# A Learning Generalization Bound with an Application to Sparse-Representation Classifiers*

YORAM GAT
*University of California, Berkeley, USA*

**Editor:** Robert Schapire

**Abstract.** A classifier is said to have good generalization ability if it performs on test data almost as well as it does on the training data. The main result of this paper provides a sufficient condition for a learning algorithm to have good finite sample generalization ability. This criterion applies in some cases where the set of all possible classifiers has infinite VC dimension. The result is applied to prove the good generalization ability of support vector machines by a exploiting a sparse-representation property.

**Keywords:** generalization ability, sparsity, support vector machines, VC dimension, perceptron algorithm

## 1. Introduction

I consider the classical problem of learning a classifier from examples which can be formalized as follows: Let $Z_i = (X_i, Y_i)$, $i = 1, 2, \ldots$ be iid random variables taking values in $\mathcal{Z} = \mathcal{X} \times \{-1, +1\}$. The problem is predicting $Y_{l+1}$ given $X_1, \ldots, X_{l+1}$ and $Y_1, \ldots, Y_l$.

The solution to the problem is a map $L : \mathcal{Z}^l \to \mathcal{F}$, where $\mathcal{F}$ is a space of classifier functions, i.e., each $f \in \mathcal{F}$ is a function $f : \mathcal{X} \to \{-1, +1\}$. Thus the prediction is $Y_{l+1}^* = f_L(X_{l+1})$ where $f_L = L(Z_1, \ldots, Z_l)$.

The quality of the classifier $f_L$ may be measured using its expected error rate (also called expected risk):

$$R = \mathbf{P}(Y_{l+1}^* \neq Y_{l+1}).$$

The solution $L$ is usually geared toward finding a function which has low empirical error rate (also called empirical risk):

$$R_{\text{emp}} = \frac{1}{2l} \sum_{i=1}^{l} |f_L(X_i) - Y_i|.$$

Therefore, it is often desirable to be able to obtain bounds for the difference between the empirical and the expected error rates. The behavior of the difference will depend on the underlying, unknown probability measure. The term generalization ability is used to describe

the worst-case behavior of the difference between the empirical and expected error rate for a specific algorithm. The smaller the probability for a large difference, the better is the generalization ability of the algorithm.

One map $L$ commonly used is

$$L(Z_1, \ldots, Z_l) = \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{l} |Y_i - f(X_i)|.$$

This is known as the Empirical Risk Minimization (ERM) method. It has been shown that the generalization ability of the algorithm can be determined by using the VC dimension of the set of functions $\mathcal{F}$ (Vapnik, 1998).

Other learning algorithms use maps of the form

$$L(Z_1, \ldots, Z_l) = \arg\min_{f \in \tilde{L}(Z_1, \ldots, Z_l)} \sum_{i=1}^{l} |Y_i - f(X_i)|,$$

where $\tilde{L}$ is an auxiliary map $\tilde{L} : \mathcal{Z}^l \to 2^{\mathcal{F}}$. I call this type of algorithms Restricted Empirical Risk Minimization (RERM) rules.

## 2. The main result

The following theorem guarantees the generalization ability of certain learning algorithms even when $\mathcal{F}$ has an infinite VC dimension:

**Theorem 1.** *Denote*

$$\bar{L}(z_1, \ldots, z_{2l}) = \left\{ L\big(z_{i(1)}, \ldots, z_{i(l)}\big) : \textit{the } i(j)\textit{'s are } l \textit{ distinct indices in the range} \atop 1, \ldots, 2l \right\}.$$

*If*

$$\sup_{z_1, \ldots, z_{2l} \in \mathcal{Z}} |\bar{L}(z_1, \ldots, z_{2l})| = c(l),$$

*then*

$$\mathbf{P}(|R - R_{\text{emp}}| > \epsilon) < 2c(l) \exp{-(l\epsilon^2 - 2\epsilon)}.$$

**Proof:** Since for any Binomial variable, $B$, $\mathbf{P}(B > \mathbf{E}B + 1) < 0.5$, it is enough to bound

$$p_{\epsilon'} = \mathbf{P}\left( \left| \frac{1}{2l} \sum_{i=l+1}^{2l} |f_L(X_i) - Y_i| - R_{\text{emp}} \right| > \epsilon' \right),$$

where $\epsilon' = \epsilon - \frac{1}{l}$. This is done by conditioning on the values of $z_i$, $i = 1, \ldots, 2l$ and then taking the expectation over the different possible orderings.

