## Why

In the statistical learning framework, since the algorithm only has access to the training set, it is natural to select a hypothesis which performs well on the training set.

## Definition

Let $((\Omega, \mathcal{A}, \mathcal{P}), \{x_i : \Omega \to \mathcal{X}\}_{i=1}^n, f : \mathcal{X} \to \mathcal{Y})$ be *probabilistic data-generation model* so that $S : \Omega \to (\mathcal{X} \times \mathcal{Y})^n$ is the training dataset.

The *training error* of a classifier $h : \mathcal{X} \to \mathcal{Y}$ is

$$(1/m)\,|\{i \in \{1, 2, \ldots, m\} \mid h(x_i) \neq y_i\}|\,.$$

Other terminology includes *empirical error* and *empirical risk*. For these reasons, the learning paradigm of selecting a predictor $h$ to minimizer the empirical risk is called *empirical risk minimization* or *ERM*.

## Overfitting

Although selecting a classifier to minimize the empirical risk seems natural, it can be foolish. Let $A \subset \mathcal{X} \subset \mathbf{R}^2$ and $\mathcal{Y} = \{0, 1\}$. Suppose that the true classifier $f : \mathcal{X} \to \mathcal{Y}$ is $f(x) = 1$ if $x \in A$ and $f(x) = 0$ otherwise. Suppose that for the underlying distribution $(\mathcal{X}, \mathcal{A}, \mathbf{P})$ we have $A \in \mathcal{A}$ and $\mathbf{P}(A) = 1/2$.

For the training set $S$ in $\mathcal{X} \times \mathcal{Y}$, the hypothesis $h : \mathcal{X} \to \mathcal{Y}$

defined by

$$h_S(x) = \begin{cases} y_i & \text{if } x_i = x \\ 0 & \text{otherwise.} \end{cases}$$

achieves zero empirical risk but has error (w.r.t. **P**) of 1/2. Such a classifier is said to be *overfit* or to exhibit *overfitting*. It is said to fit the training dataset "too well."

**Inductive Bias**

In spite of the prior example, are there conditions under which ERM does not overfit? One such formulation is to ask whether there are situations in which good performance on empirical risk implies (or makes it highly likely to achieve) good performance on the underling distribution.

One simple approach is to constrain the set of predictors considered. A *hypothesis class* is a subset of predictors $\mathcal{H} \subset \mathcal{X} \to \mathcal{Y}$. The *restricted empirical risk minimizer* for $\mathcal{H}$ is a hypothesis $h \in \mathcal{H}$ minimizing the empirical risk on the data set. Many authorities call the hypothesis class a *inductive bias* and speak of "biasing" the "learning algorithm". Since one specifies the hypothesis class prior to the data it often said to "encode prior knowledge about the problem to be learned."

Empirical Risk Minimizers

Probabilistic Data-Generation Models

Induction

Independent Identically Distributed Random Variables

Random Variable Independence

Datasets

Random Variables

Topological Sigma Algebra

Measurable Function

Outcome Variables

Topological Spaces

Sigma Algebra Independence

Metrics

Event Independence

Absolute Value

Distance

Probability Measure

Conditional Event Probability

Space Distance

Measure

Plane Distance

Event Probabilities

Sigma Algebras

Interval Length

Probability Distributions

Extended Real Numbers

Subset Algebras

Intervals

Real Limits

Cardinality

Real Line

Real Plane

Real Space

Real Sequences

Set Operations

Integral Line

Real Order

Sequences

Real Summation

Geometry

Set Numbers

Real Numbers

Uncertain Outcomes

Finite Sets

Rational Numbers

Direct Products

Integer Arithmetic

Natural Summation

Natural Multiplicative Identity

Natural Additive Identity

Family Operations

Equivalence Sets

Arithmetic

Identity Elements

Algebras

Natural Exponents

Integer Products

Function Inverses

Integer Sums

Operations

Family Unions and Intersections

Natural Products

Function Images

Integer Numbers

Function Composites

Families

Natural Sums

Equivalence Relations

Functions

Natural Order

Recursion Theorem

Relations

Peano Axioms

Cartesian Products

Subset Systems

Natural Induction

Set Powers

Ordered Pairs

Natural Numbers

Unordered Triples

Generalized Set Dualities

Set Symmetric Differences

Partitions

Successor Sets

Intersection of Empty Set

Set Dualities

Set Unions and Intersections

Set Complements

Set Intersections

Pair Unions

Set Differences

Set Unions

Pair Intersections

Empty Set

Unordered Pairs

Set Specification

Set Inclusion

Standardized Accounts

Accounts

Deductions

Quantified Statements

Logical Statements

Statements

Sets

Identities

Names

Letters

Objects