
Algorithmic Luckiness

Ralf Herbrich
Microsoft Research Ltd.
CB3 0FB Cambridge
United Kingdom
rherb@microsoft.com

Robert C. Williamson
Australian National University
Canberra 0200
Australia
Bob.Williamson@anu.edu.au

Abstract

In contrast to standard statistical learning theory which studies uniform bounds on the expected error we present a framework that exploits the specific learning algorithm used. Motivated by the luckiness framework [8] we are also able to exploit the serendipity of the training sample. The main difference to previous approaches lies in the complexity measure; rather than covering all hypotheses in a given hypothesis space it is only necessary to cover the functions which could have been learned using the fixed learning algorithm. We show how the resulting framework relates to the VC, luckiness and compression frameworks. Finally, we present an application of this framework to the maximum margin algorithm for linear classifiers which results in a bound that exploits both the margin and the distribution of the data in feature space.

1 Introduction

Statistical learning theory is mainly concerned with the study of *uniform* bounds on the expected error of hypotheses from a given hypothesis space [9, 1]. Such bounds have the appealing feature that they provide performance guarantees for classifiers found by *any* learning algorithm. However, it has been observed that these bounds tend to be overly pessimistic. One explanation is that only in the case of learning algorithms which minimise the training error it has been proven that uniformity of the bounds is equivalent to studying the learning algorithm's generalisation performance directly.

In this paper we present a theoretical framework which aims at *directly* studying the generalisation error of a learning algorithm rather than taking the detour via the uniform convergence of training errors to expected errors in a given hypothesis space. In addition, our new model of learning allows the exploitation of the fact that we serendipitously observe a training sample which is easy to learn by a given learning algorithm. In that sense, our framework is a descendant of the luckiness framework of Shawe-Taylor et al. [8]. In the present case, the luckiness is a function of a given learning algorithm and a given training sample and characterises the diversity of the algorithms solutions. The notion of luckiness allows us to study given learning algorithms at many different perspectives. For example, the maximum margin algorithm [9] can either been studied via the number of dimensions in feature space,

the margin of the classifier learned or the sparsity of the resulting classifier. Our main results are two generalisation error bounds for learning algorithms: one for the zero training error scenario and one agnostic bound (Section 2). We shall demonstrate the usefulness of our new framework by studying its relation to the VC framework, the original luckiness framework and the compression framework of Littlestone and Warmuth [6] (Section 3). Finally, we present an application of the new framework to the maximum margin algorithm for linear classifiers (Section 4). The detailed proofs of our main results can be found in [5].

We denote vectors using bold face, e.g. $\mathbf{x} = (x_1, \dots, x_m)$ and the length of this vector by $|\mathbf{x}|$, i.e. $|\mathbf{x}| = m$. In order to unburden notation we use the shorthand notation $\mathbf{z}_{[i:j]} := (z_i, \dots, z_j)$ for $i \leq j$. Random variables are typeset in sans-serif font. The symbols \mathbf{P}_X , $\mathbf{E}_X[f(X)]$ and \mathbb{I} denote a probability measure over X , the expectation of $f(\cdot)$ over the random draw of its argument x and the indicator function, respectively. The shorthand notation $\mathcal{Z}^{(\infty)} := \cup_{m=1}^{\infty} \mathcal{Z}^m$ denotes the union of all m -fold Cartesian products of the set \mathcal{Z} . For any $m \in \mathbb{N}$ we define $I_m \subset \{1, \dots, m\}^m$ as the set of all permutations of the numbers $1, \dots, m$,

$$I_m := \{(i_1, \dots, i_m) \in \{1, \dots, m\}^m \mid \forall j \neq k : i_j \neq i_k\}.$$

Given a $2m$ -vector $\mathbf{i} \in I_{2m}$ and a sample $\mathbf{z} \in \mathcal{Z}^{2m}$ we define $\pi_{\mathbf{i}} : \{1, \dots, 2m\} \rightarrow \{1, \dots, 2m\}$ by $\pi_{\mathbf{i}}(j) := i_j$ and $\Pi_{\mathbf{i}}(\mathbf{z})$ by $\Pi_{\mathbf{i}}(\mathbf{z}) := (z_{\pi_{\mathbf{i}}(1)}, \dots, z_{\pi_{\mathbf{i}}(2m)})$.