To simplify the formulas, I use below $\Delta^f(z)$ as shorthand for

$$\frac{1}{2}|f(x) - y|,$$

where $z = (x, y)$. Thus $\Delta^f(z)$ is either 0 or 1, and

$$p_{\epsilon'} = \mathbf{E}\frac{1}{(2l)!}\sum_\sigma \mathbf{1}\left(\left|\sum_{i=1}^{l}\Delta^{f_L}\left(Z_{\sigma(i)}\right) - \sum_{i=l+1}^{2l}\Delta^{f_L}\left(Z_{\sigma(i)}\right)\right| > l\epsilon'\right).$$

Here, as below, $\sum_\sigma$ means summing over all permutations of the numbers $1, \ldots, 2l$.

$$p_{\epsilon'} \leq \mathbf{E}\frac{1}{(2l)!}\sum_\sigma \mathbf{1}\left(\sup_{f\in\bar{L}(Z_1,\ldots,Z_{2l})}\left|\sum_{i=1}^{l}\Delta^f\left(Z_{\sigma(i)}\right) - \sum_{i=l+1}^{2l}\Delta^f\left(Z_{\sigma(i)}\right)\right| > l\epsilon'\right)$$

$$\leq \mathbf{E}\frac{1}{(2l)!}\sum_\sigma \sum_{f\in\bar{L}(Z_1,\ldots,Z_{2l})}\mathbf{1}\left(\left|\sum_{i=1}^{l}\Delta^f\left(Z_{\sigma(i)}\right) - \sum_{i=l+1}^{2l}\Delta^f\left(Z_{\sigma(i)}\right)\right| > l\epsilon'\right)$$

$$\leq \mathbf{E}\sum_{f\in\bar{L}(Z_1,\ldots,Z_{2l})}\frac{1}{(2l)!}\sum_\sigma \mathbf{1}\left(\left|\sum_{i=1}^{l}\Delta^f\left(Z_{\sigma(i)}\right) - \sum_{i=l+1}^{2l}\Delta^f\left(Z_{\sigma(i)}\right)\right| > l\epsilon'\right)$$

$$\leq c(l)\exp{-l\epsilon'^2}$$

$$\leq c(l)\exp{-(l\epsilon^2 - 2\epsilon)}.$$

The bound for the fraction of permutations giving a difference greater than $\epsilon'$ was calculated by Vapnik (1998, Sec. 4.13).                                              □

The proof above follows the argument of Theorem 4.1 of Vapnik (1998) which deals with the generalization ability of ERM algorithms. The main difference is the reference to the random set $\bar{L}(Z_1, \ldots, Z_{2l})$ rather than to a fixed set of functions. Two variants of the result stated in Theorem 4.1 of Vapnik (1998) are Theorem 4.2 of Vapnik (1998) and the main result of Devroye (1982). Both can be adapted and proven for the setup here in a manner similar to that of Theorem 1. The first variant gives better bounds when the empirical error rate is small, and the other gives a better rate of convergence when $c(l)$ is polynomial.

The next result follows immediately from Theorem 1:

**Corollary 1.**   *For maps L of the RERM type, the bound of Theorem* 1 *holds provided that*

$$\sup_{z_1,\ldots,z_{2l}\in\mathcal{Z}}|\bar{\bar{L}}(z_1, \ldots, z_{2l})| = c(l),$$

*with*

$$\bar{\bar{L}}(z_1, \ldots, z_{2l}) = \bigcup_{i(1),\ldots,i(l)}\tilde{L}\left(z_{i(1)}, \ldots, z_{i(l)}\right).$$

*where the $i(j)$'s are $l$ distinct indices in the range* $1, \ldots, 2l$.

*Example* ($r$-sparse rules).    Corollary 1 can be used to obtain a non-trivial generalization property for any rule of the RERM type where $\tilde{L}$ is of the form

$$\tilde{L}(z_1, \ldots, z_l) = \left\{ f_{z_{j(1)}, \ldots, z_{j(r)}} : j(i) \in \{1, \ldots, l\}, \ i = 1, \ldots, r \right\},$$

since for any map $L$ of this type, $|\tilde{\bar{L}}(z_1, \ldots, z_{2l})| \leq (2l)^r$. Below, I refer to such rules as sparse, or $r$-sparse rules.

Similar bounds for sparse rules were obtained by Littlestone and Warmuth (quoted in Floyd and Warmuth (1995)), and by Graepel et al. (2000). These bounds were obtained using an approach which is quite different than the one used here.

