First, the existence of α in the range $2\epsilon < \alpha < \frac{1}{2}$ satisfying the equation of Step 1 is guaranteed by Lemma 10. It is sufficient to show the following three claims for the values of α and k chosen by the algorithm in Steps 1 and 2:

Claim 1: Step 3 is executable. That is, for α and k as chosen in Steps 1 and 2, we can actually construct a set of k concepts such that the minimum pairwise distance is $\lceil \alpha 2^n \rceil$.

Claim 2: For any $h \in \mathcal{H}$ and any $c \in \text{Ball}(h, \epsilon)$, the probability that a sample of c of size m is mapped to an hypothesis other than h is at most δ .

Claim 3: The number of balls, k , is at least $\left\lfloor \frac{\delta}{2} 2^{\rho(\frac{m}{2^n}, \epsilon) m} \right\rfloor$.

Claim 1 follows from Corollary 1 since k is set to $\left\lfloor 2^{2^n - H(\alpha)2^n} \times \frac{\delta}{2} \right\rfloor$ and thus

$$k \leq 2^{2^n - H(\alpha)2^n} \times \frac{\delta}{2}$$

$$\begin{aligned}
&= 2^{2^n - H(\alpha)2^n - \log_2 \frac{2}{\delta}} \\
&\leq 2^{2^n - \lceil H(\alpha)2^n \rceil + 1 - \log_2 \frac{2}{\delta}} \\
&< 2^{2^n - \lceil H(\alpha)2^n \rceil}
\end{aligned}$$

since $\log_2 \frac{2}{\delta}$ is necessarily greater than 1.

For Claim 2, we use Lemma 1 which states that the probability of a mistake is bounded above by

$$k \cdot \left\{ e^{-2(1-\beta)^2 \alpha^2 m} + e^{-\beta \frac{(\alpha-2\epsilon)^2}{2\alpha} m} \right\} \quad (3)$$

for *any* β in the range $[0, 1]$. Let us choose β such that

$$e^{-2(1-\beta)^2 \alpha^2 m} = e^{-\beta \frac{(\alpha-2\epsilon)^2}{2\alpha} m}$$

which means that

$$2(1-\beta)^2 \alpha^2 m = \beta \frac{(\alpha-2\epsilon)^2}{2\alpha} m. \quad (4)$$

Solving this quadratic equation for β gives the two solutions

$$1 + \frac{(\alpha-2\epsilon)^2}{8\alpha^3} \pm \sqrt{\left[1 + \frac{(\alpha-2\epsilon)^2}{8\alpha^3}\right]^2 - 1}. \quad (5)$$

It can be checked that Equation (4) has two roots, one in the range $[0, 1]$ and the other larger than 1. Thus, the smaller root in Equation (5), namely

$$1 + \frac{(\alpha-2\epsilon)^2}{8\alpha^3} - \sqrt{\left[1 + \frac{(\alpha-2\epsilon)^2}{8\alpha^3}\right]^2 - 1}$$

is between 0 and 1, and thus, is an appropriate value for β . Note that this is exactly the quantity $\beta(\alpha)$ in Step 1 of the algorithm.

Now, substituting for k and β in the bound of Lemma 1, the probability that Multi-Balls commits a mistake is at most

$$\begin{aligned}
&\min_{0 < \beta < 1} k \times \left\{ e^{-2[1-\beta]^2 \alpha^2 m} + e^{-\beta \frac{(\alpha-2\epsilon)^2}{2\alpha} m} \right\} \\
&\leq \frac{\delta}{2} \times 2^{2^n \lceil 1-H(\alpha) \rceil} \times \left\{ e^{-2[1-\beta(\alpha)]^2 \alpha^2 m} + e^{-\beta(\alpha) \frac{(\alpha-2\epsilon)^2}{2\alpha} m} \right\} \\
&= \frac{\delta}{2} \times 2^{2^n \lceil 1-H(\alpha) \rceil} \times 2 e^{-2[1-\beta(\alpha)]^2 \alpha^2 m}.
\end{aligned}$$

The value of α is chosen in Step 1 such that

$$1 - H(\alpha) = 2[1 - \beta(\alpha)]^2 \alpha^2 \frac{m}{2^n} \log_2 e.$$

Therefore, continuing from above, the desired probability is at most

$$\begin{aligned}
& \delta \times 2^{2^n 2[1-\beta(\alpha)]^2 \alpha^2 \frac{m}{2^n} \log_2 e} e^{-2[1-\beta(\alpha)]^2 \alpha^2 m} \\
&= \delta \times e^{2[1-\beta(\alpha)]^2 \alpha^2 m} e^{-2[1-\beta(\alpha)]^2 \alpha^2 m} \\
&= \delta,
\end{aligned}$$

which shows Claim 2.

Finally, Claim 3 follows immediately from the definition of ρ . Since the value of α found in Step 1 is equivalent to $\hat{\alpha}$ in the definition of ρ , we can write

$$\begin{aligned}
\rho\left(\frac{m}{2^n}, \epsilon\right) m &= \frac{1 - H(\alpha)}{\frac{m}{2^n}} m \\
&= 2^n [1 - H(\alpha)].
\end{aligned}$$

Thus, k is equivalent to $\left\lfloor \frac{\delta}{2} 2^{\rho(\frac{m}{2^n}, \epsilon) m} \right\rfloor$, which completes the proof. \square

More specific bounds on the coverage of Multi-Balls can be obtained from Theorem 5 if upper bounds on $\frac{m}{2^n}$ and ϵ are assumed. For example, if we are interested only in the range $0 < \epsilon \leq 0.05$ and $0 \leq \frac{m}{2^n} \leq 0.25$ (which are reasonable assumptions in practice), then the coverage of Multi-Balls as given by Theorem 5 is at least

$$\left\lfloor \frac{\delta}{2} 2^{0.094 m} \right\rfloor \sum_{i=0}^{\epsilon 2^n} \binom{2^n}{i}$$

where the constant 0.094 is just the value of $\rho(0.25, 0.05)$.