2 Algorithmic Luckiness

Suppose we are given a training sample $\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in (\mathcal{X} \times \mathcal{Y})^m = \mathcal{Z}^m$ of size $m \in \mathbb{N}$ independently drawn (iid) from some unknown but fixed distribution $\mathbf{P}_{XY} = \mathbf{P}_Z$ together with a learning algorithm $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$. For a predefined loss $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ we would like to investigate the generalisation error $G_l[\mathcal{A}, \mathbf{z}] := R_l[\mathcal{A}(\mathbf{z})] - \inf_{h \in \mathcal{Y}^{\mathcal{X}}} R_l[h]$ of the algorithm where the *expected error* $R_l[h]$ of h is defined by

$$R_l[h] := \mathbf{E}_{XY}[l(h(X), Y)].$$

Since $\inf_{h \in \mathcal{Y}^{\mathcal{X}}} R_l[h]$ (which is also known as the *Bayes error*) is independent of \mathcal{A} it suffices to bound $R_l[\mathcal{A}(\mathbf{z})]$. Although we know that for any fixed hypothesis h the *training error*

$$\widehat{R}_l[h, \mathbf{z}] := \frac{1}{|\mathbf{z}|} \sum_{(x_i, y_i) \in \mathbf{z}} l(h(x_i), y_i)$$

is with high probability (over the random draw of the training sample $\mathbf{z} \in \mathcal{Z}^{(\infty)}$) close to $R_l[h]$, this might no longer be true for the *random* hypothesis $\mathcal{A}(\mathbf{z})$. Hence we would like to state that with only small probability (at most δ), the expected error $R_l[\mathcal{A}(\mathbf{z})]$ is larger than the training error $\widehat{R}_l[\mathcal{A}(\mathbf{z}), \mathbf{z}]$ plus some sample *and* algorithm dependent complexity $\varepsilon(\mathcal{A}, \mathbf{z}, \delta)$,

$$\mathbf{P}_{Z^m} \left(R_l[\mathcal{A}(\mathbf{Z})] > \widehat{R}_l[\mathcal{A}(\mathbf{Z}), \mathbf{Z}] + \varepsilon(\mathcal{A}, \mathbf{Z}, \delta) \right) < \delta. \quad (1)$$

In order to derive such a bound we utilise a modified version of the basic lemma of Vapnik and Chervonenkis [10].

Lemma 1. *For all loss functions $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$, all probability measures \mathbf{P}_Z , all algorithms \mathcal{A} and all measurable formulas $\Upsilon : \mathcal{Z}^m \rightarrow \{\text{true}, \text{false}\}$, if $m\varepsilon^2 > 2$ then*

$$\mathbf{P}_{Z^m} \left(\left(R_l[\mathcal{A}(\mathbf{Z})] > \widehat{R}_l[\mathcal{A}(\mathbf{Z}), \mathbf{Z}] + \varepsilon \right) \wedge \Upsilon(\mathbf{Z}) \right) < 2\mathbf{P}_{Z^{2m}} \left(\underbrace{\left(\widehat{R}_l[\mathcal{A}(\mathbf{Z}_{[1:m]}), \mathbf{Z}_{[(m+1):2m]}] > \widehat{R}_l[\mathcal{A}(\mathbf{Z}_{[1:m]}), \mathbf{Z}_{[1:m]}] + \frac{\varepsilon}{2} \right)}_{J(\mathbf{Z})} \wedge \Upsilon(\mathbf{Z}_{[1:m]}) \right).$$