## 3.   The support-vector setup

The support-vector machine (SVM) (Vapnik, 1998) creates a linear discriminant classifier in a ball within a high dimensional, or an infinite dimensional, Euclidean space:

$$\mathcal{X} = \{x \in \mathcal{R}^n : |x| \leq 1\},$$
$$\mathcal{F} = \{f_{a,b}(x) = \mathrm{sign}(a \cdot x + b) : a \in \mathcal{R}^n, b \in \mathcal{R}, |a| = 1\}.$$

To put the definition of an SVM into the framework presented here, I introduce the following definitions:

*Definition 1.*    Let $S(x_1, \ldots, x_l, t_1, \ldots, t_l)$, $x_i \in \mathcal{X}$, $t_i \in \{-1, +1\}$ be the set of classifiers $f_{a,b} \in \mathcal{F}$ such that for all $i = 1, \ldots, l$, $f_{a,b}(x_i) = 1$ iff $t_i = 1$.

In other words, the set $S(x_1, \ldots, x_l, t_1, \ldots, t_l)$ is the set of classifiers which predict $Y = t_i$ when presented with $X = x_i$, for all $i = 1, \ldots, l$.

*Definition 2.*    The margin of a classifier $f_{a,b} \in \mathcal{F}$ with respect to a set of points $x_1, \ldots, x_l \in \mathcal{X}$ is defined as

$$\min_{i=1,\ldots,l} |a \cdot x_i + b|.$$

The maximum margin classifier (MMC) is the member, $f_{a,b}$, of the set $S$ with the property that its margin is the largest in the set.

The margin of the MMC is denoted by $\mathbf{marg}(x_1, \ldots, x_l, t_1, \ldots, t_l)$.

Using the definitions above, the SVM can now be defined as a RERM type rule with:

$$\tilde{L}(z_1, \ldots, z_l) = \{f_{a,b} = s(x_1, \ldots, x_l, t_1, \ldots, t_l) : t_i \in \{-1, +1\}, i = 1, \ldots, l,$$
$$\mathbf{marg}(x_1, \ldots, x_l, t_1, \ldots, t_l) \geq h\},$$

where $s(x_1, \ldots, x_l, t_1, \ldots, t_l)$ is some member of $S(x_1, \ldots, x_l, t_1, \ldots, t_l)$ and $h$ is some fixed constant.

The set $\tilde{L}(z_1, \ldots, z_l)$ may or may not contain a representative from the set $S(x_1, \ldots, x_l, y_1, \ldots, y_l)$. If it does contain such a representative, $f$, then $f$ will have zero empirical error rate, and therefore

$$L(z_1, \ldots, z_l) = f$$

will hold. If such a representative is not in $\tilde{L}(z_1, \ldots, z_l)$ then

$$L(z_1, \ldots, z_l) = s(x_1, \ldots, x_l, t_1, \ldots, t_l)$$

for some $t_1, \ldots, t_l$ and the empirical error rate of the algorithm will be equal to the cardinality of the set $\{i : t_i \neq y_i\}$.

Based on heuristic appeal and experimental results, $s$ is usually chosen to be equal to the MMC. Here, however, I propose a different way to select a representative, for which the generalization ability can be determined. Note that the empirical risk achieved is the same for any choice of a representative.

The algorithm below, known as the perceptron algorithm (Minsky & Papert, 1998), may be used to obtain a member of $S(x_1, \ldots, x_l, t_1, \ldots, t_l)$. Let the representative, $s$, be the one produced by the algorithm. This algorithm had been previously considered in this context by Freund and Schapire (1998) and by Graepel et al. (2000).

- **Initialization:** Set $a \leftarrow 0, b \leftarrow 0, k \leftarrow 1$
- **Update:** If $t_k(a \cdot x_k + b) > 0$ then go to step **Loop**
- **Correction:** Set $a \leftarrow a + t_k x_k, b \leftarrow b + t_k$
- **Loop:** If a **Correction** step was not carried out in the last $l$ loops, stop. Otherwise, set $k \leftarrow k + 1 (\text{mod } l)$ and go to step **Update**

The Perceptron Convergence Theorem (Minsky & Papert, 1988) states that if the points $x_i$ all lie inside the unit sphere, and

$$\mathbf{marg}(x_1, \ldots, x_l, t_1, \ldots, t_l) \geq h,$$

then the algorithm will execute at most $\lfloor 1/h^2 \rfloor$ corrections, after which the resulting $a, b$ parameters will provide a member $f_{a,b}$ of $S(x_1, \ldots, x_l, t_1, \ldots, t_l)$. By construction the resulting classifier is $r$-sparse with $r \leq \lfloor 1/h^2 \rfloor$.

