

Binary classification

Maxim Raginsky

October 15, 2013

The problem of binary classification can be stated as follows. We have a random couple $Z = (X, Y)$, where $X \in \mathbb{R}^d$ is called the *feature vector* and $Y \in \{-1, 1\}$ is called the *label*¹. In the spirit of the model-free framework, we assume that the relationship between the features and the labels is stochastic and described by an unknown probability distribution $P \in \mathcal{P}(Z)$, where $Z = \mathbb{R}^d \times \{-1, 1\}$.

As usual, we consider the case when we are given an i.i.d. sample of length n from P . The goal is to learn a *classifier*, i.e., a mapping $g: \mathbb{R}^d \rightarrow \{-1, 1\}$, such that the probability of classification error, $\mathbb{P}(g(X) \neq Y)$, is small. As we have seen before, the optimal choice is the *Bayes classifier*

$$g^*(x) \triangleq \begin{cases} 1, & \text{if } \eta(x) > 1/2 \\ -1, & \text{otherwise} \end{cases} \quad (1)$$

where $\eta(x) \triangleq \mathbb{P}[Y = 1|X = x]$ is the *regression function*. However, since we make no assumptions on P , in general we cannot hope to learn the Bayes classifier g^* . Instead, we focus on a more realistic goal: We fix a collection \mathcal{G} of classifiers and then use the training data to come up with a hypothesis $\hat{g}_n \in \mathcal{G}$, such that

$$\mathbb{P}(\hat{g}_n(X) \neq Y) \approx \inf_{g \in \mathcal{G}} \mathbb{P}(g(X) \neq Y)$$

with high probability.

By way of notation, let us write $L(g)$ for the classification error of g , i.e., $L(g) \triangleq \mathbb{P}(g(X) \neq Y)$, and let $L^*(\mathcal{G})$ denote the smallest classification error attainable over \mathcal{G} :

$$L^*(\mathcal{G}) \triangleq \inf_{g \in \mathcal{G}} L(g).$$

We will assume that a minimizing $g^* \in \mathcal{G}$ exists. For future reference, we note that

$$L(g) = \mathbb{P}(g(X) \neq Y) \leq \mathbb{P}(Y g(X) < 0). \quad (2)$$

Warning: In what follows, we will use C or c to denote various absolute constants; their values may change from line to line.

1 Learning linear discriminant rules

One of the simplest classification rules (and one of the first to be studied) is a *linear discriminant rule*: given a nonzero vector $w = (w^{(1)}, \dots, w^{(d)}) \in \mathbb{R}^d$ and a scalar $b \in \mathbb{R}$, let

$$g(x) \equiv g_{w,b}(x) \triangleq \begin{cases} 1, & \text{if } \langle w, x \rangle + b > 0 \\ -1, & \text{otherwise} \end{cases} \quad (3)$$

¹The reason why we chose $\{-1, 1\}$, rather than $\{0, 1\}$, for the label space is merely convenience.

Let \mathcal{G} be the class of all such linear discriminant rules as w ranges over all nonzero vectors in \mathbb{R}^d and b ranges over all reals: $\mathcal{G} = \{g_{w,b} : w \in \mathbb{R}^d \setminus \{0\}, b \in \mathbb{R}\}$.

Given the training sample Z^n , let $\hat{g}_n \in \mathcal{G}$ be the output of the ERM algorithm, i.e.,

$$\hat{g}_n \triangleq \operatorname{argmin}_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{g(X_i) \neq Y_i\}}.$$

In other words, \hat{g}_n is any classifier of the form (3) that minimizes the number of misclassifications on the training sample. Then we have the following:

Theorem 1. *There exists an absolute constant $C > 0$, such that for any $n \in \mathbb{N}$ and any $\delta \in (0, 1)$, the bound*

$$L(\hat{g}_n) \leq L^*(\mathcal{G}) + C \sqrt{\frac{d+1}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}} \quad (4)$$

holds with probability at least $1 - \delta$.

