

---

# Inductive Principles

**Robert C. Williamson\***

Australian National University  
Canberra, 0200 ACT, Australia



(Bowdlerised Edition)

\* Joint work with **Ralf Herbrich**, Microsoft Research Cambridge

---



- Induction (what’s “inductive”?)
- Inductive Principles (what’s the “principle”?)
- Empirical Risk Minimization
- Key Theorem of Learning Theory
- Conditioning on the Data

After the break, we will move on to the more technical part of the talk . . .

- “Conditioning on the data” in a Frequentist (PAC) setting — The Luckiness Framework
- A new approach — Algorithmic Luckiness

# Induction

---

Hume's problem is how to justify Induction: the inference (discovery of laws) from empirical data.

If we take in our hand any volume - of divinity or school metaphysics, for instance - let us ask. Does it contain any abstract reasoning concerning quantity or number? No. Does it contain any experimental reasoning containing matter of fact and existence? No. Commit it then to the flames, for it can contain nothing but sophistry and illusion.

Impossible “for all is but a woven web of guesses”.

— David Hume

Popper's key insight: scientific theories do not lead to certain knowledge; merely approximations to the truth. Thus no “justification”



We can, however, reason logically about the process of scientific discovery. Doing so shows one should prefer a more refutable theory over a less refutable one.

# Induction



[W]e can always construct our machine so that it starts issuing probabilistic predictions only after the 1000th event, say, or after any other number  $n$  which we may choose, bearing in mind the

number of different hypotheses in our 'world'. (The problem is so trivial that it is not worth making any effort to solve it systematically; for we know, after all, that applications of the

simple inductive rule will never give us more than increasingly good approximations).

Popper argued that we could formally consider the "dimension" of a theory William Popper, *Objective Knowledge*, Oxford University Press, 1972, p321 (1956,1983)

Although it is impossible to build a general learning machine Popper clearly admitted the possibility *within a constrained framework* of building a machine and being able to probabilistically reason about its performance.

We will study learning machines, and not induction in general.



# Learning Problem

---

The players ... threw these abstract formulas at one another displaying the sequences and possibilities of their science.

— Herman Hesse: *The Glass Bead Game*

## Given:

- A **training sample**  $\mathbf{z} = (\mathbf{x}, \mathbf{y}) = (z_1, \dots, z_m) \in (\mathcal{X} \times \mathcal{Y})^m = \mathcal{Z}^m$  drawn iid from  $\mathbf{P}_{\mathcal{Z}}$  (unknown).
- A deterministic **learning algorithm**  $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$ .
- A **loss function**  $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ .



**Question:** How can one tell whether  $\mathcal{A}$  is good or not?

# One Possible Solution

---



Find a (probabilistic) **bound** on  $R_l[\mathcal{A}(z)] := \mathbf{E}_{XY}[l(\mathcal{A}(z)(X), Y)]$ , that is, a function  $\psi$  such that

$$\mathbf{P}_{Z^m}(R_l[\mathcal{A}(Z)] \leq \psi(\mathcal{A}, Z, \delta)) \geq 1 - \delta.$$

# Why do we want such bounds?

A bound such as

$$\mathbf{P}_{\mathcal{Z}^m} (R_l[\mathcal{A}(\mathbf{Z})] \leq \psi(\mathcal{A}, \mathbf{Z}, \delta)) \geq 1 - \delta$$

is not an end in itself; it suggests how to adjust the parameters (or knobs) of the learning algorithm  $\mathcal{A}$ .



Thus the *closer* the analysis is to the algorithm the more insightful we would hope it to be.

**Difficulty:** Given  $\mathbf{z} \in \mathcal{Z}^m$  how can the algorithm  $\mathcal{A}$  choose an hypothesis that achieves a small value of  $R_l[\mathcal{A}(\mathbf{z})]$ ?

Key point: **Given**  $\mathbf{z} \in \mathcal{Z}^m$ . No chance of computing  $R_l[\mathcal{A}(\mathbf{z})]$  even in principle because we do not know  $\mathbf{P}_{\mathcal{Z}}$ .

# A Recipe for Generalisation Error Bounds

1. Relate the prediction error  $R_l [h]$  to some **empirical quantity**, e.g. training error

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



that converges exponentially to  $R_l [h]$  for any  $h$ .

2. Apply the **basic lemma** to the difference of the prediction error and the empirical quantity (training error). This introduces a **ghost sample**.
3. Fully exploit the independence assumption of  $\mathbf{z}$  by using a technique known as **symmetrisation by permutation**: (probability is over **double** sample  $\mathbf{Z}^{2m}$ )

$$\mathbf{P}_{\mathbf{Z}^{2m}} (\Upsilon (\mathbf{Z})) = \mathbf{E}_{\mathbf{I}} \left[ \mathbf{P}_{\mathbf{Z}^{2m} | \mathbf{I} = \mathbf{i}} (\Pi_{\mathbf{i}} (\Upsilon (\mathbf{Z}))) \right] = \mathbf{E}_{\mathbf{Z}^{2m}} \left[ \mathbf{P}_{\mathbf{I} | \mathbf{Z}^{2m} = \mathbf{z}} (\Pi_{\mathbf{I}} (\Upsilon (\mathbf{z}))) \right].$$

4. Since  $\mathbf{z} \in \mathcal{Z}^{2m}$  is fixed, we can construct a **cover** w.r.t. the loss  $l$  and apply the **union bound**.