Note that the upper bound of Theorem 2 can be written as

$$\frac{1}{1 - \delta} 2^{(1 - 1.44\epsilon)m + 1} \binom{2^n}{\epsilon 2^n}$$

The main difference between this upper bound and the above lower bound is the coefficient of m in the exponent of 2. Therefore, for any fixed δ , this lower bound indicates that to achieve a given coverage, Multi-Balls requires a sample size that is within a constant factor of that required by an optimal learning algorithm.

5 The Large-Ball Algorithm

So far we have been trying to maximize the coverage by learning as many balls of concepts as possible, while fixing the radius of each ball at $\epsilon 2^n$. An alternative approach to increase the coverage is to learn ball(s) of concepts with larger radius. Because of the extremely high dimensionality of the space of concepts, any small increment in a ball's radius results in a huge increase in the number of concepts contained in the ball.

It turns out that learning a single ball of radius larger than $\epsilon 2^n$ is a surprisingly easy task that can be achieved by the simple algorithm “Large-Ball” given in Figure 4. This algorithm

Algorithm: Large-Ball (*Sample*)

1. Let c_1 be a constant concept.
2. Define the concept s as: $s(X) = \begin{cases} 1 & \text{if } X \in \text{Sample} \\ 0 & \text{otherwise} \end{cases}$
3. Define the concept p as: $p(X) = \begin{cases} 1 & \text{if } X \in \text{Sample and } X \text{ is a positive example} \\ 0 & \text{otherwise} \end{cases}$
4. Return $h = (\neg s \wedge c_1) \vee p$.

Figure 4: The Large-Ball Learning Algorithm.

works by modifying a *default* hypothesis c_1 so that the final hypothesis fully agrees with the training sample. For example, suppose c_1 is the *nil* concept. Then this algorithm classifies all examples as negative unless they appeared as positive examples in the training sample! The coverage of this algorithm is computed by the following theorem.

Theorem 6 *For sample size m , accuracy parameter ϵ and confidence parameter δ , the coverage of the Large-Ball algorithm under the uniform distribution is*

$$\sum_{i=0}^{\epsilon 2^n + \beta} \binom{2^n}{i}$$

where β is the largest integer such that

$$\sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k} \left(\frac{2^n - \epsilon 2^n - \beta + k}{2^n} \right)^m \left\{ 1 - \sum_{i=1}^k \binom{k}{i} \left(\frac{2^n - \epsilon 2^n - \beta + k - i}{2^n - \epsilon 2^n - \beta + k} \right)^m (-1)^{i+1} \right\} \leq \delta.$$

Proof: Without loss of generality, assume that c_1 of the algorithm is the *nil* concept. Suppose c is a concept such that $\text{distance}(c, c_1) = \epsilon 2^n + \beta$ for some integer β . That is, c has exactly $\epsilon 2^n + \beta$ distinct positive examples. The algorithm classifies all examples as negative except those appearing in the training sample as positive examples. Therefore, any sample that has β or more positive examples is mapped by the algorithm to an hypothesis within ϵ of c . Thus, to prove the theorem, it is enough to show that the probability of getting a sample with less than β *distinct* positive examples is just the left-hand side of the inequality of the theorem.

Let

- F be the set of negative examples of the concept c . Thus, $|F| = 2^n - \epsilon 2^n - \beta$.
- T be the set of positive examples of the concept of c . Thus, $|T| = \epsilon 2^n + \beta$.

- D be the (random) number of distinct positive examples in a random sample of size m of the concept c .

What we need to compute is just

$$Pr[D < \beta] = \sum_{k=0}^{\beta-1} Pr[D = k] .$$

Now, let $S \subset T$ be some *specific* set of k distinct positive examples of the concept c , and let R be the event of getting a sample of the concept c of size m that contains exactly S and any number of negative examples, but no positive examples not in S . Then

$$Pr[D = k] = \binom{\epsilon 2^n + \beta}{k} Pr[R] .$$

Now, R occurs if and only if both

- event $R1$: No positive examples from $T - S$ are in the sample, and
- event $R2$: Each positive example in S is present in the sample

are true. Therefore, $Pr[R] = Pr[R1 \cdot R2] = Pr[R1]Pr[R2|R1]$, and so finding $Pr[R1]$ and $Pr[R2|R1]$ is all what we need.