Proof (Sketch). The probability on the r.h.s. is lower bounded by the probability of the conjunction of event on the l.h.s. and $Q(\mathbf{z}) \equiv R_l[\mathcal{A}(\mathbf{z}_{[1:m]})] - \widehat{R}_l[\mathcal{A}(\mathbf{z}_{[1:m]}), \mathbf{z}_{(m+1):2m}] < \frac{\varepsilon}{2}$. Note that this probability is over $\mathbf{z} \in \mathcal{Z}^{2m}$. If we now condition on the first m examples, $\mathcal{A}(\mathbf{z}_{[1:m]})$ is fixed and therefore by an application of Hoeffding's inequality (see, e.g. [1]) and since $m\varepsilon^2 > 2$ the additional event Q has probability of at least $\frac{1}{2}$ over the random draw of (z_{m+1}, \dots, z_{2m}) . \square

Use of Lemma 1 — which is similar to the approach of classical VC analysis — reduces the original problem (1) to the problem of studying the deviation of the training errors on the first and second half of a double sample $\mathbf{z} \in \mathcal{Z}^{2m}$ of size $2m$. It is of utmost importance that the hypothesis $\mathcal{A}(\mathbf{z}_{[1:m]})$ is always learned from the first m examples. Now, in order to fully exploit our assumptions of the mutual independence of the double sample $\mathbf{z} \in \mathcal{Z}^{2m}$ we use a technique known as symmetrisation by permutation: since $\mathbf{P}_{\mathcal{Z}^{2m}}$ is a product measure, it has the property that $\mathbf{P}_{\mathcal{Z}^{2m}}(J(\mathbf{Z})) = \mathbf{P}_{\mathcal{Z}^{2m}}(J(\Pi_i(\mathbf{Z})))$ for any $i \in I_{2m}$. Hence, it suffices to bound the probability of permutations π_i such that $J(\Pi_i(\mathbf{z}))$ is true for a given and *fixed* double sample \mathbf{z} . As a consequence thereof, we only need to count the number of different hypotheses that can be learned by \mathcal{A} from the first m examples when permuting the double sample.

Definition 1 (Algorithmic luckiness). Any function L that maps an algorithm $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$ and a training sample $\mathbf{z} \in \mathcal{Z}^{(\infty)}$ to a real value is called an *algorithmic luckiness*. For all $m \in \mathbb{N}$, for any $\mathbf{z} \in \mathcal{Z}^{2m}$, the *lucky set* $\mathcal{H}_{\mathcal{A}}(L, \mathbf{z}) \subseteq \mathcal{Y}^{\mathcal{X}}$ is the set of all hypotheses that are learned from the first m examples $(z_{\pi_i(1)}, \dots, z_{\pi_i(m)})$ when permuting the whole sample \mathbf{z} whilst not decreasing the luckiness, i.e.

$$\mathcal{H}_{\mathcal{A}}(L, \mathbf{z}) := \{ \mathcal{A}(z_{\pi_i(1)}, \dots, z_{\pi_i(m)}) \mid i \in \mathcal{I}_{\mathcal{A}}(L, \mathbf{z}) \}, \quad (2)$$

where

$$\mathcal{I}_{\mathcal{A}}(L, \mathbf{z}) := \{ i \in I_{2m} \mid L(\mathcal{A}, (z_{\pi_i(1)}, \dots, z_{\pi_i(m)})) \geq L(\mathcal{A}, (z_1, \dots, z_m)) \}. \quad (3)$$

Given a fixed loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ the *induced loss function set* $\mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z}))$ is defined by

$$\mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z})) := \{(x, y) \mapsto l(h(x), y) \mid h \in \mathcal{H}_{\mathcal{A}}(L, \mathbf{z})\}.$$

For any luckiness function L and any learning algorithm \mathcal{A} , the complexity of the double sample \mathbf{z} is the minimal number $\mathcal{N}_1(\tau, \mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z})), \mathbf{z})$ of hypotheses $\hat{h} \in \mathcal{Y}^{\mathcal{X}}$ needed to cover $\mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z}))$ at some predefined scale τ , i.e. for any hypothesis $h \in \mathcal{H}_{\mathcal{A}}(L, \mathbf{z})$ there exists a $\hat{h} \in \mathcal{Y}^{\mathcal{X}}$ such that

$$\frac{1}{2m} \sum_{i=1}^{2m} |l(h(x_i), y_i) - l(\hat{h}(x_i), y_i)| \leq \tau. \quad (4)$$