Proof. A standard argument leads to the bound

$$L(\hat{g}_n) \leq L^*(\mathcal{G}) + 2\Delta_n(Z^n), \quad (5)$$

where

$$\Delta_n(Z^n) \triangleq \sup_{g \in \mathcal{G}} |L(g) - L_n(g)|$$

is the uniform deviation and $L_n(g)$ denotes the *empirical classification error* of g on Z^n :

$$L_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{g(X_i) \neq Y_i\}},$$

which is the fraction of incorrectly labeled points in the training sample Z^n . Consider a classifier $g \in \mathcal{G}$ and define the set

$$C_g \triangleq \{(x, y) \in \mathbb{R}^d \times \{-1, 1\} : y \cdot (\langle w, x \rangle + b) \leq 0\}.$$

Then it is easy to see that

$$L(g) = P(C_g) \quad \text{and} \quad L_n(g) = P_n(C_g),$$

where, as before,

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Z_i} = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$$

is the empirical distribution of the sample Z^n . Let \mathcal{C} denote the collection of all sets of the form $C = C_g$ for some $g \in \mathcal{G}$. Then

$$\Delta_n(Z^n) = \sup_{C \in \mathcal{C}} |P_n(C) - P(C)|.$$

Let $\mathcal{F} = \mathcal{F}_{\mathcal{C}}$ denote the class of indicator functions of the sets in \mathcal{C} : $\mathcal{F}_{\mathcal{C}} = \{\mathbf{1}_{\{C \in \mathcal{C}\}} : C \in \mathcal{C}\}$. Then we know that, with probability at least $1 - \delta$,

$$\Delta_n(Z^n) \leq 2\mathbb{E}R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{\log(1/\delta)}{2n}}, \quad (6)$$

where $R_n(\mathcal{F}(Z^n))$ is the Rademacher average of the projection of \mathcal{F} onto the sample Z^n . Now,

$$\begin{aligned}\mathcal{F}(Z^n) &= \{(f(Z_1), \dots, f(Z_n)) : f \in \mathcal{F}\} \\ &= \{\mathbf{1}_{\{Z_1 \in C\}}, \dots, \mathbf{1}_{\{Z_n \in C\}} : C \in \mathcal{C}\}.\end{aligned}$$

Therefore, if we prove that \mathcal{C} is a VC class, then

$$R_n(\mathcal{F}(Z^n)) \leq C \sqrt{\frac{V(\mathcal{C})}{n}}.$$

But this follows from the fact that any $C \in \mathcal{C}$ has the form

$$C = \left\{ (x, y) \in \mathbb{R}^d \times \{-1, 1\} : \sum_{j=1}^d w^{(j)} y x^{(j)} + b y \leq 0 \right\}$$

for some $w \in \mathbb{R}^d \setminus \{0\}$ and some $b \in \mathbb{R}$, and the functions $(x, y) \mapsto y$ and $(x, y) \mapsto y x^{(j)}$, $1 \leq j \leq d$, span a linear space of dimension no greater than $d + 1$. Hence, $V(\mathcal{C}) \leq d + 1$, so that

$$R_n(\mathcal{F}(Z^n)) \leq C \sqrt{\frac{V(\mathcal{C})}{n}} \leq C \sqrt{\frac{d+1}{n}}.$$

Combining this with (5) and (6), we see that (4) holds with probability at least $1 - \delta$. \square

1.1 Generalized linear discriminant rules

In the same vein, we may consider classification rules of the form

$$g(x) = \begin{cases} 1, & \text{if } \sum_{j=1}^k w^{(j)} \psi_j(x) + b > 0 \\ -1, & \text{otherwise} \end{cases} \quad (7)$$

where k is some positive integer (not necessarily equal to d), $w = (w^{(1)}, \dots, w^{(k)}) \in \mathbb{R}^k$ is a nonzero vector, $b \in \mathbb{R}$ is an arbitrary scalar, and $\Psi = \{\psi_j : \mathbb{R}^d \rightarrow \mathbb{R}\}_{j=1}^k$ is some fixed “dictionary” of real-valued functions on \mathbb{R}^d . For a fixed Ψ , let \mathcal{G} denote the collection of all classifiers of the form (7) as w ranges over all nonzero vectors in \mathbb{R}^k and b ranges over all reals. Then the ERM rule is, again, given by

$$\hat{g}_n \triangleq \inf_{g \in \mathcal{G}} L_n(g) \equiv \inf_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{g(X_i) \neq Y_i\}}.$$

The following result can be proved pretty much along the same lines as Theorem 1:

Theorem 2. *There exists an absolute constant $C > 0$, such that for any $n \in \mathbb{N}$ and any $\delta \in (0, 1)$, the bound*

$$L(\hat{g}_n) \leq L^*(\mathcal{G}) + C \sqrt{\frac{k+1}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}} \quad (8)$$

holds with probability at least $1 - \delta$.