First, it is not difficult to see that

$$Pr[R1] = \left(\frac{2^n - \epsilon 2^n - \beta + k}{2^n} \right)^m .$$

Second, to find $Pr[R2|R1]$, let $A_i, i = 1, 2, \dots, k$, be the event that the positive example number i of S is *not* included in a sample of size m drawn from $S \cup F$. Then

$$Pr[R2|R1] = 1 - Pr \left[\bigcup_{i=1}^k A_i \right] .$$

By the principle of inclusion and exclusion (e.g. [Ross 88]), we get

$$Pr \left[\bigcup_{i=1}^k A_i \right] = \sum_i Pr[A_i] - \sum_{i \neq j} Pr[A_i A_j] + \dots + (-1)^{k+1} Pr[A_1 A_2 \dots A_k] .$$

Now,

$$\begin{aligned} Pr[A_i] &= \left(\frac{2^n - \epsilon 2^n - \beta + k - 1}{2^n - \epsilon 2^n - \beta + k} \right)^m, 1 \leq i \leq k \\ Pr[A_i A_j] &= \left(\frac{2^n - \epsilon 2^n - \beta + k - 2}{2^n - \epsilon 2^n - \beta + k} \right)^m, 1 \leq i, j \leq k \text{ and } i \neq j \\ &\vdots \\ &\vdots \\ Pr[A_1 A_2 \dots A_k] &= \left(\frac{2^n - \epsilon 2^n - \beta}{2^n - \epsilon 2^n - \beta + k} \right)^m . \end{aligned}$$

Therefore,

$$Pr[\bigcup_{i=1}^k A_i] = \sum_{i=1}^k \binom{k}{i} \left(\frac{2^n - \epsilon 2^n - \beta + k - i}{2^n - \epsilon 2^n - \beta + k} \right)^m (-1)^{i+1}.$$

Assembling, we get

$$\begin{aligned} Pr[D < \beta] &= \sum_{k=0}^{\beta-1} P(D = k) \\ &= \sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k} Pr[R] \\ &= \sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k} \left(\frac{2^n - \epsilon 2^n - \beta + k}{2^n} \right)^m \times \\ &\quad \left\{ 1 - \sum_{i=1}^k \binom{k}{i} \left(\frac{2^n - \epsilon 2^n - \beta + k - i}{2^n - \epsilon 2^n - \beta + k} \right)^m (-1)^{i+1} \right\}, \end{aligned}$$

and the theorem follows. \square

It can be easily shown that a sample of size $\frac{1}{\epsilon} \ln \frac{1}{\delta}$ is sufficient to make β at least 1. The following theorem gives a lower bound on the value of β in general.

Theorem 7 *If $\epsilon < \frac{1}{2}$, $\epsilon 2^n > 1$, and $m \leq \frac{1}{2 \log_2 e} 2^n \approx 0.347 \times 2^n$, then setting*

$$\beta = \left\lfloor \frac{\epsilon m \log_2 e - \log \frac{1}{\delta} - 1}{n - \log \frac{1}{\epsilon}} \right\rfloor$$

satisfies the inequality of Theorem 6.

Proof: The left-hand side of the inequality of Theorem 6 is at most

$$\sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k} \left(\frac{2^n - \epsilon 2^n - \beta + k}{2^n} \right)^m \quad (6)$$

since $\{1 - \sum_{i=1}^k \binom{k}{i} \left(\frac{2^n - \epsilon 2^n - \beta + k - i}{2^n - \epsilon 2^n - \beta + k} \right)^m (-1)^{i+1}\}$ is just the probability $Pr[R1|R2]$ as defined in the proof of Theorem 6, and hence, can be at most 1. Also, since k is at most $\beta - 1$,

$$\begin{aligned} \left(\frac{2^n - \epsilon 2^n - \beta + k}{2^n} \right)^m &\leq \left(\frac{2^n - \epsilon 2^n - \beta + (\beta - 1)}{2^n} \right)^m \\ &= \left(\frac{2^n - \epsilon 2^n - 1}{2^n} \right)^m \\ &< (1 - \epsilon)^m \\ &\leq e^{-\epsilon m}. \end{aligned}$$

This is independent of k and thus can be moved outside the summation of Equation (6).

Next, we need to bound $\sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k}$. Here, we can use Lemma 8 in addition to the fact that

$$\binom{\epsilon 2^n + \beta}{\beta - 1} = \frac{\beta(\epsilon 2^n + \beta)}{\epsilon 2^n(\epsilon 2^n + 1)} \binom{\epsilon 2^n + \beta - 1}{\beta} . \quad (7)$$

These give

$$\begin{aligned} \sum_{k=0}^{\beta-1} \binom{\epsilon 2^n + \beta}{k} &= \sum_{k=0}^{\frac{\beta-1}{\epsilon 2^n + \beta}(\epsilon 2^n + \beta)} \binom{\epsilon 2^n + \beta}{k} \\ &\leq \frac{1 - \frac{\beta-1}{\epsilon 2^n + \beta}}{1 - 2 \frac{\beta-1}{\epsilon 2^n + \beta}} \binom{\epsilon 2^n + \beta}{\beta - 1} \\ &\quad \text{(by Lemma 8, provided that } \beta < \epsilon 2^n + 2 \text{)} \end{aligned} \quad (8)$$

$$\begin{aligned} &= \frac{\epsilon 2^n + 1}{\epsilon 2^n - \beta + 2} \binom{\epsilon 2^n + \beta}{\beta - 1} \\ &= \frac{\beta(\epsilon 2^n + \beta)}{\epsilon 2^n(\epsilon 2^n - \beta + 2)} \binom{\epsilon 2^n + \beta - 1}{\beta} \quad \text{(by Equation (7))} \\ &\leq \frac{\beta(\epsilon 2^n + \beta)}{\epsilon 2^n(\epsilon 2^n - \beta + 2)} (\epsilon 2^n)^\beta \\ &\leq \frac{3}{2} (\epsilon 2^n)^\beta \quad \text{(provided that } \beta \leq \frac{1}{2} \epsilon 2^n \text{).} \end{aligned} \quad (9)$$

Therefore, the left-hand side of the inequality of Theorem 6 is at most

$$\frac{3}{2} (\epsilon 2^n)^\beta (1 - \epsilon)^m \leq \frac{3}{2} (\epsilon 2^n)^\beta e^{-\epsilon m}$$