To see this note that whenever $J(\Pi_i(\mathbf{z}))$ is true (over the random draw of permutations) then there exists a function \hat{h} which has a difference in the training errors on the double sample of at least $\frac{\varepsilon}{2} + 2\tau$. By an application of the union bound we see that the number $\mathcal{N}_1(\tau, \mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z})), \mathbf{z})$ is of central importance. Hence, if we are able to bound this number over the random draw of the double sample \mathbf{z} only using the luckiness on the first m examples we can use this bound in place of the worst case complexity $\sup_{\mathbf{z} \in \mathcal{Z}^{2m}} \mathcal{N}_1(\tau, \mathcal{L}_l(\mathcal{H}_{\mathcal{A}}(L, \mathbf{z})), \mathbf{z})$ as usually done in the VC framework (see [9]).

Definition 2 (ω -smallness of L). Given an algorithm $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$ and a loss $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ the algorithmic luckiness function L is ω -small at scale $\tau \in \mathbb{R}^+$ if for all $m \in \mathbb{N}$, all $\delta \in (0, 1]$ and all $\mathbf{P}_{\mathbf{Z}}$

$$\mathbf{P}_{\mathbf{Z}^{2m}} \left(\underbrace{\mathcal{N}_1 \left(\tau, \mathcal{L}_l \left(\mathcal{H}_{\mathcal{A}} \left(L, \mathbf{Z} \right), \mathbf{z} \right) \right)}_{S(\mathbf{Z})} > \omega \left(L \left(\mathcal{A}, \mathbf{Z}_{[1:m]} \right), l, m, \delta, \tau \right) \right) < \delta.$$

Note that if the range of l is $\{0, 1\}$ then $\mathcal{N}_1 \left(\frac{1}{2m}, \mathcal{L}_l \left(\mathcal{H}_{\mathcal{A}} \left(L, \mathbf{z} \right), \mathbf{z} \right) \right)$ equals the number of dichotomies on \mathbf{z} incurred by $\mathcal{L}_l \left(\mathcal{H}_{\mathcal{A}} \left(L, \mathbf{z} \right) \right)$.

Theorem 1 (Algorithmic luckiness bounds). *Suppose we have a learning algorithm $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$ and an algorithmic luckiness L that is ω -small at scale τ for a loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$. For any probability measure $\mathbf{P}_{\mathbf{Z}}$, any $d \in \mathbb{N}$ and any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random draw of the training sample $\mathbf{z} \in \mathcal{Z}^m$ of size m , if $\omega \left(L \left(\mathcal{A}, \mathbf{z} \right), l, m, \delta/4, \tau \right) \leq 2^d$ then*

$$R_l [\mathcal{A} (\mathbf{z})] \leq \widehat{R}_l [\mathcal{A} (\mathbf{z}), \mathbf{z}] + \sqrt{\frac{8}{m} \left(d + \log_2 \left(\frac{4}{\delta} \right) \right)} + 4\tau. \quad (5)$$

Furthermore, under the above conditions if the algorithmic luckiness L is ω -small at scale $\frac{1}{2m}$ for a binary loss function $l (\cdot, \cdot) \in \{0, 1\}$ and $\widehat{R}_l [\mathcal{A} (\mathbf{z}), \mathbf{z}] = 0$ then

$$R_l [\mathcal{A} (\mathbf{z})] \leq \frac{2}{m} \left(d + \log_2 \left(\frac{4}{\delta} \right) \right). \quad (6)$$

