

No-free-lunch theorem for machine learning

Michikazu Hirata

February 6, 2026

Abstract

This entry is a formalization of the no-free-lunch theorem for machine learning following Section 5.1 of the book *Understanding Machine Learning: From Theory to Algorithms* [1] by Shai Shalev-Shwartz and Shai Ben-David. The theorem states that for binary classification prediction tasks, there is no universal learner, meaning that for every learning algorithms, there exists a distribution on which it fails.

Contents

1	No-Free-Lunch Theorem for ML	1
1.1	Preliminaries	1
1.2	No-Free-Lunch Theorem	2

1 No-Free-Lunch Theorem for ML

```
theory No-Free-Lunch-ML
imports
  HOL-Probability.Probability
begin
```

1.1 Preliminaries

```
lemma sum-le-card-Max-of-nat:finite A
   $\implies \text{sum } f \ A \leq (\text{of-nat } :: - \implies - :: \{\text{semiring-1, ordered-comm-monoid-add}\}) (\text{card } A) * \text{Max } (f \ ' \ A)$ 
  <proof>
```

```
lemma card-Min-le-sum-of-nat: finite A
   $\implies (\text{of-nat } :: - \implies - :: \{\text{semiring-1, ordered-comm-monoid-add}\}) (\text{card } A) * \text{Min } (f \ ' \ A) \leq \text{sum } f \ A$ 
  <proof>
```

The following lemma is used to show the last equation of the proof of the no-free-lunch theorem in the book [1].

Let A be a finite set. If A is divided into the pairs $(x_1, y_1), \dots, (x_n, y_n)$ such that $f(x_i) + f(y_i) = k$ for all $i = 1, \dots, n$. Then, we have $\sum_{x \in A} f(x) = k * |A|/2$.

lemma *sum-of-const-pairs*:

fixes $k :: \text{real}$
assumes $A::\text{finite } A$
and $\text{fst } ' B \cup \text{snd } ' B = A \text{ fst } ' B \cap \text{snd } ' B = \{\}$
and $\text{inj-on } \text{fst } B \text{ inj-on } \text{snd } B$
and $\text{sum}: \bigwedge x y. (x, y) \in B \implies f x + f y = k$
shows $(\sum_{x \in A}. f x) = k * \text{real } (\text{card } A) / 2$
<proof>

lemma(**in** *prob-space*) *Markov-inequality-measure-minus*:

assumes $u \in \text{borel-measurable } M$ **and** $A E x \text{ in } M. 0 \leq u x \ A E x \text{ in } M. 1 \geq u x$
and $[\text{arith}]: 0 < (a::\text{real})$
shows $\mathcal{P}(x \text{ in } M. u x > 1 - a) \geq ((\int x. u x \ \partial M) - (1 - a)) / a$
<proof>

1.2 No-Free-Lunch Theorem

In our implementation, a learning algorithm of binary classification is represented as a function $A : \text{nat} \Rightarrow (\text{nat} \Rightarrow 'a \times \text{bool}) \Rightarrow 'a \Rightarrow \text{bool}$ where the first argument is the number of training data, the second argument is the training data ($S \ n = (x_n, y_n)$ denotes the n th data for a training data S), and $A \ m \ S$ is a predictor. The first argument, which denotes the number of training data, is normally used to specify the number of loop executions in learning algorithm. In this formalization, we omit the first argument because we do not need the concrete definitions of learning algorithms.

Let X be the domain set. In order to analyze the error of predictors, we assume that each data (x, y) is obtained from a distribution \mathcal{D} on $X \times \mathbb{B}$. The error of a predictor f with respect to \mathcal{D} is defined as follows.

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(f) &\stackrel{\text{def}}{=} \mathbb{P}_{(x,y) \sim \mathcal{D}} (f(x) \neq y) \\ &= \mathcal{D}(\{(x, y) \in X \times \mathbb{B} \mid f(x) \neq y\}) \end{aligned}$$

In these settings, the no-free-lunch theorem states that for any learning algorithm A and $m < |X|/2$, there exists a distribution \mathcal{D} on $X \times \mathbb{B}$ and a predictor f such that

- $\mathcal{L}_{\mathcal{D}}(f) = 0$, and
- $\mathbb{P}_{S \sim \mathcal{D}^m} \left(\mathcal{L}_{\mathcal{D}}(A(S)) > \frac{1}{8} \right) \geq \frac{1}{7}$.

theorem *no-free-lunch-ML*:

fixes $X :: 'a \text{ measure}$ **and** $m :: \text{nat}$

and $A :: (\text{nat} \Rightarrow 'a \times \text{bool}) \Rightarrow 'a \Rightarrow \text{bool}$

assumes $X1:\text{finite}(\text{space } X) \Longrightarrow 2 * m < \text{card}(\text{space } X)$

and $X2[\text{measurable}]: \bigwedge x. x \in \text{space } X \Longrightarrow \{x\} \in \text{sets } X$

and $m[\text{arith}]: 0 < m$

and $A[\text{measurable}]: (\lambda(s,x). A s x) \in (\text{PiM } \{..<m\} (\lambda i. X \otimes_M \text{count-space} (UNIV :: \text{bool set}))) \otimes_M X$

$\rightarrow_M \text{count-space} (UNIV :: \text{bool set})$

shows $\exists \mathcal{D} :: ('a \times \text{bool}) \text{ measure. sets } \mathcal{D} = \text{sets } (X \otimes_M \text{count-space} (UNIV :: \text{bool set})) \wedge$

$\text{prob-space } \mathcal{D} \wedge$

$(\exists f. f \in X \rightarrow_M \text{count-space} (UNIV :: \text{bool set}) \wedge \mathcal{P}((x, y) \text{ in } \mathcal{D}. f x \neq y) = 0) \wedge$

$\mathcal{P}(s \text{ in } \text{PiM } \{..<m\} (\lambda i. \mathcal{D}). \mathcal{P}((x, y) \text{ in } \mathcal{D}. A s x \neq y) > 1 / 8) \geq 1 / 7$
<proof>

end

References

- [1] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.