# Inductive Principle — Empirical Risk Minimization



Possibility of computing such bounds motivates:  
**E**mpirical **R**isk **M**inimization Algorithm

This is a great algorithm to analyse.

— Ralf Herbrich



$$\mathcal{A}_{\text{erm}}^{\mathcal{H}} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$$

$$\mathcal{A}_{\text{erm}}^{\mathcal{H}} : \mathbf{z} \mapsto \arg \min_{h \in \mathcal{H}} \widehat{R}_l [h, \mathbf{z}].$$

The “principle” is to minimize the empirical surrogate

$\widehat{R}_l [h, \mathbf{z}]$  of  $R_l [h, \mathbf{z}]$ . Thus, depending on the choice of sets  $\mathcal{H}$ , get a family of *empirical risk minimization algorithms*.

Note that the algorithm  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  has one “knob”: the class of functions  $\mathcal{H}$ .

How to choose  $\mathcal{H}$ ? Want  $\mathcal{H}$  as large as possible to ensure a good approximation of the underlying data generating process. Pay a price . . .



# Consistency and Strict Consistency

---

**Question:** Does  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  “work”?

**More precise question:** Is  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  consistent?

Assume that  $\mathcal{H}$  and  $l$  are such that for any  $\mathbf{z} \in \mathcal{Z}^m$

$$\widehat{R}_l \left[ \mathcal{A}_{\text{erm}}^{\mathcal{H}}(\mathbf{z}), \mathbf{z} \right] = \inf_{h \in \mathcal{H}} \widehat{R}_l [h(\mathbf{z}), \mathbf{z}]$$

and that for all  $h \in \mathcal{H}$ ,  $A \leq R_l[h] \leq B$ . Let

$$\mathcal{H}(c) := \{h \in \mathcal{H} : R_l[h] \geq c\}.$$

Say that  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  is *strictly (nontrivially) consistent* if for all  $c \geq 0$ , for all  $\varepsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \mathbf{P}_{\mathcal{Z}^m} \left( \left| \widehat{R}_l \left[ \mathcal{A}_{\text{emp}}^{\mathcal{H}(c)}(\mathbf{z}), \mathbf{z} \right] - c \right| > \varepsilon \right) = 0.$$

Need a definition like this to rule out “coding” the identity of a function into one observation: can construct such artificial function classes of arbitrary complexity which can be learned using  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  with only one observation.

---



# Consistency of ERM

$\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  is strictly consistent

$\Leftrightarrow$

$$\forall \varepsilon > 0 \quad \lim_{m \rightarrow \infty} \mathbf{P}_{\mathbf{Z}^m} \left\{ \sup_{h \in \mathcal{H}} (R[h] - \widehat{R}_l[h, \mathbf{z}]) > \varepsilon \right\} = 0 \quad \star$$

$$\forall \varepsilon > 0 \quad \lim_{m \rightarrow \infty} \mathbf{P}_{\mathbf{Z}^m} \left\{ \sup_{h \in \mathcal{H}} |R_l[h] - \widehat{R}_l[h, \mathbf{z}]| > \varepsilon \right\} = 0 \quad \star\star$$

$\Leftrightarrow$

$$\forall \varepsilon > 0 \quad \lim_{m \rightarrow \infty} \frac{1}{m} \mathbf{E}_{\mathbf{Z}^m} \log \underbrace{\mathcal{N}(\varepsilon, \mathcal{H}, \ell_1(\mathbf{z}))}_{\text{Covering number of } \mathcal{H} \text{ at scale } \varepsilon} = 0$$

Covering number of  $\mathcal{H}$  at scale  $\varepsilon$   
w.r.t. to the  $\ell_1(\mathbf{z})$  metric: for  $h \in \mathcal{H}$ ,  $\|h\|_{\ell_1(\mathbf{z})} := \frac{1}{m} \sum_{z \in \mathbf{z}} |h(z)|$ .

The effective gap in the reasoning implicit in the difference between  $\star$  and  $\star\star$  can be plugged using a more complex notion of cover — a one sided bracket cover. I am unaware of any results on the relative sizes of such covering numbers compared to  $\mathcal{N}(\varepsilon, \mathcal{H}, \ell_1(\mathbf{z}))$ .



## So What?

---

The big deal is that (modulo the small gap mentioned)

$$\mathbf{E}_{\mathbf{Z}^m} \log \mathcal{N}(\varepsilon, \mathcal{H}, \ell_1(\mathbf{z}))$$

is the “right” quantity to study for understanding the effect of the  $\mathcal{H}$  knob on  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$ . (Why it is worth fussing with strict consistency.)

Thus we know how to understand the effect of the “knob”  $\mathcal{H}$ .

Note it is impossible to compute (even in principle) since we do not know  $\mathbf{P}_{\mathbf{Z}^m}$  (the distribution from which  $\mathbf{z}$  is drawn).

Can upper bound by  $\sup_{\mathbf{z} \in \mathcal{Z}^m} \log \mathcal{N}(\varepsilon, \mathcal{H}, \ell_1(\mathbf{z}))$  which can be effectively bounded.