Proof (Compressed Sketch). We will only sketch the proof of equation (5); the proof of (6) is similar and can be found in [5]. First, we apply Lemma 1 with $\Upsilon (\mathbf{z}) \equiv \omega \left(L \left(\mathcal{A}, \mathbf{z} \right), l, m, \delta/4, \tau \right) \leq 2^d$. We now exploit the fact that

$$\begin{aligned} \mathbf{P}_{\mathbf{Z}^{2m}} (J (\mathbf{Z})) &= \underbrace{\mathbf{P}_{\mathbf{Z}^{2m}} (J (\mathbf{Z}) \wedge S (\mathbf{Z}))}_{\leq \mathbf{P}_{\mathbf{Z}^{2m}} (S (\mathbf{Z}))} + \mathbf{P}_{\mathbf{Z}^{2m}} (J (\mathbf{Z}) \wedge \neg S (\mathbf{Z})) \\ &< \frac{\delta}{4} + \mathbf{P}_{\mathbf{Z}^{2m}} (J (\mathbf{Z}) \wedge \neg S (\mathbf{Z})), \end{aligned}$$

which follows from Definition 2. Following the above-mentioned argument it suffices to bound the probability of a random permutation $\Pi_1 (\mathbf{z})$ that $J (\Pi_1 (\mathbf{z})) \wedge \neg S (\Pi_1 (\mathbf{z}))$ is true for a *fixed* double sample \mathbf{z} . Noticing that $\Upsilon (\mathbf{z}) \wedge \neg S (\mathbf{z}) \Rightarrow \mathcal{N}_1 \left(\tau, \mathcal{L}_l \left(\mathcal{H}_{\mathcal{A}} \left(L, \mathbf{z} \right), \mathbf{z} \right) \right) \leq 2^d$ we see that we only consider swappings π_i for which $\mathcal{N}_1 \left(\tau, \mathcal{L}_l \left(\mathcal{H}_{\mathcal{A}} \left(L, \Pi_i (\mathbf{z}) \right), \Pi_i (\mathbf{z}) \right) \right) \leq 2^d$. Thus let us consider such a cover of size not more than 2^d . By (4) we know that whenever $J (\Pi_i (\mathbf{z})) \wedge \neg S (\Pi_i (\mathbf{z}))$ is true for a swapping i then there exists a hypothesis $\hat{h} \in \mathcal{Y}^{\mathcal{X}}$ in the cover such that $\widehat{R}_l \left[\hat{h}, (\Pi_1 (\mathbf{z}))_{[(m+1):2m]} \right] - \widehat{R}_l \left[\hat{h}, (\Pi_1 (\mathbf{z}))_{[1:m]} \right] > \frac{\epsilon}{2} + 2\tau$. Using the union bound and Hoeffding's inequality for a particular choice of \mathbf{P}_1 shows that $\mathbf{P}_1 (J (\Pi_1 (\mathbf{z})) \wedge \neg S (\Pi_1 (\mathbf{z}))) \leq \frac{\delta}{4}$ which finalises the proof. \square

A closer look at (5) and (6) reveals that the essential difference to uniform bounds on the expected error is within the definition of the covering number: rather than covering all hypotheses h in a given hypothesis space $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ for a given double sample it suffices to cover all hypotheses that can be learned by a given learning algorithm from the first half when permuting the double sample. Note that the usage of permutations in the definition of (2) is not only a technical matter; it fully exploits all the assumptions made for the training sample, namely the training sample is drawn iid.

3 Relationship to Other Learning Frameworks

In this section we present the relationship of algorithmic luckiness to other learning frameworks (see [9, 8, 6] for further details of these frameworks).

VC Framework If we consider a binary loss function $l(\cdot, \cdot) \in \{0, 1\}$ and assume that the algorithm \mathcal{A} selects functions from a given hypothesis space $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ then $L(\mathcal{A}, \mathbf{z}) = -\text{VCDim}(\mathcal{H})$ is a ω -small luckiness function where

$$\omega\left(L_0, l, m, \delta, \frac{1}{2m}\right) \leq \left(\frac{2em}{-L_0}\right)^{-L_0}. \quad (7)$$

This can easily be seen by noticing that the latter term is an upper bound on $\max_{\mathbf{z} \in \mathcal{Z}^{2m}} |\{(l(h(x_1), y_1), \dots, l(h(x_{2m}), y_{2m})) : h \in \mathcal{H}\}|$ (see also [9]). Note that this luckiness function neither exploits the particular training sample observed nor the learning algorithm used.

