

The Lack of A Priori Distinctions Between Learning Algorithms aka a No Free Lunch Theorem for Learning

Nikita Kazeev, based on David H. Wolpert

March 14, 2017

Abstract

Machine learning is about genralization. The perormance measured on not yet seen data is how we compare algorithms. In this paper we show that a perfect universal learner is impossible. If there are no restrictions on the strucutre of the problem, then for any two algorithms there are “as many” targets on which each outperforms the other.

1 Introduction

TODO paper link

We have a nice array of impossibility theorems. Goedel, halting, Arrow. NFL is one more.

We make practical climes. Empirically ML works) However theory-side, a universal learner is impossible.

2 Formalism

Glossary

F f value space in. 2

OTS error

$$P(q|d) = \frac{\delta(q \notin d_X)\pi(q)}{\sum_q [\delta(q \notin d_X)\pi(q)]}, \quad (1)$$

where $\delta(z) \equiv 1$ if z is true and 0 otherwise. . 3

vertical $P(d|f)$: iff $P(d|f)$ is independent of the values $f(x, y_F)$ for $x \notin d_X$. . 3

Integral

$$\int A(f)df = \int_{\mathbf{R}^{rn}} A(f)df \prod_{i=0}^n \left[\delta \left(\sum_{j=i \times r}^{(i+1) \times r - 1} f_j - 1 \right) \prod_{j=i \times r}^{(i+1) \times r - 1} \theta(f_j) \right] \quad (2)$$

3 No Free Lunch

Lemma 1.

$$P(c|d, f) = \sum_{y_H, y_F, q} \delta[c, L(y_H, y_F)] P(y_H|q, d) P(y_F|q, f) P(q|d) \quad (3)$$

- $P(q|d)$ – conditional probability of test set q given training set d
- $P(y_F|q, f)$ – conditional probability of given target sample y_F for the given the target distribution f and test set q
- $P(y_H|q, d)$ – conditional probability of given predicted sample y_H given the test and training sets. Is a function of the learning algorithm.
- $P(c|f, d)$ – conditional probability of given cost c given target f and training set d .

Proof.

$$c = L(y_H, y_F) \quad (4)$$

$$\begin{aligned} P(c|q, d, f) &= \sum_{y_H, y_F} \delta[c, L(y_H, y_F)] P(y_H, y_F|q, d, f) \\ P(c|d, f) &= \sum_{y_H, y_F, q} \delta[c, L(y_H, y_F)] P(y_H, y_F|q, d, f) P(q|d) \\ &= \sum_{y_H, y_F, q} \delta[c, L(y_H, y_F)] P(y_H|q, d) P(y_F|q, f) P(q|d) \end{aligned} \quad (5)$$

□

CONSIDER example with random guessing.

Theorem 3.1. *For homogenous loss L , the uniform average over all f of $P(c|d, f)$ equals $\Lambda(c)/r$.*

Proof. The uniform average over all targets f of $P(c|f, d)$ equals to

$$E_f[P(c|f, d)] = \sum_{y_H, y_F, q} \delta[c, L(y_H, y_F)] P(y_H|q, d) E_f[P(y_F|q, f)] P(q|d) \quad (6)$$

$$E_f[P(y_F|q, f)] = E_f f(q, y_F) \quad (7)$$

Because F is symmetric, the average is a constant that does not depend on q and y_F . Also, $\sum_{y_F} E_f[f(y_F, q)] = E_f\left[\sum_{y_F} f(y_F, q)\right] = 1$, therefore $E_f[f(y_F, q)] = 1/r$.

Using the homogeneity property of L :

$$E_f P(c|f, d) = \sum_{y_H, q} \Lambda(c) P(y_H|q, d) P(q|d) / r = \Lambda(c)/r \quad (8)$$

□

Theorem 3.2. For OTS error, a vertical $P(d|f)$, and a homogeneous loss L , the uniform average over all targets f of $P(c|f, m) = \Lambda(c)/r$

Proof.

$$P(c|f, m) = \sum_{d:|d|=m} P(c|f, d) P(d|f) \quad (9)$$

□

From 3:

$$P(c|d, f) = \sum_{y_H, y_F, q} \delta[c, L(y_H, y_F)] P(y_H|q, d) P(y_F|q, f) P(q|d). \quad (10)$$

Because we consider OTS error, $P(q|d)$ will be non-zero only for $q \notin d_X$, so $P(y_F|q, f)$ only depends on components of $f(x, y)$ that correspond to $x \notin d_X$.

We also know that $P(d|f)$ is vertical, so it is independent of the values $f(x, y_F)$ for $x \notin d_X$.

Therefore the integral can be split into two parts, over dimensions corresponding to d_X and $\mathbf{X} \setminus d_X$:

$$E_f [P(c|f, m)] = \frac{\sum_{d:|d|=m} [\int df_{x \notin d_X} P(c|d, f) \int df_{x \in d_X} P(d|f)]}{\int df_{x \notin d_X} df_{x \in d_X} 1} \quad (11)$$

Again using $P(c|d, f)$ independence from $x \in d_X$ and theorem 3.1:

$$E_{f_{x \notin d_X}} [P(c|d, f)] = E_f [P(c|d, f)] = \Lambda(c)/r \quad (12)$$

$$E_f [P(c|f, m)] = \Lambda(c)/r \frac{\sum_{d:|d|=m} [\int df_{x \in d_X} P(d|f)]}{\int df_{x \in d_X} 1} = \Lambda(c)/r \quad (13)$$

Theorem 3.3. For OTS error, a vertical $P(d|f)$, uniform $P(f)$, and a homogeneous loss L , $P(c|d) = \Lambda(c)/r$.

Proof.

$$P(c|d) = \frac{\int df P(c|d, f) P(f|d)}{\int df} \quad (14)$$

Using the Bayes theorem:

$$P(f) = P(f|d) P(d)/P(d|f) \quad (15)$$

$$P(c|d) = \frac{\int df P(c|d, f) P(f) P(d|f)/P(d)}{\int df} = \alpha(d) \int df P(c|d, f) P(d|f), \quad (16)$$

where $\alpha(d)$ is some function. □

4 Implications

[mine?] Empirism philosophy?