

# Complexity Measures

Vladimir Feinberg

May 30, 2017

Complexity measures evaluate the expressiveness of a hypothesis class; they are useful to the extent with which they relate sample and generalization error.

Notation. TODO move to common place in repo. Big- $O$  notation (and  $\Theta, \Omega$ ) is used in the standard sense for single variables:  $f = O(g)$  if there exists a  $c$  such that  $f(x) \leq cg(x)$  for all  $x$  sufficiently large. In multiple variables, big- $O$  and similarly  $\Theta, \Omega$  require  $f(\mathbf{x}) \leq cg(\mathbf{x})$  for all  $\|x\|_\infty$  sufficiently large. To capture the essence of asymptotics, a tilde will crudely capture asymptotic behavior up to logarithmic factors:  $f = \tilde{O}(g)$  if  $f = O(g \log^k g)$  for some  $k \in \mathbb{N}$ .

Abbreviations. TODO move to common place.

## 1 Setup

We suppose that our data comes in the form of ordered pairs from  $\mathcal{X} \times \mathcal{Y}$ . Samples follow a particular distribution  $(x, y) \sim D$ . A hypothesis class  $\mathcal{H}$  is set of functions  $\mathcal{X} \rightarrow \mathcal{Y}$ .

A common approach to supervised learning is empirical risk minimization, where  $m$  iid samples from  $D$ ,  $S$ , are used to find the  $h \in \mathcal{H}$  minimizing a specified loss  $\ell : \mathcal{Y}^2 \rightarrow \mathbb{R}$  over this set. Complexity measures then let us quantify exactly how much loss we can expect when sampling from  $D$  again.

Thus, it is useful to find bounds on  $\varepsilon$ , the difference between the generalization loss  $\mathbb{E}[\ell(h(x), y)]$ , where  $(x, y) \sim D$ , and sample loss, where the loss is the expectation before taken for  $(x, y)$  is uniform over  $S$ .

Let the gap between the generalization and sample error be  $\varepsilon$ .

## 2 Complexity Measures

Rademacher complexity. TODO. relates complexity through noise correlation. equivalently, gives a consistent view of “error” through a “uniform” noise model.

The growth function. TODO.

VC-dimension. TODO.

## 3 Overview of Results

Proofs can be found in a cogent write-up by Prof. Beckage from the University of Kansas,<sup>1</sup> copied into this repository.<sup>2</sup> Recall our generalization gap definition. For a fixed  $h \in \mathcal{H}$ :

$$\varepsilon = \mathbb{E}[\ell(h(x), y) | (x, y) \sim D] - \mathbb{E}[\ell(h(x), y) | (x, y) \sim \text{Uniform}(S)]$$

---

<sup>1</sup>[http://ittc.ku.edu/~beckage/ml800/VC\\_dim.pdf](http://ittc.ku.edu/~beckage/ml800/VC_dim.pdf)

<sup>2</sup><https://github.com/vlad17/shallow-ml-notes/raw/7c3db7b7c924fbd2ae2891d4ecfef84e5647dcc8/computational-learning/complexity-measures-beckage.pdf>

### 3.1 VC Generalization Bounds

Upper bound. If  $d$  is the VC-dimension of  $\mathcal{H}$ , then for any  $D$  wp  $1 - \delta$ :

$$\varepsilon \leq \tilde{O} \left( \sqrt{\frac{d - \log \delta}{m}} \right)$$

The above inequality is random since it depends on  $S$ , the  $D^m$ -valued rv. TODO. find source removing tilde?

Agnostic lower bound. We may find a  $D$  such that with a fixed nonzero probability (a nonnegligible set of candidate samples  $S$ ), the following holds:

$$\varepsilon \geq \Omega \left( \sqrt{\frac{d}{m}} \right)$$

The above implies that in the common case of agnostic hypothesis learning, where we do not know distribution  $D$ , VC-dimension is, *up to logarithmic factors, asymptotically efficient* in quantifying the generalization gap.

Realizability. Suppose  $D$  is realizable wrt  $\mathcal{H}$ , so that almost surely there exists an  $f \in \mathcal{H}$  such that for any  $(x, y)$  sampled from  $D$ ,  $f(x) = y$ . Then all statements above hold but with  $\sqrt{\varepsilon}$  instead of  $\varepsilon$ .

### 3.2 Rademacher bounds

With  $R_m$  either the empirical or expected Rademacher complexity over the sample for a given  $h, \ell$  we have again wp  $1 - \delta$ :

$$\varepsilon \leq 2R_m + O \left( \frac{\log 1/\delta}{m} \right)$$

$R_m$  may be NP-hard to compute, depending on  $\mathcal{H}$ . This tells us Rademacher complexity could only be a useful improvement over VC-bounds, asymptotically, if we have an efficient approximation for the empirical Rademacher complexity or some knowledge of  $D$  as required to compute the true Rademacher complexity.

TODO. NP-hardness? More computational learning theory.