1.2 Two fundamental issues

As Theorems 1 and 2 show, the ERM algorithm applied to the collection of all (generalized) linear discriminant rules is guaranteed to work well in the sense that the classification error of the output hypothesis will, with high probability, be close to the optimum achievable by any discriminant rule with the given structure. The same argument extends to any collection of classifiers \mathcal{G} , for which the “error sets” $\{(x, y) : y \cdot g(x) \leq 0\}$, $g \in \mathcal{G}$, form a VC class of dimension much smaller than the sample size n . In other words, with high probability the difference

$$L(\hat{g}_n) - L^*(\mathcal{G}) = L(\hat{g}_n) - \inf_{g \in \mathcal{G}} L(g)$$

will be small. However, precisely because the VC dimension of \mathcal{G} cannot be too large, the approximation properties of \mathcal{G} will be limited. Another problem is computational. For instance, the problem of finding an empirically optimal linear discriminant rule is NP-hard. In other words, unless P is equal to NP, there is no hope of coming up with an efficient ERM algorithm for linear discriminant rules that would work for all feature space dimensions d . If d is fixed, then it is possible to enumerate all projections of a given sample Z^n onto the class of indicators of all halfspaces in $O(n^{d-1} \log n)$ time, which allows for an exhaustive search for an ERM solution, but the usefulness of this naive approach is limited to $d < 5$.

2 Risk bounds for combined classifiers via surrogate loss functions

One way to sidestep the above approximation-theoretic and computational issues is to replace the 0–1 Hamming loss function that gives rise to the probability of error criterion with some other loss function. What we gain is the ability to bound the performance of various complicated classifiers built up by combining simpler *base classifiers* in terms of the complexity (e.g. the VC dimension) of the collection of the base classifiers, as well as considerable computational advantages, especially if the problem of minimizing the empirical surrogate loss turns out to be a convex programming problem. What we lose, though, is that, in general, we will not be able to compare the generalization error of the learned classifier to the minimum classification risk. Instead, we will have to be content with the fact that the generalization error will be close to the smallest *surrogate loss*.

We will consider classifiers of the form

$$g_f(x) = \text{sgn } f(x) \equiv \begin{cases} 1, & \text{if } f(x) \geq 0 \\ -1, & \text{otherwise} \end{cases} \quad (9)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is some function. From (2) we have

$$L(g_f) = \mathbb{P}(g_f(X) \neq Y) \leq \mathbb{P}(Y g_f(X) < 0) = \mathbb{P}(Y f(X) < 0).$$

From now on, when dealing with classifiers of the form (9), we will write $L(f)$ instead of $L(g_f)$ to keep the notation simple. Now we introduce the notion of a surrogate loss function.

Definition 1. A surrogate loss function is any function $\varphi : \mathbb{R} \rightarrow \mathbb{R}_+$, such that

$$\varphi(x) \geq \mathbf{1}_{\{x > 0\}}. \quad (10)$$

Some examples of commonly used surrogate losses:

1. Exponential, $\varphi(x) = e^x$
2. Logit, $\varphi(x) = \log_2(1 + e^x)$
3. Hinge loss, $\varphi(x) = (x + 1)_+ \equiv \max\{x + 1, 0\}$

Let φ be a surrogate loss. Then for any $(x, y) \in \mathbb{R}^d \times \{-1, 1\}$ and any $f : \mathbb{R}^d \rightarrow \mathbb{R}$ we have

$$yf(x) < 0 \quad \Rightarrow \quad \varphi(-yf(x)) \geq \mathbf{1}_{\{-yf(x) > 0\}} = \mathbf{1}_{\{yf(x) < 0\}}. \quad (11)$$