Leads to “generalization bounds” of the form: for  $\mathbf{z} \in \mathcal{Z}^m$

$$\mathbf{P}_{\mathbf{Z}^m} \left( R_l \left[ \mathcal{A}_{\text{erm}}^{\mathcal{H}}(\mathbf{Z}) \right] \leq \psi_l(\mathcal{H}, \mathbf{Z}, \delta) \right) \geq 1 - \delta.$$



# The “Key Theorem” in Learning Theory

---

“ERM is strictly consistent iff covering numbers behave nicely”

Observe that whilst we set out to understand the behaviour of  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  our bounds are in fact for

$$\mathcal{A}_{\text{worst}}^{\mathcal{H}} := \mathbf{z} \mapsto \arg \max_{h \in S(\mathcal{H}, \mathbf{z})} R_l[h]$$

where

$$S(\mathcal{H}, \mathbf{z}) = \left\{ h \in \mathcal{H} : \widehat{R}_l[h, \mathbf{z}] = \widehat{R}_l[\mathcal{A}_{\text{erm}}^{\mathcal{H}}, \mathbf{z}] \right\}.$$

Consequently the bounds are very loose.

Furthermore  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  could perform as poorly as  $\mathcal{A}_{\text{worst}}^{\mathcal{H}}$  (what is there to stop it?).

**Conclusion:** behaviour of covering numbers is the crucial quantity for this inductive principle (algorithm). Suggests to make  $\mathcal{H}$  as small as possible.



## Another Algorithm: SRM

---

An obvious difficulty with  $\mathcal{A}_{\text{erm}}^{\mathcal{H}}$  is that if one chooses  $\mathcal{H}$  badly, the algorithm has no hope of approximating the data.

Suppose for  $\mathbf{z} \in \mathcal{Z}^m$ , we know

$$\mathbf{P}_{\mathcal{Z}^m} \left( R_l \left[ \mathcal{A}_{\text{erm}}^{\mathcal{H}}(\mathbf{Z}) \right] \leq \psi_l(\mathcal{H}, \mathbf{Z}, \delta) \right) \geq 1 - \delta.$$

Given a sequence of nonnegative numbers  $\boldsymbol{\delta} = (\delta_i)_{i \in \mathbb{N}}$  such that  $\sum_i \delta_i = \delta$  and a sequence of hypothesis classes  $\mathcal{H} = (\mathcal{H}_i)_{i \in \mathbb{N}}$

$$i^* = i^*(\mathbf{z}, \mathcal{H}, \boldsymbol{\delta}, \psi) := \arg \min_{i \in \mathbb{N}} \psi_l(\mathcal{H}_i, \mathbf{z}, \delta_i)$$

$$\mathcal{A}_{\text{srm}}^{\mathcal{H}, \boldsymbol{\delta}}(\mathbf{z}) := \mathcal{A}_{\text{erm}}^{\mathcal{H}_{i^*}}(\mathbf{z}).$$



By the definition of  $\mathcal{A}_{\text{srm}}^{\mathcal{H}, \delta}$  it comes with a performance bound already. For  $i \in \mathbb{N}$ , with probability at least  $1 - \delta_i$  over a random draw of  $\mathbf{z}$ ,

$$R \left[ \mathcal{A}_{\text{erm}}^{\mathcal{H}_i}(\mathbf{z}) \right] \leq \psi_l(\mathcal{H}_i, \mathbf{z}, \delta_i)$$

Thus the union bound ensures that with probability at least  $1 - \delta$  over a random draw of  $\mathbf{z}$ , for all  $i \in \mathbb{N}$

$$R \left[ \mathcal{A}_{\text{erm}}^{\mathcal{H}_i}(\mathbf{z}) \right] \leq \psi_l(\mathcal{H}_i, \mathbf{z}, \delta_i)$$

and thus with probability at least  $1 - \delta$  over a random draw of  $\mathbf{z}$ ,

$$R \left[ \mathcal{A}_{\text{srm}}^{\mathcal{H}, \delta}(\mathbf{z}) \right] \leq \psi_l(\mathcal{H}_{i^*}, \mathbf{z}, \delta_{i^*}) \quad \spadesuit$$

# Algorithm Independence of Bound

Classical bound takes form: with probability at least  $1 - \delta$  over a random draw of  $\mathbf{z} \in \mathcal{Z}^m$  according to  $\mathbf{P}_{\mathcal{Z}^m}$ ,

$$R_l[\mathcal{A}_{\text{erm}}^{\mathcal{H}}(\mathbf{z})] \leq \psi(\mathcal{H}, \mathbf{z}, \delta).$$

Thus *any* algorithm

$$\mathcal{A}_{\text{any}}^{\mathcal{H}} : \mathcal{Z}^m \rightarrow \mathcal{H}$$

for which  $\widehat{R}_l[\mathcal{A}_{\text{any}}^{\mathcal{H}}(\mathbf{z}), \mathbf{z}] = 0$  has the same bound on performance.

This is *good* because one gets a general theory.

It is *bad* because the same bound holds for the *worst* algorithm.

Bayesians would say the problem is that we are not *conditioning on the data*.

“You must *condition on the data!*”

— (Spirit of) Thomas Bayes





# Why Conditioning on the Data is Important

---

That bayesians and frequentists are willing to discuss these matters is an important first step toward developing a theory that synthesizes both unconditional and conditional inference.

Why condition — a simple example.

Suppose  $X = (X_1, X_2)$ ,  $X_1, X_2$  iid according to

$$\mathbf{P}_\theta(X_i = \theta - 1) = \mathbf{P}_\theta(X_i = \theta + 1) = \frac{1}{2}$$

for  $-\infty < \theta < \infty$ .

Consider the “confidence procedure”

$$C(x) := \begin{cases} \frac{x_1 + x_2}{2} & \text{if } |x_1 - x_2| = 2 \\ x_1 - 1 & \text{if } |x_1 - x_2| = 0 \end{cases}$$

where  $x = (x_1, x_2)$ . Can check that

$$\mathbf{P}_\theta(C(X) \text{ contains } \theta) = 0.75 \quad \forall \theta$$

so we would be happy using  $C(X)$  according to standard frequentist notions of acceptability.