Luckiness Framework Firstly, the luckiness framework of Shawe-Taylor et al. [8] only considered binary loss functions l and the zero training error case. In this work, the luckiness \tilde{L} is a function of hypothesis and training samples and is called $\tilde{\omega}$ -small if the probability over the random draw of a $2m$ sample \mathbf{z} that there exists a hypothesis h with $\tilde{\omega}(\tilde{L}(h, (z_1, \dots, z_m)), \delta) < \mathcal{N}_1(\frac{1}{2m}, \{(x, y) \mapsto l(g(x), y) \mid \tilde{L}(g, \mathbf{z}) \geq \tilde{L}(h, \mathbf{z})\}, \mathbf{z})$, is smaller than δ . Although similar in spirit, the classical luckiness framework does not allow exploitation of the learning algorithm used to the same extent as our new luckiness. In fact, in this framework not only the covering number must be estimable but also the variation of the luckiness \tilde{L} itself. These differences make it very difficult to formally relate the two frameworks.

Compression Framework In the compression framework of Littlestone and Warmuth [6] one considers learning algorithms \mathcal{A} which are compression schemes, i.e. $\mathcal{A}(\mathbf{z}) = \mathcal{R}(\mathcal{C}(\mathbf{z}))$ where $\mathcal{C}(\mathbf{z})$ selects a subsample $\bar{\mathbf{z}} \subseteq \mathbf{z}$ and $\mathcal{R} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$ is a permutation invariant reconstruction function. For this class of learning algorithms, the luckiness $L(\mathcal{A}, \mathbf{z}) = -|\mathcal{C}(\mathbf{z})|$ is ω -small where ω is given by (7). In order to see this we note that (3) ensures that we only consider permutations π_i where $\mathcal{C}(\Pi_i(\mathbf{z})) \leq |\mathcal{C}(\mathbf{z})|$, i.e. we use not more than $-L$ training examples from $\mathbf{z} \in \mathcal{Z}^{2m}$. As there are exactly $\binom{2m}{d}$ distinct choices of d training examples from $2m$ examples the result follows by application of Sauer's lemma [9]. Disregarding constants, Theorem 1 gives exactly the same bound as in [6].

4 A New Margin Bound For Support Vector Machines

In this section we study the maximum margin algorithm for linear classifiers, i.e. $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{H}_\phi$ where $\mathcal{H}_\phi := \{x \mapsto \langle \phi(x), \mathbf{w} \rangle \mid \mathbf{w} \in \mathcal{K}\}$ and $\phi : \mathcal{X} \rightarrow \mathcal{K} \subseteq \ell_2^n$ is known as the feature mapping. Let us assume that $l(h(x), y) = l_{0-1}(h(x), y) := \mathbb{I}_{yh(x) \leq 0}$. Classical VC generalisation error bounds exploit the fact that $\text{VCDim}(\mathcal{H}_\phi) = n$ and (7). In the luckiness framework of Shawe-Taylor et al. [8] it has been shown that we can use $\text{fat}_{\mathcal{H}_\phi}(\gamma_{\mathbf{z}}(\mathbf{w})) \leq (\gamma_{\mathbf{z}}(\mathbf{w}))^{-2}$ (at the price of an extra $\log_2(32m)$ factor) in place of $\text{VCDim}(\mathcal{H}_\phi)$ where $\gamma_{\mathbf{z}}(\mathbf{w}) = \min_{(x_i, y_i) \in \mathbf{z}} y_i \langle \phi(x_i), \mathbf{w} \rangle / \|\mathbf{w}\|$ is known as the margin. Now, the maximum margin algorithm finds the weight vector \mathbf{w}_{MM} that maximises $\gamma_{\mathbf{z}}(\mathbf{w})$. It is known that \mathbf{w}_{MM} can be written as a linear combination of the $\phi(x_i)$. For notational convenience, we shall assume that $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathbb{R}^{(\infty)}$ maps

to the expansion coefficients α such that $\|\mathbf{w}_\alpha\| = 1$ where $\mathbf{w}_\alpha := \sum_{i=1}^{|z|} \alpha_i \phi(x_i)$. Our new margin bound follows from the following theorem together with (6).