Therefore, defining the φ -risk of f by

$$A_\varphi(f) \triangleq \mathbb{E}[\varphi(-Yf(X))]$$

and its empirical version

$$A_{\varphi,n}(f) \triangleq \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i)),$$

we see from (11) that

$$L(f) \leq A_\varphi(f) \quad \text{and} \quad L_n(f) \leq A_{\varphi,n}(f). \quad (12)$$

Now that these preliminaries are out of the way, we can state and prove the basic surrogate loss bound:

Theorem 3. Consider any learning algorithm $\mathcal{A} = \{\mathcal{A}_n\}_{n=1}^\infty$, where, for each n , the mapping \mathcal{A}_n receives the training sample $Z^n = (Z_1, \dots, Z_n)$ as input and produces a function $\hat{f}_n : \mathbb{R}^d \rightarrow \mathbb{R}$ from some class \mathcal{F} . Suppose that \mathcal{F} and the surrogate loss φ are chosen so that the following conditions are satisfied:

1. There exists some constant $B > 0$ such that

$$\sup_{(x,y) \in \mathbb{R}^d \times \{-1,1\}} \sup_{f \in \mathcal{F}} \varphi(-yf(x)) \leq B$$

2. There exists some constant $M_\varphi > 0$ such that φ is M_φ -Lipschitz, i.e.,

$$|\varphi(u) - \varphi(v)| \leq M_\varphi |u - v|, \quad \forall u, v \in \mathbb{R}.$$

Then for any n and any $\delta \in (0, 1)$ the following bound holds with probability at least $1 - \delta$:

$$L(\hat{f}_n) \leq A_{\varphi,n}(\hat{f}_n) + 4M_\varphi \mathbb{E}R_n(\mathcal{F}(X^n)) + B \sqrt{\frac{\log(1/\delta)}{2n}}. \quad (13)$$

Proof. Using (12), we can write

$$\begin{aligned} L(\hat{f}_n) &\leq A_\varphi(\hat{f}_n) \\ &= A_{\varphi,n}(\hat{f}_n) + A_\varphi(\hat{f}_n) - A_{\varphi,n}(\hat{f}_n) \\ &\leq A_{\varphi,n}(\hat{f}_n) + \sup_{f \in \mathcal{F}} |A_\varphi(f) - A_{\varphi,n}(f)|. \end{aligned}$$

Now let \mathcal{H} denote the class of functions $h : \mathbb{R}^d \times \{-1, 1\} \rightarrow \mathbb{R}$ of the form $h(x, y) = -yf(x)$, $f \in \mathcal{F}$. Then

$$\begin{aligned} \sup_{f \in \mathcal{F}} |A_\varphi(f) - A_{\varphi, n}(f)| &= \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\varphi(-Yf(X))] - \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i)) \right| \\ &= \sup_{h \in \mathcal{H}} |P(\varphi \circ h) - P_n(\varphi \circ h)|, \end{aligned}$$

where $\varphi \circ h(z) \triangleq \varphi(h(z))$ for every $z = (x, y) \in \mathbb{R}^d \times \{-1, 1\}$. Let

$$\Delta_n(Z^n) \triangleq \sup_{h \in \mathcal{H}} |P(\varphi \circ h) - P_n(\varphi \circ h)|.$$

We can now use the familiar symmetrization argument to bound

$$\mathbb{E}\Delta_n(Z^n) \leq 2\mathbb{E}R_n(\varphi \circ \mathcal{H}(Z^n)), \quad (14)$$

where $\varphi \circ \mathcal{H}$ denotes the class of all functions of the form $\varphi \circ h$, $h \in \mathcal{H}$, and

$$\varphi \circ \mathcal{H}(Z^n) \triangleq \{(\varphi \circ h(Z_1), \dots, \varphi \circ h(Z_n)) : h \in \mathcal{H}\} = \{(\varphi(h(Z_1)), \dots, \varphi(h(Z_n))) : h \in \mathcal{H}\} \quad (15)$$

is the projection of $\varphi \circ \mathcal{H}$ onto the random sample Z^n . We now use a very powerful result about the Rademacher averages called the *contraction principle*, which states the following: If $\mathcal{A} \subset \mathbb{R}^n$ is a bounded set and $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is an M_φ -Lipschitz function, then