— George Casella (1988)





# Why Conditioning on the Data is Important

---

**But** after one sees the data:

If  $|x_1 - x_2| = 2$  know *for certain* that  $\theta \in C(X)$

If  $|x_1 - x_2| = 0$ , equally unsure whether  $\theta = x_1 - 1$  or  $x_1 + 1$ .

Statisticians have expended considerable effort to develop procedures that have frequentist guarantees of performance *and* which can condition on the data to exploit a lucky observation.

Bayesian methods intrinsically condition on the data, but offer no frequentist guarantees of performance (most Bayesians would say this is no problem because such guarantees are neither necessary nor useful).

Many subtleties. To date really only for simple parameter estimation.

Something like this is needed in order to provide frequentist guarantees of performance for learning algorithms that do more than merely minimize and empirical risk functional.

---

# Maximum Margin Algorithm

Maximum Margin Classifier:  $\mathcal{H}_i(\mathbf{z})$  comprises hyperplanes  $h_{\mathbf{w}}$  achieving margin  $\gamma_{\mathbf{z}}(h_{\mathbf{w}}) = \gamma_i$  on  $\mathbf{z}$ . Here

If there exists one separating hyperplane then there exist many others. Why not choose the optimal one?  
— Vladimir Vapnik

$$\gamma_{\mathbf{z}}(h_{\mathbf{w}}) := \max_{(x_i, y_i) \in \mathbf{z}} y_i \langle \mathbf{w}, x_i \rangle / \|\mathbf{w}\|.$$

For linear hyperplanes, the risk bound is of the form

$$\psi(i) \leq \frac{c}{\gamma_i^2} \log^2(m) + c \log(1/\delta).$$



Maximum Margin algorithm: For  $\mathcal{H}$  the set of linear hyperplanes.

$$\mathcal{A}_{\text{MM}} := \mathbf{z} \mapsto \arg \max_{h \in \mathcal{H}} \gamma_{\mathbf{z}}(h)$$

(“Optimal” only in the sense that it optimizes the particular  $\psi$  function used.)

Try to understand  $\mathcal{A}_{\text{MM}}$  as an instance of  $\mathcal{A}_{\text{SRM}}^{\mathcal{H}, \delta}$ .



# Data Dependent SRM: $\mathcal{A}_{\text{dsrm}}^{\mathcal{H}(\mathbf{z}), \delta}$

---

New algorithm: *Penalize complexity of  $\mathcal{H}(\mathbf{z})$  as if independent of  $\mathbf{z}$ .* Consider  $\mathcal{H}(\mathbf{z}) = (\mathcal{H}_i(\mathbf{z}))_i$ . Suppose for data-independent  $\mathcal{H}_i$  and  $h_{\text{emp}}^i = \mathcal{A}_{\text{erm}}^{\mathcal{H}_i}(\mathbf{z})$  have a bound

$$R_l[h_{\text{emp}}^i] \leq \psi(\mathcal{H}, \mathbf{z}, \delta_i) =: \chi(i)$$

Let

$$i^* := \arg \min_i \chi(i) \qquad \mathcal{A}_{\text{dsrm}}^{\mathcal{H}(\mathbf{z}), \delta}(\mathbf{z}) := \mathcal{A}_{\text{erm}}^{\mathcal{H}_{i^*}}(\mathbf{z})$$

Gist: penalize complexity ignoring data-dependence; apply SRM.

Problem: how to rigorously justify?

# How to conceptualize what's going on?

---

Can not strictly justify the algorithm as an application of SRM.

Adherents of the various non-NPW schools take advantage of "lucky observations" to make more conclusive sounding statements than they would for "unlucky outcomes".

— Jack Kieffer: Conditional Confidence Statements . . . 1977

If as well as  $\hat{R}_l[h, \mathbf{z}] = 0$  we have  $\gamma_{\mathbf{z}}(h) = \gamma \gg 0$ , then  $R_l[h]$  is small.



We are *lucky* if our data is like this.

We want to **condition on the data** like Bayesians do.

Would like to capture this notion in a general formal way.



## BREAK

Have covered

- Induction (what's "inductive"?)
- Inductive Principles (what's the "principle"?)
- Empirical Risk Minimization
- Key Theorem of Learning Theory
- Conditioning on the Data

Yet to come:

- "Conditioning on the data" in a Frequentist (PAC) setting — The Luckiness Framework
- A new approach — Algorithmic Luckiness

# Luckiness

The margin  $\gamma_{\mathbf{z}}(h_{\mathbf{w}})$  measures how *lucky*  $h_{\mathbf{w}}$  is on  $\mathbf{z}$ .

In general  $L: \mathcal{H} \times \mathcal{Z}^m \rightarrow \mathbb{R}$ .

Would like a bound that says with probability at least  $1 - \delta$  over a random draw of  $\mathbf{z}$  according to  $\mathbf{P}_{\mathcal{Z}^m}$  if  $\widehat{R}_l[h, \mathbf{z}] = 0$  and  $\omega(L(h, \mathbf{z}), \delta) \leq 2^d$  then

$$R_l[h] \leq \psi(m, d) \quad \diamond$$

The parameter  $d$  is an *effective complexity*.

There needs to be some restrictions on  $L$ : if we “use up” all of the information in the sample estimating its luckiness there is “none left” to estimate  $R_l[h]$ .

