

Bounding generalization error by a sample-conditioned count of classifiers

Yoram Gat

June 18, 2015

Abstract

This talk will present a result useful for bounding the generalization error of certain classification algorithms. Like VC theory, the bound relies on symmetrization. However, the bound applies to algorithms with high VC-dimension ranges.

One application of the bound is for showing the good generalization behavior of the “support point machine” - a variant of the support vector machine which does not require an inner-product space structure.

Definitions

1. Classifiers and classifications algorithms

- \mathcal{X} - a **feature space**, \mathcal{Y} - a **label space**. Their product is the **labeled feature space** $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$.
- A **classifier** is a map $f : \mathcal{X} \rightarrow \mathcal{Y}$. The **space of classifiers** is \mathcal{F} .
- A **classification algorithm** or **learning algorithm** is a map from a **training set** to a classifier $L : \mathcal{Z}^n \rightarrow \mathcal{F}$.

Definitions

2. Classifier prediction error

- An **error function** corresponding to a classifier f is a map $e_f : \mathcal{Z} \rightarrow \{0, 1\}$:

$$e_f(z) = e_f(x, y) = 1(f(x) \neq y).$$

The **space of error functions** is \mathcal{E} .

- Shorthand for the average error of a sample:

$$e_f^n(\mathbf{Z}) = e_f^n(z_1, \dots, z_n) = \frac{1}{n} \sum_1^n e_f(z_i).$$

Definitions

3. Algorithm generalization error

The **generalization error of a classification algorithm**:

$$g(L) = \mathbb{E}_{\mathbf{Z}} \mathbb{E}_{Z'} e_{L(\mathbf{Z})}(Z') - e_{L(\mathbf{Z})}^n(\mathbf{Z}).$$

Symmetrization

$$\begin{aligned}g(L) &= \mathbb{E}_{\mathbf{Z}} \mathbb{E}_{\mathbf{Z}'} \left[e_{L(\mathbf{Z})}(\mathbf{Z}') - e_{L(\mathbf{Z})}^n(\mathbf{Z}) \right] \\&= \mathbb{E}_{\mathbf{Z}} \mathbb{E}_{\mathbf{Z}'} \left[e_{L(\mathbf{Z})}^{n'}(\mathbf{Z}') - e_{L(\mathbf{Z})}^n(\mathbf{Z}) \right] \\&= \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} \mathbb{E}_{\text{Partition}} \left[e_{L(\bar{\mathbf{Z}})}^{n'}(\bar{\mathbf{Z}}') - e_{L(\bar{\mathbf{Z}})}^n(\bar{\mathbf{Z}}) \right].\end{aligned}$$

Let the pair \mathbf{z}, \mathbf{z}' be the **train-test set** and define the **partition generalization error**:

$$g(L, \mathbf{z} \cup \mathbf{z}') = \mathbb{E}_{\text{Partition}} \left[e_{L(\bar{\mathbf{z}})}^{n'}(\bar{\mathbf{z}}') - e_{L(\bar{\mathbf{z}})}^n(\bar{\mathbf{z}}) \right],$$

giving

$$g(L) = \mathbb{E}_{\mathbf{Z}, \mathbf{Z}'} g(L, \mathbf{Z} \cup \mathbf{Z}').$$

- Given \mathbf{z}, \mathbf{z}' , define $\mathcal{E} \downarrow \mathbf{z} \cup \mathbf{z}'$ - the **sample-domain space of error functions**. This is \mathcal{E} where the domain of the error functions has been limited to $\mathbf{z} \cup \mathbf{z}'$.
- The partition generalization error can be bounded in terms of the size of the sample-domain space of error functions:

$$g(L, \mathbf{z} \cup \mathbf{z}') \leq |\mathcal{E} \downarrow \mathbf{z} \cup \mathbf{z}'| \epsilon(n, n').$$

VC theory

Sufficient condition

- If \mathcal{F} has low VC dimension then the size of $\mathcal{E} \downarrow \mathbf{z} \cup \mathbf{z}'$ is small and therefore $g(L, \mathbf{z} \cup \mathbf{z}')$ is small for any L with range \mathcal{F} .
- This works well if \mathcal{F} is, say, the space of planes, or balls, or boxes in d dimensions.