$$R_n(\varphi \circ \mathcal{A}) \leq 2M_\varphi R_n(\mathcal{A}), \quad (16)$$

where $\varphi \circ \mathcal{A} \triangleq \{(\varphi(a_1), \dots, \varphi(a_n)) : a = (a_1, \dots, a_n) \in \mathcal{A}\}$. (The proof of the contraction principle is somewhat involved, and we do not give it here.) From (15) we immediately see that

$$\varphi \circ \mathcal{H}(Z^n) = \varphi \circ [\mathcal{H}(Z^n)].$$

Therefore, applying (16) to $\mathcal{A} = \mathcal{H}(Z^n)$ and then using the resulting bound in (14), we obtain

$$\mathbb{E}\Delta_n(Z^n) \leq 4M_\varphi \mathbb{E}R_n(\mathcal{H}(Z^n)).$$

Furthermore, letting σ^n be an i.i.d. Rademacher tuple independent of Z^n , we have

$$\begin{aligned} R_n(\mathcal{H}(Z^n)) &= \frac{1}{n} \mathbb{E}_{\sigma^n} \left[\sup_{h \in \mathcal{H}} \left| \sum_{i=1}^n \sigma_i h(Z_i) \right| \right] \\ &= \frac{1}{n} \mathbb{E}_{\sigma^n} \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i Y_i f(X_i) \right| \right] \\ &= \frac{1}{n} \mathbb{E}_{\sigma^n} \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(X_i) \right| \right] \\ &\equiv R_n(\mathcal{F}(X^n)), \end{aligned}$$

which leads to

$$\mathbb{E}\Delta_n(Z^n) \leq 4M_\varphi \mathbb{E}R_n(\mathcal{F}(X^n)). \quad (17)$$

Now, since every function $\varphi \circ h$ is bounded between 0 and B , the function $\Delta_n(Z^n)$ has bounded differences with $c_1 = \dots = c_n = B/n$. Therefore, from (17) and from McDiarmid's inequality, we have for every $t > 0$ that

$$\mathbb{P}\left(\Delta_n(Z^n) \geq 4M_\varphi \mathbb{E}R_n(\mathcal{F}(X^n)) + t\right) \leq \mathbb{P}\left(\Delta_n(Z^n) \geq \mathbb{E}\Delta_n(Z^n) + t\right) \leq e^{-2nt^2/B^2}.$$

Choosing $t = B\sqrt{(2n)^{-1}\log(1/\delta)}$, we see that

$$\Delta_n(Z^n) \leq 4M_\varphi \mathbb{E}R_n(\mathcal{F}(X^n)) + B\sqrt{\frac{\log(1/\delta)}{2n}}$$

with probability at least $1 - \delta$. Therefore, since

$$L(\hat{f}_n) \leq A_{\varphi,h}(\hat{f}_n) + \Delta_n(Z^n),$$

we see that (13) holds with probability at least $1 - \delta$. \square

What the above theorem tells us is that the performance of the learned classifier \hat{f}_n is controlled by the Rademacher average of the class \mathcal{F} , and we can always arrange it to be relatively small. We will now look at several specific examples.

3 Weighted linear combination of classifiers

Let $\mathcal{G} = \{g : \mathbb{R}^d \rightarrow \{-1, 1\}\}$ be a class of *base classifiers* (not to be confused with *Bayes* classifiers!), and consider the class

$$\mathcal{F}_\lambda \triangleq \left\{ f = \sum_{j=1}^N c_j g_j : N \in \mathbb{N}, \sum_{j=1}^N |c_j| \leq \lambda; g_1, \dots, g_N \in \mathcal{G} \right\},$$

where $\lambda > 0$ is a tunable parameter. Then for each $f = \sum_{j=1}^N c_j g_j \in \mathcal{F}_\lambda$ the corresponding classifier g_f of the form (9) is given by

$$g_f(x) = \operatorname{sgn}\left(\sum_{j=1}^N c_j g_j(x)\right).$$

A useful way of thinking about g_f is that, upon receiving a feature $x \in \mathbb{R}^d$, it computes the outputs $g_1(x), \dots, g_N(x)$ of the N base classifiers from \mathcal{G} and then takes a weighted “majority vote” – indeed, if we had $c_1 = \dots = c_N = \lambda/N$, then $\operatorname{sgn}(g_f(x))$ would precisely correspond to taking the majority vote among the N base classifiers. Note, by the way, that the number of base classifiers is not fixed, and can be learned from the data.