A bound like  $\diamond$  is a bound for the algorithm

$$\mathcal{A}_{\text{lucky}}^{L, \mathcal{H}} := \mathbf{z} \mapsto \arg \min_{h \in \mathcal{H}} \psi(m, \log \omega(L(h, \mathbf{z}), \delta))$$

I should be so lucky; lucky, lucky, lucky.  
I should be so lucky in love.

— Kylie Minogue



# Luckiness (continued)

Given a *luckiness function*  $L: \mathcal{H} \times \mathcal{Z}^m \rightarrow \mathbb{R}$ , the *level* is

$$\ell_L(h, \mathbf{z}) := |\{(l(g(x_i), y_i))_{i=1}^m : L(g, \mathbf{z}) \geq L(h, \mathbf{z})\}|,$$

the number of dichotomies induced on  $\mathbf{z}$  by functions at least as lucky as  $h$ . Require  $L$  to be well behaved:



$L$  is *probably smooth* w.r.t.  $\omega: \mathbb{R} \times (0, 1] \rightarrow \mathbb{N}$  if for all  $m \in \mathbb{N}$  all distributions  $\mathbf{P}_{\mathbf{Z}}$  and all  $\delta \in (0, 1]$

$$\mathbf{P}_{\mathbf{Z}^{2m}}(\exists h \in \mathcal{H} : \ell_L(h, \mathbf{Z}_{[1:2m]}) > \omega(L(h, \mathbf{Z}_{[1:m]}), \delta)) \leq \delta.$$

If  $L$  is probably smooth w.r.t.  $\omega$ ,  $\delta = (\delta_i)_i$ ,  $\sum_i \delta_i = \delta$ , with probability at least  $1 - \delta$  over a random draw of  $\mathbf{z}$ , if  $\widehat{R}_l[h, \mathbf{z}] = 0$  and  $\omega(L(h, \mathbf{z}), \delta_d/4) \leq 2^d$  then

$$R_l[h] \leq \frac{2}{m} (d + \log_2(4/\delta_d)) \quad \heartsuit$$

Effectively  $\mathcal{H}(\mathbf{z}) = (\mathcal{H}_i(\mathbf{z}))_i$  with  $\mathcal{H}_i(\mathbf{z}) = \{h \in \mathcal{H} : \omega(L(h, \mathbf{z}), \delta_i/4) \leq 2^i\}$ .



## Comments on Luckiness

- Can put  $\mathcal{A}_{\text{MM}}$  into this framework.
- The luckiness function  $L$  is how we *encode our prior knowledge*. We weight with  $\delta_i$  the  $i$ th data-dependent hypothesis class

$$\mathcal{H}_i(\mathbf{z}) = \{h \in \mathcal{H} : \omega(L(h, \mathbf{z}), \delta_i/4) \leq 2^i\}$$

- Key practical difficulty is showing  $L$  is probably smooth w.r.t. a “good”  $\omega$  — the smaller the  $\omega$  the tighter  $\heartsuit$  is.
- **Problem:** Still do not pay enough attention to the *algorithm*, which motivates . . .



# Algorithmic Luckiness — Foundation

It is possible to prove a **basic lemma** for algorithms which means that the **symmetrisation by permutation** step only considers all hypotheses that can be learned using  $\mathcal{A}$ .

Basic lemma says:

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

$$2 \cdot \mathbf{P}_{\mathbf{Z}^{2m}} \left( \widehat{R}_l [\mathcal{A}(\mathbf{Z}_{[1:m]}), \mathbf{Z}_{[(m+1):2m]}] - \widehat{R}_l [\mathcal{A}(\mathbf{Z}_{[1:m]}), \mathbf{Z}_{[1:m]}] > \frac{\varepsilon}{2} \right)$$

Again, we introduce an ordering between the at most  $(2m)!$  hypotheses using an **algorithmic luckiness**  $L(\mathcal{A}, \mathbf{z})$ . This gives