VC theory

Necessary condition

A low VC dimension for \mathcal{F} is also a necessary condition for low generalization error in the sense that if \mathcal{F} has a high VC dimension then there exists a learning algorithm L (namely ERM) which will have high generalization error.

Learning in high VC dimensions

- By abandoning ERM, good generalizability can be guaranteed in high VC dimensions.
- E.g., SVM: in a high or infinite dimensional linear space, only select separating planes that have a large margin.

Bounding the partition generalization error in high VC dimensions

The sample-conditioned space of classifiers

Define the **sample-conditioned space of classifiers**:

$$\mathcal{F}_L(\mathbf{z} \cup \mathbf{z}') = \{L(\bar{\mathbf{z}}) : \bar{\mathbf{z}} \subset \mathbf{z} \cup \mathbf{z}'\}.$$

If $|\mathcal{F}_L(\mathbf{z} \cup \mathbf{z}')|$ is small then $g(L, \mathbf{z} \cup \mathbf{z}')$ is small:

$$g(L, \mathbf{z} \cup \mathbf{z}') \leq |\mathcal{F}_L(\mathbf{z} \cup \mathbf{z}')| \epsilon(n, n').$$

Applications

d -determined classification algorithm

A learning algorithm L is d -**determined** if there exists some function $L_0 : \mathcal{Z}^d \rightarrow \mathcal{F}$ such that

$$L(\mathbf{z}) = L_0(z_{i_1}(\mathbf{z}), \dots, z_{i_d}(\mathbf{z})).$$

When L is d -determined,

$$|\mathcal{F}_L(\mathbf{z} \cup \mathbf{z}')| = \binom{n+n'}{d} \ll \binom{n+n'}{n}.$$

Applications

Edited nearest neighbors

Of n training points, pick d (in any way the learner chooses) and use those for $[k-]$ nearest neighbor classification.

Applications

Support point machines

Let $s(x_1, x_2)$ be an arbitrary similarity function mapping pairs in \mathcal{X} to $[-1, 1]$.

Let

$$L(\mathbf{z})(x) = \operatorname{sgn} \sum_{i=1}^n y_i \alpha_i s(x_i, x),$$

where $\sum_i \alpha_i = 1, \alpha_i > 0$.

Applications

Support point machines, continued

Define a d -determined L' as follows:

$$L'(\mathbf{z})(x) = \operatorname{sgn} \frac{1}{d} \sum_{j=1}^d y_{K_j} s(x_{K_j}, x),$$

where K_1, \dots, K_d are IID sampled from $\{1, \dots, n\}$ with probabilities $\alpha_1, \dots, \alpha_n$.

When L has a large margin, it is well approximated by L' . The low generalization error of L' then implies low generalization error of L .

Bounding the partition generalization error

A finite classifier space

When \mathcal{E} is finite,

$$\begin{aligned}g(L, \mathbf{z} \cup \mathbf{z}') &\leq \mathbb{E}_{\text{Partition}} \max_{e \in \mathcal{E}} e^{n'}(\bar{\mathbf{z}}') - e^n(\bar{\mathbf{z}}) \\ &\leq \sum_{e \in \mathcal{E}} \mathbb{E} \left| e^{n'}(\bar{\mathbf{z}}') - e^n(\bar{\mathbf{z}}) \right| \\ &\leq |\mathcal{E}| \epsilon(n, n').\end{aligned}$$

Combining VC theory with sample-set conditioning

Sample-set conditioning can be combined with VC theory.

Define the **sample-conditioned space of error functions**:

$$\mathcal{E}_L(\mathbf{z}, \mathbf{z}') = \{e_f : f \in \mathcal{F}_L(\mathbf{z} \cup \mathbf{z}')\}.$$

Then the partition generalization error can be bounded in terms of the **sample-domain, sample-conditioned space of error functions**:

If $\mathcal{E}_L(\mathbf{z}, \mathbf{z}') \downarrow \mathbf{z}, \mathbf{z}'$ is small then $g(L, \mathbf{z} \cup \mathbf{z}')$ is small.