Now, Theorem 3 tells us that the performance of any learning algorithm that accepts a training sample Z^n and produces a function $\hat{f}_n \in \mathcal{F}_\lambda$ is controlled by the Rademacher average $R_n(\mathcal{F}_\lambda(X^n))$. It turns out, moreover, that we can relate it to the Rademacher average of the base class \mathcal{G} . To start, note that

$$\mathcal{F}_\lambda = \lambda \cdot \operatorname{absconv} \mathcal{G},$$

where

$$\operatorname{absconv} \mathcal{G} = \left\{ \sum_{j=1}^N c_j g_j : N \in \mathbb{N}; \sum_{j=1}^N |c_j| \leq 1; g_1, \dots, g_N \in \mathcal{G} \right\}$$

is the absolute convex hull of \mathcal{G} . Therefore

$$R_n(\mathcal{F}_\lambda(X^n)) = \lambda \cdot R_n(\mathcal{G}(X^n)).$$

Now note that the functions in \mathcal{G} are binary-valued. Therefore, assuming that the base class \mathcal{G} is a VC class, we will have

$$R_n(\mathcal{G}(X^n)) \leq C \sqrt{\frac{V(\mathcal{G})}{n}}.$$

Combining these bounds with the bound of Theorem 3, we conclude that for any \hat{f}_n selected from \mathcal{F}_λ based on the training sample Z^n , the bound

$$L(\hat{f}_n) \leq A_{\varphi,n}(\hat{f}_n) + CM_\varphi \sqrt{\frac{V(\mathcal{G})}{n}} + B \sqrt{\frac{\log(1/\delta)}{2n}}$$

will hold with probability at least $1 - \delta$, where B is the uniform upper bound on $\varphi(-yf(x))$, $f \in \mathcal{F}_\lambda$, $(x, y) \in \mathbb{R}^d \times \{-1, 1\}$ and M_φ is the Lipschitz constant of the surrogate loss φ .

Note that the above bound involves only the VC dimension of the *base class*, which is typically small. On the other hand, the class \mathcal{F}_λ obtained by forming weighted combinations of classifiers from \mathcal{G} is extremely rich, and will generally have infinite VC dimension! But there is a price we pay: The first term is the empirical surrogate loss $A_{\varphi,n}(\hat{f}_n)$, rather than the empirical classification error $L_n(\hat{f}_n)$. However, it is possible to choose the surrogate φ in such a way that $A_{\varphi,n}(\cdot)$ can be bounded in terms of a quantity *related* to the number of misclassified training examples. Here is an example.

Fix a positive parameter $\gamma > 0$ and consider

$$\varphi(x) = \begin{cases} 0, & \text{if } x \leq -\gamma \\ 1, & \text{if } x \geq 0 \\ 1 + x/\gamma, & \text{otherwise} \end{cases}$$

This is a valid surrogate loss with $B = 1$ and $M_\varphi = 1/\gamma$, but in addition we have $\varphi(x) \leq \mathbf{1}_{\{x > -\gamma\}}$, which implies that $\varphi(-yf(x)) \leq \mathbf{1}_{\{yf(x) < \gamma\}}$. Therefore, for any f we have

$$A_{\varphi,n}(f) = \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i)) \leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i f(X_i) < \gamma\}}. \quad (18)$$

The quantity

$$L_n^\gamma(f) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i f(X_i) < \gamma\}} \quad (19)$$

is called the *margin error* of f . Notice that:

- For any $\gamma > 0$, $L_n^\gamma(f) \geq L_n(f)$
- The function $\gamma \mapsto L_n^\gamma(f)$ is increasing.

Notice also that we can write

$$L_n^\gamma(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i f(X_i) < 0\}} + \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{0 \leq Y_i f(X_i) < \gamma\}},$$

where the first term is just $L_n(f)$, while the second term is the number of training examples that were classified correctly, but only with small “margin” (the quantity $Yf(X)$ is often called the *margin* of the classifier f).

Theorem 4 (Margin-based risk bound for weighted linear combinations). *For any $\gamma > 0$, the bound*

$$L(\widehat{f}_n) \leq L_n^\gamma(\widehat{f}_n) + \frac{C\lambda}{\gamma} \sqrt{\frac{V(\mathcal{G})}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} \quad (20)$$

holds with probability at least $1 - \delta$.