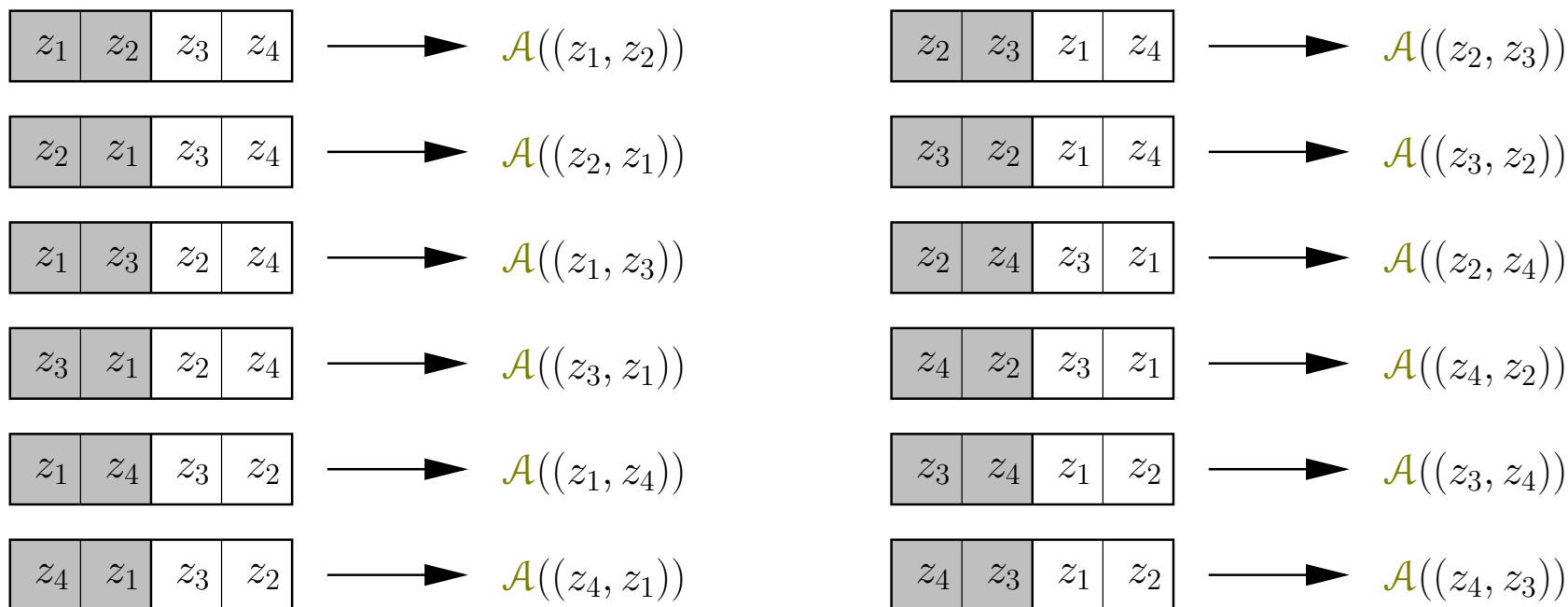
$$\mathcal{H}(\mathcal{A}, L, \mathbf{z}) := \left\{ \mathcal{A} \left( \Pi_{\mathbf{i}}(\mathbf{z})_{[1:m]} \right) \mid \mathbf{i} \in \mathcal{J}(\mathcal{A}, L, \mathbf{z}) \right\},$$

$$\mathcal{J}(\mathcal{A}, L, \mathbf{z}) := \left\{ \mathbf{i} \mid L \left( \mathcal{A}, \Pi_{\mathbf{i}}(\mathbf{z})_{[1:m]} \right) \geq L(\mathcal{A}, \mathbf{z}_{[1:m]}) \right\}.$$





# Illustration of Basic Lemma when $m = 2$



Consider the simple case of  $m = 2$ . Consider all the hypotheses generated by  $\mathcal{A}$  and take  $\mathcal{H}(\mathcal{A}, L, \mathbf{z})$  to be those so generated that are at least as lucky as  $\mathcal{A}((z_1, z_2))$  where the luckiness is measured on  $(z_1, z_2)$ .

# Algorithmic Luckiness — $\omega$ -smallness

Algorithmic luckiness  $\mathcal{H}(\mathcal{A}, L, \mathbf{z})$  is a **function of  $\mathcal{A}$  directly**.

Need to be able to bound the covering number  $\mathcal{N}$  of hypotheses  $h \in \mathcal{H}(\mathcal{A}, L, \mathbf{z})$  on the double sample  $\mathbf{z}$  only using the luckiness on the first half, i.e. the training sample.



**$\omega$ -smallness of  $L$ :** Given an algorithm  $\mathcal{A}$  and a loss  $l$ , the algorithmic luckiness  $L$  is  $\omega$ -small at scale  $\tau$ , if for all  $\delta$

$$\mathbf{P}_{\mathbf{Z}^{2m}}(\mathcal{N}(\mathcal{T}, \mathcal{H}(\mathcal{A}, L, \mathbf{Z}), \ell_{l,1}(\mathbf{Z})) > \omega(L(\mathcal{A}(\mathbf{Z}_{[1:m]})), \delta, \tau)) < \delta.$$

To prove this property we can only exploit that  $\mathbf{P}_{\mathbf{Z}^{2m}}$  is a product measure.

# Algorithmic Luckiness — Main Result

**Algorithmic Luckiness Bound:** For all  $[0, 1]$ -valued loss functions  $l$ , for all  $\omega$ -small algorithmic luckiness functions  $L$  w.r.t.  $\mathcal{A}$ , for all  $\tau$ , with probability at least  $1 - \delta$  over  $\mathbf{z} \in \mathcal{Z}^m$ ,

And the winner is ... Lucky!  
— Britney Spears



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

where  $d = \log \left( \omega \left( L(\mathcal{A}, \mathbf{z}), \frac{\delta}{4m}, \tau \right) \right)$ . If  $l$  is  $\{0, 1\}$ -valued, then whenever  $\mathcal{A}(\mathbf{z})$  has zero training error,  $\widehat{R}_l[\mathcal{A}(\mathbf{z}), \mathbf{z}] = 0$ , for  $\tau = 1/2m$

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

## Application — Classical VC Setting

If  $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^{\mathcal{X}}$  maps to a hypothesis space  $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ , we know that

$$\mathcal{H}(\mathcal{A}, L, \mathbf{z}) \subseteq \mathcal{H}$$

regardless of  $\mathbf{z}$  and  $L$ .

VC-dimension is not most important quantity.

— Vladimir Vapnik: Dagstuhl, Germany (July 2001)



Thus, for the zero-one loss  $l(\hat{y}, y) = \mathbb{I}_{\hat{y} \neq y}$  the **growth function** is an upper bound on  $\mathcal{N}\left(|\mathbf{z}|^{-1}, \mathcal{H}(\mathcal{A}, L, \mathbf{z}), \ell_{l,1}(\mathbf{z})\right)$  and can thus serve as a  $\omega$  function.