Theorem 2. Let $\epsilon_i(\mathbf{x})$ be the smallest $\epsilon > 0$ such that $\{\phi(x_1), \dots, \phi(x_m)\}$ can be covered by at most i balls of radius less than or equal ϵ . Let $\Gamma_{\mathbf{z}}(\mathbf{w})$ be defined by $\Gamma_{\mathbf{z}}(\mathbf{w}) := \min_{(x_i, y_i) \in \mathbf{z}} \frac{y_i \langle \phi(x_i), \mathbf{w} \rangle}{\|\phi(x_i)\| \cdot \|\mathbf{w}\|}$. For the zero-one loss l_{0-1} and the maximum margin algorithm \mathcal{A} , the luckiness function

$$L(\mathcal{A}, \mathbf{z}) = - \min \left\{ i \in \mathbb{N} \mid i \geq \left(\frac{\epsilon_i(\mathbf{x}) \sum_{j=1}^m |\mathcal{A}(\mathbf{z})_j|}{\Gamma_{\mathbf{z}}(\mathbf{w}_{\mathcal{A}(\mathbf{z})})} \right)^2 \right\}, \quad (8)$$

is ω -small at scale $1/2m$ w.r.t. the function

$$\omega \left(L_0, l, m, \delta, \frac{1}{2m} \right) = \left(\frac{2em}{-L_0} \right)^{-2L_0}. \quad (9)$$

Proof (Sketch). First we note that by a slight refinement of a theorem of Makovoz [7] we know that for any $\mathbf{z} \in \mathcal{Z}^m$ there exists a weight vector $\tilde{\mathbf{w}} = \sum_{i=1}^m \tilde{\alpha}_i \phi(x_i)$ such that

$$\|\tilde{\mathbf{w}} - \mathbf{w}_{\mathcal{A}(\mathbf{z})}\|^2 \leq \Gamma_{\mathbf{z}}^2(\mathbf{w}_{\mathcal{A}(\mathbf{z})}) \quad (10)$$

and $\tilde{\alpha} \in \mathbb{R}^m$ has no more than $-L(\mathcal{A}, \mathbf{z})$ non-zero components. Although only $\mathbf{w}_{\mathcal{A}(\mathbf{z})}$ is of unit length, one can show that (10) implies that

$$\langle \mathbf{w}_{\mathcal{A}(\mathbf{z})}, \tilde{\mathbf{w}} / \|\tilde{\mathbf{w}}\| \rangle \geq \sqrt{1 - \Gamma_{\mathbf{z}}^2(\mathbf{w}_{\mathcal{A}(\mathbf{z})})}.$$

Using equation (10) of [4] this implies that $\tilde{\mathbf{w}}$ correctly classifies $\mathbf{z} \in \mathcal{Z}^m$. Consider a fixed double sample $\mathbf{z} \in \mathcal{Z}^{2m}$ and let $k_0 := L(\mathcal{A}, (z_1, \dots, z_m))$. By virtue of (3) and the aforementioned argument we only need to consider permutations π_i such that there exists a weight vector $\tilde{\mathbf{w}} = \sum_{j=1}^m \tilde{\alpha}_j \phi(x_j)$ with no more than k_0 non-zero $\tilde{\alpha}_j$. As there are exactly $\binom{2m}{d}$ distinct choices of $d \in \{1, \dots, k_0\}$ training examples from the $2m$ examples \mathbf{z} there are no more than $(2em/k_0)^{k_0}$ different subsamples to be used in $\tilde{\mathbf{w}}$. For each particular subsample $\bar{\mathbf{z}} \subseteq \mathbf{z}$ the weight vector $\tilde{\mathbf{w}}$ is a member of the class of linear classifiers in a k_0 (or less) dimensional space. Thus, from (7) it follows that for the given subsample $\bar{\mathbf{z}}$ there are no more $(2em/k_0)^{k_0}$ different dichotomies induced on the double sample $\mathbf{z} \in \mathcal{Z}^{2m}$. As this holds for any double sample, the theorem is proven. \square