Setting this to be at most δ , and taking the logarithm of both sides of the inequality gives

$$\log_2 \frac{3}{2} + \beta(\log_2 \epsilon + n) - \epsilon m \log_2 e \leq \log_2 \delta .$$

This is certainly satisfied when

$$\beta = \left\lceil \frac{\epsilon m \log_2 e - \log_2 \frac{3}{2\delta}}{n - \log_2 \frac{1}{\epsilon}} \right\rceil .$$

Clearly, this is less than $\epsilon m \log_2 e$. Given the assumption that $m < \frac{1}{2 \log_2 e} 2^n$, the above value of β is at most $\frac{1}{2} \epsilon 2^n$, which satisfies the conditions of Equations (8) and (9). This completes the proof. \square

For fixed δ , this theorem says that the radius of the ball of concepts learned by the Large-Ball algorithm grows linearly in ϵm . Note that a unit increment in the radius of the ball corresponds to a very large increment in coverage. The coverage of Large-Ball, therefore, grows quite rapidly as the sample size increases.

The coverage lower bound obtained for Large-Ball appears to overlap the bound for Multi-Balls, although Multi-Balls gives higher coverage for small values of ϵ (e.g. 0.01). In any case, the fact that a trivial algorithm like Large-Ball achieves such high coverage suggests that coverage analysis alone is not strong enough to derive good inductive biases. Before considering this point further, let us measure the coverage of some popular learning algorithms.

6 Coverage of Current Learning Algorithms

The last problem investigated in this paper is the evaluation of the coverage of existing learning algorithms. Due to the difficulty of performing coverage analysis for empirical algorithms (e.g, ID3 [Quinlan 86]), one may consider measuring the coverage experimentally by running learning algorithms on every possible training sample (as done in [Dietterich 89]). However, this involves an immense amount of computation. Specifically, if the number of features is n and the sample size is m , then we have to run an algorithm on as many as $2^m \binom{2^n}{m}$ samaples and test the learnability of 2^{2^n} concepts. This is doable when $n = 3$ [Dietterich 89] or 4, but soon becomes unaffordable when $n = 5$.

To reduce these computational costs, we can employ the following two techniques: First, we can exploit the fact that most of the learning algorithms are symmetric with respect to permutations and/or negations of input features. More precisely, if an algorithm learns a concept represented by a Boolean function $f(x_1, x_2, \dots, x_i, \dots, x_j, \dots, x_n)$, then the same algorithm also learns the concepts represented by $f(x_1, x_2, \dots, x_j, \dots, x_i, \dots, x_n)$, $f(x_1, x_2, \dots, \bar{x}_i, \dots, x_j, \dots, x_n)$ and so on for all functions obtained by permuting and/or negating the features in f . These symmetry properties partition the space of concepts into equivalence classes such that it suffices to test one representative concept in each equivalence class to determine learnability for all concepts in the class. It turns out that for $n = 5$, the number of representative concepts one needs to consider is 1,228,158. Also, if the algorithm being tested is also symmetric with respect to complementing the target concept then this number is further reduced to 698,635. This is a considerable reduction since one has to consider more than 4 billion concepts in the exhaustive approach.

Second, we can measure the learnability of each concept statistically by running the learning algorithm on a *large* number of randomly-chosen samples, rather than *all* the possible samples. The goal is to estimate the ratio of the samples on which the algorithm returns an ϵ -close hypothesis, and check that this is at least $1 - \delta$ in order to consider a concept learned.

The details of the above cost-reduction techniques are not central to this study. For a thorough discussion of these, we refer the interested reader to [Almuallim 92].

6.1 Experimental Work

Using the above cost-reduction techniques, we experimentally measured the coverage figures for three algorithms: ID3 [Quinlan 86], FRINGE [Pagallo and Haussler 90] and MDT, which is an exhaustive algorithm that finds a decision tree with fewest nodes consistent with the

Table 1: Coverage Figures For Various Algorithms.

Algorithm	Sample size				
	8	10	12	14	16
ID3	12 \pm 0	332 \pm 0	396 \pm 0	1,756 \pm 0	4,954 \pm 640
FRINGE	12 \pm 0	332 \pm 0	396 \pm 0	1,756 \pm 0	5,284 \pm 970
MDT	12 \pm 0	12 \pm 0	116 \pm 40	496 \pm 0	3,694 \pm 0
Large-Ball	5,489	5,489	5,489	41,449	41,449
Two-Large-Balls	10,978	10,978	10,978	82,898	82,898
Upper Bound	661,333	2,041,173	6,148,551	17,985,991	50,753,991

training sample. The coverage of these algorithms was measured for $n = 5$, $\epsilon = \frac{3}{32}$, $\delta = 0.1$ and $m = 8, 10, 12, 14$, and 16. Sampling in these experiments was done *without* replacement.

FRINGE is symmetric with respect to permuting and negating the features. Thus, the number of representative concepts we tested for this algorithm was 1,228,158. ID3 and MDT are in addition symmetric with respect to complementing the target concept as well, and thus the number of representative concepts for these algorithms was only 698,635 concepts.

Determining whether or not a concept is learned by an algorithm in the above setting was done in two passes:

- In the first pass, only 100 randomly-chosen samples per concept were tested. Concepts for which the number of samples that resulted in ϵ -close hypothesis is less than 60 out of 100, were considered not learned and thus excluded.
- In the second pass, each algorithm was run on 10,000 randomly-chosen samples for each of the remaining concepts. Three outcomes are then considered: (i) If the ratio of the samples on which the algorithm returned an ϵ -close hypothesis is less than 0.893 then the concept is considered unlearned. (ii) If the ratio is greater than 0.907 then the concept is considered learned. (iii) If the ratio is between 0.893 and 0.907, then we *hesitate* to make either decisions.