Neither the serendipity of the sample nor the properties of the algorithm  $\mathcal{A}$  have been exploited!

As is widely known, this results in bounds which are quite loose.

# Application — Compression Bounds

---

I used to be God...

— Manfred Warmuth: Dagstuhl, Germany (2001)

**Sparsity luckiness:** If  $\mathcal{A} : \mathcal{Z}^{(\infty)} \rightarrow \mathcal{Y}^x$  is a compression scheme, that is,  $\mathcal{A}(\mathbf{z}) = \mathcal{R}(\mathcal{C}(\mathbf{z}))$ , then

$$L_{\text{sparse}}(\mathcal{A}, \mathbf{z}) := -|\mathcal{C}(\mathbf{z})|$$

is  $\omega$ -small at any scale  $\tau$ , where

$$\omega(L, \delta, \tau) = \left( \frac{2em}{-L} \right)^{-L}.$$

Plugging this result into ♣ gives a new compression result for regression as well as resembling the original result of Littlestone and Warmuth (1986).



# Sparsity Luckiness — Proof

Since we only have to consider permutations  $\Pi_i$  where

$$\left| \mathcal{C} \left( \Pi_i(\mathbf{z})_{[1:m]} \right) \right| \leq \left| \mathcal{C}(\mathbf{z}_{[1:m]}) \right| =: -L_0$$

we know that the permutation invariant reconstruction function  $\mathcal{R}$  never uses more than  $-L_0$  examples.

The number of different choices of no more than  $-L_0$  examples out of  $2m$  (double sample size) is given by

$$\sum_{i=0}^{-L_0} \binom{2m}{i} \leq \left( \frac{2em}{-L_0} \right)^{-L_0} .$$

You must be my lucky star . . .  
But I'm the luckiest by far.  
— Madonna





## Application — Kernel Classifiers

Consider learning algorithms for **kernel classifiers**, that is,

$$\mathcal{H}_\phi := \{ \mathbf{x} \mapsto \langle \phi(\mathbf{x}), \mathbf{w} \rangle \mid \mathbf{w} \in \mathcal{K} \}, \quad \phi : \mathcal{X} \rightarrow \mathcal{K} \subseteq \ell_2^n.$$

Assume that the learning algorithms  $\mathcal{A}$  have the property that

$$\mathcal{A} : \mathbf{z} = (\mathbf{x}, \mathbf{y}) \mapsto \langle \phi(\mathbf{x}), \mathbf{w}_z \rangle \text{ where } \mathbf{w}_z = \sum_{x_i \in \mathbf{x}} \hat{\alpha}_i \phi(x_i) \text{ and } \|\mathbf{w}_z\| = 1.$$

Examples are SVMs, BPMs and the perceptron algorithm.

Let the (normalised) **margin**  $\Gamma(\mathbf{z})$  be defined by

$$\Gamma(\mathbf{z}) := \min_{(x_i, y_i) \in \mathbf{z}} \frac{y_i \langle \phi(x_i), \mathbf{w}_z \rangle}{\|\phi(x_i)\| \cdot \|\mathbf{w}_z\|}.$$



## Application — Kernel Classifiers (cont.)

**Margin Luckiness:** Let  $\varepsilon_i(\mathbf{x})$  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 to  $\epsilon$ . For the loss  $l(\hat{y}, y) = \mathbb{I}_{y\hat{y} \leq 0}$ , the luckiness function



$$L_{\text{margin}}(\mathcal{A}, \mathbf{z}) = - \min \left\{ i \in \mathbb{N} \mid i \geq \left( \frac{\varepsilon_i(\mathbf{x}) \sum_{j=1}^m |\hat{\alpha}_j|}{\Gamma(\mathbf{z})} \right)^2 \right\}$$

is  $\omega$ -small at scale  $1/2m$  where

$$\omega(L, \delta, 1/2m) = \left( \frac{2em}{-L} \right)^{-2L}.$$

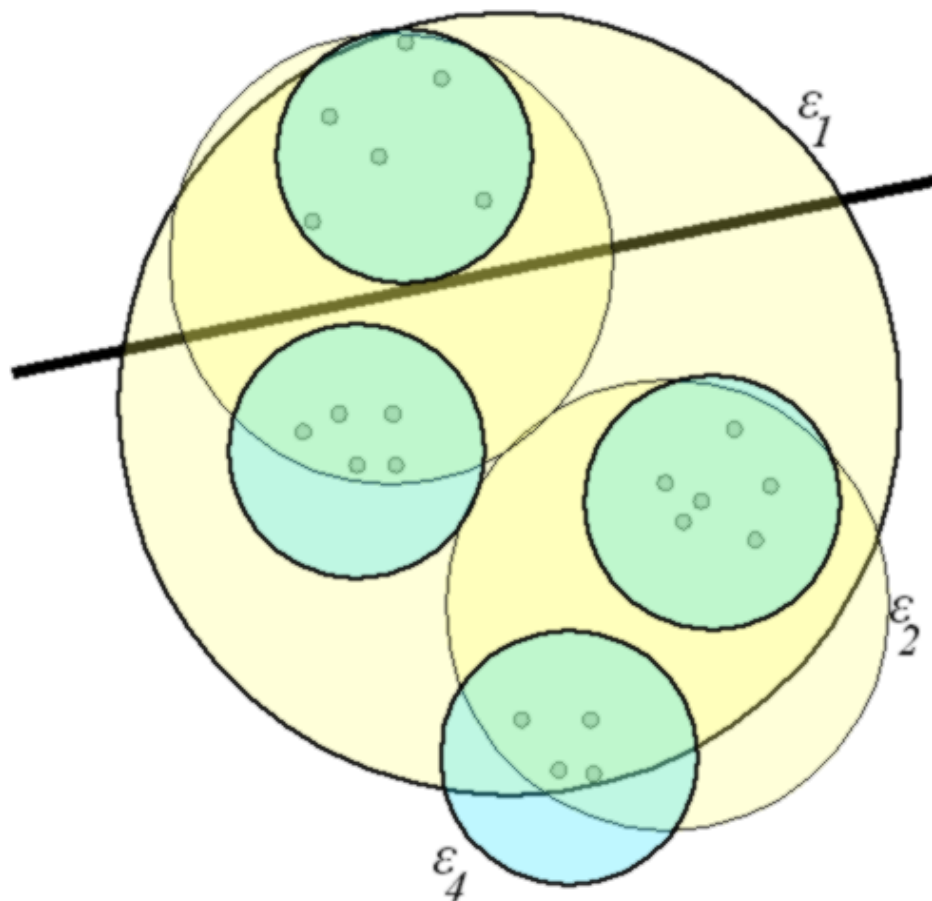
The bound comprises 3 main terms: *margin*  $\Gamma(\mathbf{z})$ , *sparsity surrogate*  $\sum_{j=1}^m |\hat{\alpha}_j|$  and a factor depending on the *distribution of the data*  $\varepsilon_i(\mathbf{x})$ .