There are several interesting features about this margin bound. Firstly, observe that $\sum_{j=1}^m |\mathcal{A}(\mathbf{z})_j|$ is a measure of sparsity of the solution found by the maximum margin algorithm which, in the present case, is combined with margin. Note that for normalised data, i.e. $\|\phi(\cdot)\| = \text{constant}$, the two notion of margins coincide, i.e. $\Gamma_{\mathbf{z}}(\mathbf{w}) = \gamma_{\mathbf{z}}(\mathbf{w})$. Secondly, the quantity $\epsilon_i(\mathbf{x})$ can be considered as a measure of the distribution of the mapped data points in feature space. Note that for all $i \in \mathbb{N}$, $\epsilon_i(\mathbf{x}) \leq \epsilon_1(\mathbf{x}) \leq \max_{j \in \{1, \dots, m\}} \|\phi(x_j)\|$. Supposing that the two class-conditional probabilities $\mathbf{P}_{\mathbf{X}|\mathbf{Y}=y}$ are highly clustered, $\epsilon_2(\mathbf{x})$ will be very small. An extension of this reasoning is useful in the multi-class case; binary maximum margin classifiers are often used to solve multi-class problems [9]. There appears to be also a close relationship of $\epsilon_i(\mathbf{x})$ with the notion of kernel alignment recently introduced in [3]. Finally, one can use standard entropy number techniques to bound $\epsilon_i(\mathbf{x})$ in terms of eigenvalues of the inner product matrix or its centred variants. It is worth mentioning that although our aim was to study the maximum margin algorithm the

above theorem actually holds for any algorithm whose solution can be represented as a linear combination of the data points.

5 Conclusions

In this paper we have introduced a new theoretical framework to study the generalisation error of learning algorithms. In contrast to previous approaches, we considered specific learning algorithms rather than specific hypothesis spaces. We introduced the notion of algorithmic luckiness which allowed us to devise data dependent generalisation error bounds. Thus we were able to relate the compression framework of Littlestone and Warmuth with the VC framework. Furthermore, we presented a new bound for the maximum margin algorithm which not only exploits the margin but also the distribution of the *actual* training data in feature space. Perhaps the most appealing feature of our margin based bound is that it naturally combines the three factors considered important for generalisation with linear classifiers: margin, sparsity and the distribution of the data. Further research is concentrated on studying Bayesian algorithms and the relation of algorithmic luckiness to the recent findings for stable learning algorithms [2].

Acknowledgements This work was done while RCW was visiting Microsoft Research Cambridge. This work was also partly supported by the Australian Research Council. RH would like to thank Olivier Bousquet for stimulating discussions.

References

- [1] M. Anthony and P. Bartlett. *A Theory of Learning in Artificial Neural Networks*. Cambridge University Press, 1999.
- [2] O. Bousquet and A. Elisseeff. Algorithmic stability and generalization performance. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*, pages 196–202. MIT Press, 2001.
- [3] N. Cristianini, A. Elisseeff, and J. Shawe-Taylor. On optimizing kernel alignment. Technical Report NC2-TR-2001-087, NeuroCOLT, <http://www.neurocolt.com>, 2001.
- [4] R. Herbrich and T. Graepel. A PAC-Bayesian margin bound for linear classifiers: Why SVMs work. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*, pages 224–230, Cambridge, MA, 2001. MIT Press.
- [5] R. Herbrich and R. C. Williamson. Algorithmic luckiness. Technical report, Microsoft Research, 2002.
- [6] N. Littlestone and M. Warmuth. Relating data compression and learnability. Technical report, University of California Santa Cruz, 1986.
- [7] Y. Makovoz. Random approximants and neural networks. *Journal of Approximation Theory*, 85:98–109, 1996.
- [8] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.
- [9] V. Vapnik. *Statistical Learning Theory*. John Wiley and Sons, New York, 1998.
- [10] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–281, 1971.