Due to the third outcome in the above process, the coverage figures are eventually given as a range $a \pm b$. The number of samples tested per concept and the margin ± 0.007 were chosen so that the probability of a wrong decision (considering an unlearned concept as learned and vice versa) becomes within 0.01 [Almuallim 92], that is, to achieve 99% level of significance.

The final results of our experiments are summarized in Table 1.

6.2 Coverage of the “Balls” Approach

For comparison, let us compute the coverage obtained by the balls approach under the same conditions of our experiments. We look at two algorithms:

Large-Ball: Without loss of generality, assume that c_1 in Large-Ball is just the *nil* concept.

Algorithm: Two-Large-Balls(*Sample*)

1. If $\text{disagreement}(\text{Sample}, c_1) < \text{disagreement}(\text{Sample}, \neg c_1)$, then $h = c_1$,
 else, if $\text{disagreement}(\text{Sample}, c_1) > \text{disagreement}(\text{Sample}, \neg c_1)$, then $h = \neg c_1$.
 Otherwise, arbitrarily set h to c_1 or $\neg c_1$.

2. Define the function s as follows:

$$s(X) = \begin{cases} 1 & \text{if } X \in \text{Sample} \\ 0 & \text{otherwise} \end{cases}$$

3. Define the function p as follows:

$$p(X) = \begin{cases} 1 & \text{if } X \in \text{Sample} \text{ and } X \text{ is a positive example} \\ 0 & \text{otherwise} \end{cases}$$

4. Return $(\neg s \wedge h) \vee p$.

Figure 5: The Two-Large-Balls algorithm.

This means that the algorithm guesses *negative* for all the examples not included in the sample. Let us define the *weight* of a concept as the number of positive examples in that concept. Then, the algorithm trivially learns all the concepts of weight 0 to 3 (i.e., $\text{Ball}(0, \frac{3}{32})$).

Consider a concept c of weight 4. It should be obvious that any sample of c having one or more positive examples is mapped within ϵ from c . The only samples of c that are mapped to an ϵ -far hypothesis (the *nil* hypothesis) are those consisting of negative examples only. The probability of getting such a sample is just

$$\frac{\binom{32-4}{m}}{\binom{32}{m}}$$

which is less than 0.1 when $m \geq 14$. Therefore, the coverage is $\sum_{i=0}^3 \binom{32}{i} = 5,489$ when $m = 8, 10$ or 12 , and $\sum_{i=0}^4 \binom{32}{i} = 41,449$ when $m = 14$ or 16 .

The Two-Large-Balls algorithm: Consider the algorithm given in Figure 5. This algorithm is just a consistent version of the Two-Balls algorithm (Figure 1)—it modifies the final hypothesis so that there is no disagreement with the training sample.

Again, without loss of generality, assume that c_1 is the *nil* concept. This means that the algorithm classifies all the examples that are not in the training sample as positive if the majority of the examples in the sample are positive, or as negative otherwise (breaking ties arbitrarily). For those examples included in the training sample, the algorithm gives the same class as given in the sample.

If $m \geq 7$ then Two-Large-Balls behaves exactly like Large-Ball (with $c_1 = \text{nil}$) for all the concepts of weight up to 3, since at the end of Step 1, h will definitely be the *nil* concept. If $m \geq 9$, then this can also be said about the concepts of weight 4.

Similarly, h will be the *true* concept with certainty if the target concept has weight 29

to 32 and $m \geq 7$, or if the target concept has weight 28 and $m \geq 9$. Thus, for these cases, Two-Large-Balls will again behave exactly as Large-Ball but with $c_1 = \text{true}$.

As a result, the coverage of Two-Large-Balls is twice as that of Large-Ball, that is 10,978 when $m = 8, 10$ or 12 and 82,898 when $m = 14$ or 16.

6.3 Discussion

Three important points can be seen in the coverage figures of Table 1:

1. MDT does not give better coverage than the heuristic algorithms ID3 and FRINGE.
2. The coverage of ID3 and FRINGE is disappointingly smaller than that of Large-Ball and Two-Large-Balls.
3. The coverage of all these algorithms is far below the upper bound of Theorem 1.

7 Conclusion and Future Extensions

We began this paper by suggesting that an important design criterion for learning algorithms should be the coverage of the algorithm. We presented the Multi-Balls algorithm and showed that it can achieve optimal coverage with a sample size that is within a constant factor of optimal. However, we then showed that a fairly trivial algorithm, Large-Ball, can also achieve very large coverage—larger than Multi-Balls in cases where ϵ is reasonably big. Experimental tests confirm that Large-Ball and BALLS, a variation of Multi-Balls, have much better coverage than the popular ID3 algorithm and its relatives.

Why does Large-Ball strike us as trivial? Because it merely memorizes the training sample—it does not attempt to find any regularity in the data. Furthermore, the concepts it learns, while they are very numerous, are all located near a single concept. In short, the bias of Large-Ball is unlikely to be appropriate in real-world learning situations. This argument shows that coverage analysis alone is not sufficient to find a practically-useful inductive bias.