## Kernel Classifiers — Meaning of $\varepsilon_i(\mathbf{x})$

The sequence  $(\varepsilon_i(\mathbf{x}))_i$  measures how clumpy the data is.

A small number of small clumps means  $\varepsilon_i(\mathbf{x})$  is small for small  $i$ .

Compare with the idea of “kernel alignment”.



## Application — Kernel Classifiers (proof)

**Makovoz theorem** shows that for all  $\mathbf{z} \in \mathcal{Z}^m$  there exists a weight vector  $\tilde{\mathbf{w}} = \sum_{i=1}^m \tilde{\alpha}_i \phi(x_i)$  such that

$$\|\tilde{\mathbf{w}} - \mathbf{w}_{\mathbf{z}}\|^2 \leq \Gamma^2(\mathbf{z})$$

and  $\|\tilde{\alpha}\|_0 \leq -L_{\text{margin}}(\mathcal{A}, \mathbf{z}) =: -L_0$ .



It follows that  $\langle \mathbf{w}_{\mathbf{z}}, \tilde{\mathbf{w}} / \|\tilde{\mathbf{w}}\| \rangle \geq \sqrt{1 - \Gamma^2(\mathbf{z})}$ ; that is,  $\tilde{\mathbf{w}}$  still correctly classifies  $\mathbf{z}$ .

For every of the no more than  $\binom{2em}{-L_0}^{-L_0}$  many subsamples  $\tilde{\mathbf{z}} \subseteq \mathbf{z}$ ,  $\tilde{\mathbf{w}}$  lives in a space of dimension no more than  $-L_0$ .

By an application of the **growth function bound**, each  $\tilde{\mathbf{w}}$  can achieve no more than  $\binom{2em}{-L_0}^{-L_0}$  many dichotomies on  $\mathbf{z}$ .



---

Algorithmic luckiness framework differs from classical statistical learning theory approaches in that it does not use the crude step of viewing algorithms just in terms of their hypothesis space.

Generalization of standard VC results. Get agnostic and realizable bounds.

Example of maximum margin algorithm illustrates that the framework has the power to develop new insights into what makes algorithms perform well.

Main point is that it provides new theoretical tools for understanding algorithms “smarter” than Empirical Risk Minimization.

Hope is that by analysing algorithms in this (or related) ways, will be able to better discern the features about particular learning problems that make them easy or difficult.



## References and Slides

Jack Kiefer, “Conditional Confidence Statements and Confidence Estimators,” *Journal of the American Statistical Association*, **72**(360), pp. 789–827, (1977).

James Berger, “The Frequentist Viewpoint and Conditioning”, pp 15–44 in *Proceedings of the Berkeley Conference in Honour of Jerzy Neyman and Jack Kiefer*, Volume 1, Wadsworth, (1985).

George Casella, “Conditionally Acceptable Frequentist Solutions,” pp. 73–117 in *Statistical Decision Theory and Related Topics IV*, Springer, New York (1988).

John Shawe-Taylor, Peter L. Bartlett, Robert C. Williamson and Martin Anthony, “Structural Risk Minimization over Data-Dependent Hierarchies,” *IEEE Transactions on Information Theory*, **44**(5), 1926–1940, (1998). <http://axiom.anu.edu.au/~williams/papers/P85.ps>

Ralf Herbrich, *Learning Kernel Classifiers*, MIT Press, (2002).

Ralf Herbrich and Robert C. Williamson, “Algorithmic Luckiness” submitted to *Journal of Machine Learning Research* (December 2001) <http://axiom.anu.edu.au/~williams/papers/P159.ps.gz>

Ralf Herbrich and Robert C. Williamson, “Learning and Generalization: Theoretical Bounds” to appear in Michael Arbib (Ed.) *Handbook of Brain Theory and Neural Networks*, 2nd Edition, MIT Press, (2002). <http://axiom.anu.edu.au/~williams/papers/P158.ps.gz>

Slides at <http://axiom.anu.edu.au/~williams/papers/P156.pdf>