Applying the bound for $r$-sparse rules leads to the following conclusion: For any fixed $h$, if a support-vector method is employed and a classifier with a margin of $h$ and empirical error rate $R_{\text{emp}}$ is found, then there exists an $r(h)$-sparse classifier for which the following statement holds:

$$\mathbf{P}(R > R_{\text{emp}} + \epsilon) < 2(2l)^{\lfloor 1/h^2 \rfloor} \exp -(l\epsilon^2 - 2\epsilon). \tag{1}$$

The perceptron algorithm can be used to obtain such a classifier.

An important point about the perceptron algorithm is that it can be executed without reference to the training vectors themselves but rather making use only of the inner products between training vectors. The importance of this property stems from the fact that often in applications of the support vector machine calculating inner products between training vectors is feasible, but any explicit representation of the vectors is prohibitively expensive.

Equation (1) can be converted into a $1 - \delta$ upper confidence bound. With probability of at least $1 - \delta$, the following inequality holds:

$$ R < R_{\text{emp}} + \sqrt{\frac{1}{l}\left(\frac{\log 2l}{h^2} + \log \frac{1}{\delta} + \log 2e^2\right)}. \tag{2} $$

The upper confidence bound (2) holds under the assumption that $h$ is fixed in advance. It is common practice, however, to have $h$ random. This is, for example, the case when the empirical error rate is pre-specified (e.g. zero).

A result suitable for the case of a random $h$ will have the form of simultaneous upper confidence bounds for $r = \frac{1}{h^2} = 1, \ldots, l$. This is obtained by simply replacing $\delta$ by $\delta/l$ in (2), obtaining a $1 - \delta$ an upper confidence bound of the following form:

$$ R < R_{\text{emp}} + \sqrt{\frac{1}{l}\left(\frac{\log 2l}{h^2} + \log \frac{l}{\delta} + \log 2e^2\right)}. \tag{3} $$

Since insisting on a pre-specified empirical error rate may lead to a large upper confidence bound, different procedures may be followed. One such procedure would be an adaptation of the perceptron algorithm:

– **Initialization:** Set $a(0) \leftarrow 0, b(0) \leftarrow 0, k \leftarrow 1, j \leftarrow 0, R(0) \leftarrow l$
– **Update:** If $y_k(a(j) \cdot x_k + b(j)) > 0$ go to step **Loop**
– **Correction:** Set $a(j+1) \leftarrow a(j) + y_k x_k, b(j+1) \leftarrow b(j) + y_k, j \leftarrow j+1, R(j) \leftarrow |\{i : y_k(a \cdot x_k + b) \leq 0\}|$
– **Loop:** If $R(j) = 0$ or $j = l$, go to step **Optimization**. Otherwise, set $k \leftarrow k + 1 (\text{mod } l)$ and go to step **Update**
– **Optimization:** Set

$$ j^* = \arg\min_{0 \leq i \leq j} R(i) + \sqrt{\frac{1}{l}\left(i \log 2l + \log \frac{l}{\delta} + \log 2e^2\right)}. $$

Set $f_L = f_{a(j^*),b(j^*)}$. Stop

At the termination of the algorithm, $f_L$ is a classifier with empirical error rate $R(j^*)$, and which with probability of at least $1 - \delta$ has expected error rate no greater than

$$ R(j^*) + \sqrt{\frac{1}{l}\left(j^* \log 2l + \log \frac{l}{\delta} + \log 2e^2\right)}. $$

## 4. Experimental results

The use of variants of the perceptron algorithm in the support vector context had been previously suggested and implemented by Freund and Schapire (1998). They carried out experiments using the perceptron algorithm for classifying images of handwritten digits and report error rates which are somewhat larger than those obtained with maximum margin classifiers.

## Acknowledgments

I thank Peter Bickel for pointing out the problem which resulted in this paper, for his valuable comments and for having reviewed the paper.

## References

Devroye, L. (1982). Bounds for the uniform deviation of empirical measures. *Journal of Multivariate Analysis, 12*, 72–79.

Floyd, S. & Warmuth M. (1995). Sample compression, learnability, and the Vapnik-Chervonenekis dimension. *Machine Learning Journal, 21*, 269–304.

Freund, Y. & Schapire, R. E. (1998). Large margin classification using the perceptron algorithm. In *COLT '98: Proceedings of the Eleventh Annual Conference on Computational Learning Theory*.

Graepel, T., Herbrich, R., & Shawe-Taylor, J. (2000). Generalisation error bounds for sparse linear classifiers. In *COLT 2000: Proceedings of the Thirteenth Annual Conference on Computational Learning Theory*.

Minsky, M. L. & Papert, S. A. (1988). *Perceptrons*. The MIT Press, Cambridge, MA.

Vapnik, V. N. (1998). *Statistical Learning Theory*. John Wiley & Sons, Inc., New York City, NY.