This suggests that we combine coverage analysis with other methods for choosing inductive bias. For example, in [Almuallim and Dietterich 91], we described learning situations in which the MIN-FEATURES bias—the bias that prefers consistent concepts definable over fewer features—is appropriate. However, the MIN-FEATURES bias does not uniquely define a learning algorithm, because, given a training sample, there are typically many consistent hypotheses that have the same, minimum, number of features. Hence, *within* the MIN-FEATURES bias, we could apply coverage analysis to design a learning algorithm that has the largest coverage among all algorithms that implement MIN-FEATURES.

In general, let $Pref(c_1, c_2)$ be a preference bias that prefers c_1 to c_2 in all cases where both concepts are consistent with the training sample. Let $Learns(L, c, m, \epsilon, \delta)$ be true if algorithm L can learn concept c from a sample of size m with accuracy and confidence parameters ϵ and δ . The *coverage within bias Pref* for L (with respect to m , ϵ , and δ), is the size of the

set

$$C = \{c \mid \text{Learns}(L, c, m, \epsilon, \delta) \text{ and} \\ \forall c' \text{ Pref}(c', c) \Rightarrow \text{Learns}(L, c', m, \epsilon, \delta)\}$$

That is, a concept is “covered” only if all concepts preferred to it are also covered.

In conclusion, the results from this paper suggest that an important problem for future research is to design and analyze algorithms that have optimal coverage-within-bias for many of the popular biases. This will be particularly important for biases that are so weak that they do not have polynomial sample complexity.

Acknowledgements

Hussein Almuallim was supported by a scholarship from the University of Petroleum and Minerals, Saudi Arabia. The authors gratefully acknowledge the support of the NSF under grant number IRI-86-57316. Thanks to Bella Bose for clarifying some issues on the Gilbert-Varshamov bound, to Rob Holte for useful discussions and for providing [Holte 91], and to Prasad Tadepalli for comments on an earlier draft of the paper.

References

- [Almuallim and Dietterich 91] Almuallim, H. and Dietterich, T. G. Learning With Many Irrelevant Features. Proceedings of the 9th National Conference on Artificial Intelligence (AAAI-91), 547–552, 1991.
- [Almuallim 91] Almuallim, H. Exploiting Symmetry Properties in the Evaluation of Inductive Learning Algorithms: An Empirical Domain-Independent Comparative Study. Technical Report, 1991-30-09, Dept. of Computer Science, Oregon State University, Corvallis, OR 97331-3202.
- [Almuallim 92] Almuallim, H. Concept Coverage and Its Application to Two Learning Tasks. Ph.D. Thesis. Department of Computer Science, Oregon State University, Corvallis, Oregon. Forthcoming, 1992.
- [Blumer et.al. 87] Blumer, A.; Ehrenfeucht, A.; Haussler, D.; and Warmuth, M. Learnability and the Vapnik-Chervonenkis Dimension, Technical Report UCSC-CRL-87-20, Department of Computer and Information Sciences, University of California, Santa Cruz, Nov. 1987. Also in *Journal of ACM*, 36(4):929–965, 1990.
- [Chernoff 52] Chernoff, H. A Measure of Asymptotic Efficiency for Tests of Hypothesis Based on the Sums of Observations. *Annals of Mathematical Statistics*, 23, pp. 493–509, 1952.
- [COLT 88] Proceedings of the 1988 Workshop on Computational Learning Theory. Haussler, D. and Pitt, L., Editors. Morgan Kaufmann Publishers, 1988.

- [COLT 89] Proceedings of the Second Annual Workshop on Computational Learning Theory. Rivest, R., Haussler, D. and Warmuth, M.K., Editors. Morgan Kaufmann Publishers, 1989.
- [COLT 90] Proceedings of the Third Annual Workshop on Computational Learning Theory. Fulk, M.A. and Case, J., Editors. Morgan Kaufmann Publishers, 1990.
- [COLT 91] Proceedings of the Fourth Annual Workshop on Computational Learning Theory. Valiant, L.G. and Warmuth, M., Editors. Morgan Kaufmann Publishers, 1991.
- [Dietterich 89] Dietterich, T. G. Limitations on inductive learning. In Proceedings of the Sixth International Conference on Machine Learning, 124–128. Ithaca, NY: Morgan Kaufmann, 1989.
- [Hoeffding 63] Hoeffding, W. Probability Inequalities for Sums of Bounded Random Variables. *Journal of The American Statistical Association*, 58, 13-3-, 1963.
- [Holte 91] Holte, R. C. Machine Learning as Error-Correction. Unpublished note, 1991.
- [Natarajan 87] Natarajan, B.K. On learning Boolean Functions. In Proceedings of the 19th Annual ACM Symposium on Theory of Computing 296–304. New York, NY, 1987.
- [Pagallo and Haussler 90] Pagallo, G.; and Haussler, D. Boolean feature discovery in empirical learning. *Machine Learning*, 5(1):71–100, 1990.
- [Peterson and Weldon 72] W. W. Peterson and E. J. Weldon. Error Correcting Codes, The MIT Press. p.86, 1972.
- [Quinlan 86] Quinlan, J. R. Induction of Decision Trees, *Machine Learning*, 1(1):81–106, 1986.
- [Rivest 87] Rivest, R. Learning Decision Lists. *Machine Learning*, 2(3):229–246, 1987.
- [Ross 88] Ross, S. A First Course in Probability. Macmillan Publishing Company, New York. 3rd edition, pp 111, 1988.
- [Valiant 84] L. G. Valiant. A Theory of the Learnable. *Communications of ACM*, 27(11):1134–1142, 1984.

Appendix

Lemma 3 [Hoeffding 63] *Let Y be a binomial random variable with t trials and p as the ratio of success. Then*

$$Pr[Y \geq rt] = \sum_{i=\lceil rt \rceil}^t \binom{t}{i} p^i (1-p)^{t-i} \leq e^{-2(r-p)^2 t}$$

for any r such that $p < r \leq 1$, and

$$Pr[Y \leq rt] = \sum_{i=0}^{\lfloor rt \rfloor} \binom{t}{i} p^i (1-p)^{t-i} \leq e^{-2(p-r)^2 t}$$

for any r such that $0 \leq r < p$.

Lemma 4 [Chernoff 52] *Let Y be a binomial random variable with t trials and p as the ratio of success. Then*

$$Pr[Y \geq rt] = \sum_{i=\lceil rt \rceil}^t \binom{t}{i} p^i (1-p)^{t-i} \leq e^{-\frac{(r-p)^2}{3p} t}$$

for any r such that $p < r \leq 1$, and

$$Pr[Y \leq rt] = \sum_{i=0}^{\lfloor rt \rfloor} \binom{t}{i} p^i (1-p)^{t-i} \leq e^{-\frac{(p-r)^2}{2p} t}$$

for any r such that $0 \leq r < p$.

Lemma 5 *For any $0 < \epsilon < \frac{1}{2}$,*

$$\sum_{i=0}^{\epsilon k} \binom{k}{i} \leq 2^{kH(\epsilon)}$$

where $H(\epsilon) = -\epsilon \log_2 \epsilon - (1-\epsilon) \log_2 (1-\epsilon)$.

Proof: For any integer r , it must be true that

$$\begin{aligned} 2^{r(1-\epsilon)k} \sum_{i=0}^{\epsilon k} \binom{k}{i} &\leq \sum_{i=0}^{\epsilon k} 2^{r(k-i)} \binom{k}{i} \\ &\leq \sum_{i=0}^k 2^{r(k-i)} \binom{k}{i} \\ &= \sum_{j=0}^k 2^{rj} \binom{k}{k-j} \\ &= \sum_{j=0}^k 2^{rj} \binom{k}{j} \\ &= (1 + 2^r)^k \end{aligned}$$

Thus

$$\begin{aligned} \sum_{i=0}^{\epsilon k} \binom{k}{i} &\leq \frac{(1 + 2^r)^k}{2^{r(1-\epsilon)k}} \\ &= \left[2^{-(1-\epsilon)r} + 2^{\epsilon r} \right]^k \end{aligned}$$

Now, letting $r = \log_2 \frac{1-\epsilon}{\epsilon}$ (which is clearly positive for $0 < \epsilon < \frac{1}{2}$), the last quantity is equal to

$$\begin{aligned}
& \left[2^{-(1-\epsilon)\log_2(1-\epsilon) + (1-\epsilon)\log_2 \epsilon} + 2^{\epsilon \log_2(1-\epsilon) - \epsilon \log_2 \epsilon} \right]^k \\
&= \left[2^{-(1-\epsilon)\log_2(1-\epsilon) - \epsilon \log_2 \epsilon + \log_2 \epsilon} + 2^{-\epsilon \log_2 \epsilon - (1-\epsilon)\log_2(1-\epsilon) + \log_2(1-\epsilon)} \right]^k \\
&= \left[2^{H(\epsilon)} \left\{ 2^{\log_2 \epsilon} + 2^{\log_2(1-\epsilon)} \right\} \right]^k \\
&= 2^{kH(\epsilon)} [\epsilon + 1 - \epsilon]^k \\
&= 2^{kH(\epsilon)}
\end{aligned}$$

which shows the lemma. \square

Lemma 6 *For any set S containing k distinct training examples and any concept $h \in 2^{U_n}$, the number of concepts that are consistent with the training examples in S and within ϵ of h is at most*

$$\sum_{i=0}^{\epsilon 2^n} \binom{2^n - k}{i}.$$

Proof: In how many ways can we construct a concept c such that c is consistent with S and within ϵ of h ? For the k bits of c that correspond to the k examples of S , we have no choice but to follow the classification of these examples as given in S . For the remaining $2^n - k$ bits, we can afford to disagree with h on at most $\epsilon 2^n$ bits. Thus, there are at most

$$\sum_{i=0}^{\epsilon 2^n} \binom{2^n - k}{i}$$

ways to construct such c , and the lemma follows. \square

Lemma 7 *For any $0 < \epsilon < \frac{1}{2}$ and any integer j such that $0 \leq j \leq \epsilon 2^n$, the quantity*

$$2^k \binom{2^n - k}{j}$$

is monotonically increasing in k in the range $1 \leq k < (1 - 2\epsilon)2^n$.

Proof:

$$\begin{aligned}
\frac{2^{k+1} \binom{2^n - k - 1}{j}}{2^k \binom{2^n - k}{j}} &= 2 \cdot \frac{(2^n - k - 1)(2^n - k - 2) \cdots (2^n - k - j)}{(2^n - k)(2^n - k - 1) \cdots (2^n - k - j + 1)} \\
&= 2 \cdot \frac{2^n - k - j}{2^n - k} \\
&= 2 \cdot \left(1 - \frac{j}{2^n - k} \right) \\
&> 2 \cdot \left(1 - \frac{\epsilon 2^n}{2\epsilon 2^n} \right) \quad (\text{by substituting } j = \epsilon 2^n \text{ and } k = (1 - 2\epsilon)2^n) \\
&= 1
\end{aligned}$$

This shows the lemma. \square

Lemma 8 For any ϵ such that $\epsilon < \frac{1}{2}$ and such that ϵk is an integer

$$\sum_{i=0}^{\epsilon k} \binom{k}{i} \leq \frac{1-\epsilon}{1-2\epsilon} \binom{k}{\epsilon k}.$$

Proof: The left-hand side of the inequality is

$$\sum_{i=0}^{\epsilon k} \binom{k}{i} = \binom{k}{\epsilon k} + \binom{k}{\epsilon k-1} + \cdots + \binom{k}{1} + \binom{k}{0}.$$

If we divide the second term by the first, the third by the second, the fourth by the third, and so on, we get

$$\frac{\epsilon k}{k - \epsilon k + 1}, \frac{\epsilon k - 1}{k - \epsilon k + 2}, \frac{\epsilon k - 2}{k - \epsilon k + 3}, \dots, \frac{2}{k-1}, \frac{1}{k}.$$

Note that these ratios are monotonically decreasing. Therefore, we can write

$$\begin{aligned} \sum_{i=0}^{\epsilon k} \binom{k}{i} &\leq \binom{k}{\epsilon k} \left[1 + \frac{\epsilon k}{k - \epsilon k + 1} + \left(\frac{\epsilon k}{k - \epsilon k + 1} \right)^2 + \cdots + \right. \\ &\quad \left. \left(\frac{\epsilon k}{k - \epsilon k + 1} \right)^{\epsilon k - 1} + \left(\frac{\epsilon k}{k - \epsilon k + 1} \right)^{\epsilon k} \right] \\ &= \binom{k}{\epsilon k} \sum_{i=0}^{\epsilon k} \left[\frac{\epsilon k}{k - \epsilon k + 1} \right]^i \\ &\leq \binom{k}{\epsilon k} \sum_{i=0}^{\infty} \left[\frac{\epsilon k}{k - \epsilon k + 1} \right]^i \\ &\leq \binom{k}{\epsilon k} \sum_{i=0}^{\infty} \left[\frac{\epsilon}{1 - \epsilon} \right]^i \\ &= \frac{1 - \epsilon}{1 - 2\epsilon} \binom{k}{\epsilon k}. \end{aligned}$$

The last step follows by viewing $\sum_{i=0}^{\infty} \left[\frac{\epsilon}{1 - \epsilon} \right]^i$ as an infinite sum of a geometric series. This evaluates to $\frac{1-\epsilon}{1-2\epsilon}$ provided that $\epsilon < \frac{1}{2}$. \square

Lemma 9 Provided that $\epsilon < \frac{1}{4}$ and $m < \frac{k}{4}$,

$$\sum_{i=0}^{\epsilon k} \binom{k-m}{i} \leq 2e^{-\epsilon m} \binom{k}{\epsilon k}.$$

Proof: First, note that

$$\begin{aligned}
\frac{\binom{k-m}{\epsilon k}}{\binom{k}{\epsilon k}} &= \frac{(k-m)!}{(k-\epsilon k-m)!(\epsilon k)!} \cdot \frac{(k-\epsilon k)!(\epsilon k)!}{k!} \\
&= \frac{(k-\epsilon k)(k-\epsilon k-1)\cdots(k-\epsilon k-m+1)}{k(k-1)\cdots(k-m+1)} \\
&\leq (1-\epsilon)^m \\
&\leq e^{-\epsilon m}
\end{aligned} \tag{10}$$

Now, the quantity $\sum_{i=0}^{\epsilon k} \binom{k-m}{i}$ can be rewritten as

$$\begin{aligned}
\sum_{i=0}^{\left(\frac{\epsilon k}{k-m}\right)(k-m)} \binom{k-m}{i} &\leq \frac{1 - \frac{\epsilon k}{k-m}}{1 - 2\frac{\epsilon k}{k-m}} \binom{k-m}{\epsilon k} \quad (\text{by Lemma 8}) \\
&= \frac{k-m-\epsilon k}{k-m-2\epsilon k} \binom{k-m}{\epsilon k} \\
&\leq \frac{k-m-\epsilon k}{k-m-2\epsilon k} e^{-\epsilon m} \binom{k}{\epsilon k} \quad (\text{by Equation (10)}).
\end{aligned}$$

The lemma follows by verifying that

$$\frac{k-m-\epsilon k}{k-m-2\epsilon k} \leq 2$$

when $\epsilon < \frac{1}{4}$ and $m < \frac{k}{4}$. This is equivalent to

$$k-m-\epsilon k \leq 2k-2m-4\epsilon k$$

or

$$k \geq m + 3\epsilon k,$$

which is satisfied by the assumptions of the lemma. \square

Lemma 10 *There exists a value for α in the range $2\epsilon < \alpha < \frac{1}{2}$ that satisfies Equation 2, provided that $0 < \epsilon < \frac{1}{4}$.*

Proof: First, note that $H(\alpha)$ is monotonically increasing in the range $0 < \alpha < \frac{1}{2}$. Moreover, H is always strictly less than 1 except when $\alpha = \frac{1}{2}$, where H becomes exactly 1.

Now, call the left-hand side of Equation 2 $L(\alpha)$ and the right-hand side $R(\alpha)$. The reader can confirm the following:

- $L(2\epsilon) > 0$, since $2\epsilon < \frac{1}{2}$ by assumption.
- $R(2\epsilon) = 0$, since $\beta(2\epsilon) = 1$.
- $L(\frac{1}{2}) = 0$, since $H(\frac{1}{2}) = 1$.
- $R(\frac{1}{2}) > 0$, since this is a product of positive quantities.

Therefore, $L(\alpha)$ and $R(\alpha)$ must intersect somewhere in the range $2\epsilon < \alpha < \frac{1}{2}$ and the lemma holds. \square