Remark 1. Note that the first term on the right-hand side of (20) increases with γ , while the second term decreases with γ . Hence, if the learned classifier \widehat{f}_n has a small margin error for a large γ , i.e., it classifies the training samples well and with high “confidence,” then its generalization error will be small.

4 Kernel machines

Another powerful way of building complicated classifiers from simple functions is by means of *kernels*. Kernel methods are popular in machine learning for a variety of reasons, not the least of which is that any algorithm that operates in a Euclidean space and relies only on the computation of inner products between feature vectors can be modified to work with any suitably well-behaved kernel.

To start with, let us define what we mean by a kernel. We will stick to Euclidean feature spaces, although everything works out for arbitrary separable metric spaces.

Definition 2. *Let X be a closed subset of \mathbb{R}^d . A real-valued function $K : X \times X \rightarrow \mathbb{R}$ is called a Mercer kernel provided the following conditions are met:*

1. *It is symmetric, i.e., $K(x, x') = K(x', x)$ for any $x, x' \in X$.*
2. *It is continuous, i.e., if $\{x_n\}$ is a sequence of points in X converging to a point x , then*

$$\lim_{n \rightarrow \infty} K(x_n, x') = K(x, x'), \quad \forall x' \in X.$$

3. *It is positive semidefinite, i.e., for all $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ and all $x_1, \dots, x_n \in X$,*

$$\sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) \geq 0. \quad (21)$$

Remark 2. Another way to interpret the positive semidefiniteness condition is as follows. For any n -tuple $x^n = (x_1, \dots, x_n) \in X^n$, define the $n \times n$ kernel Gram matrix

$$G_K(x^n) \triangleq [K(x_i, x_j)]_{i,j=1}^n.$$

Then (21) is equivalent to saying that $G_K(x^n)$ is positive semidefinite in the usual sense, i.e., for any vector $v \in \mathbb{R}^n$ we have

$$\langle v, G_K(x^n) v \rangle \geq 0.$$

Remark 3. From now on, we will just say “kernel,” but always mean “Mercer kernel.”

Here are some examples of kernels:

1. With $X = \mathbb{R}^d$, $K(x, x') = \langle x, x' \rangle$, the usual Euclidean inner product.

2. A more general class of kernels based on the Euclidean inner product can be constructed as follows. Let $X = \{x \in \mathbb{R}^d : \|x\| \leq R\}$; choose any sequence $\{a_j\}_{j=0}^{\infty}$ of nonnegative reals such that

$$\sum_{j=0}^{\infty} a_j R^{2j} < \infty.$$

Then

$$K(x, x') = \sum_{j=0}^{\infty} a_j \langle x, x' \rangle^j$$

is a kernel.

3. Let $X = \mathbb{R}^d$, and let $k : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function, which is *reflection-symmetric*, i.e., $k(-x) = k(x)$ for all x . Then $K(x, x') \triangleq k(x - x')$ is a kernel provided the Fourier transform of k ,

$$\widehat{k}(\xi) \triangleq \int_{\mathbb{R}^d} e^{-i\langle \xi, x \rangle} k(x) dx,$$

is nonnegative. A prime example is the *Gaussian kernel*, induced by the function $k(x) = e^{-\gamma \|x\|^2}$.

In all of the above cases, the first two properties of a Mercer kernel are easy to check. The third, i.e., positive semidefiniteness, requires a bit more work. For details, consult Section 2.5 of the book by Cucker and Zhou [CZ07].

The importance of kernels in machine learning stems from the fact that we can use them to represent (or approximate) arbitrarily complicated continuous functions on the feature space X . In order to take full advantage of this representational power, we must take a detour into the theory of *Hilbert spaces*.

4.1 A crash course on Hilbert spaces

Hilbert spaces are a powerful generalization of the usual Euclidean space with an inner product; once we have an inner product, we can introduce the notion of an *angle* and, consequently, orthogonality. Moreover, a Hilbert space has certain favorable convergence properties, so we can speak about (unique) linear projections of their elements onto closed linear subspaces. Let us make these ideas precise.

Definition 3. A real vector space V is an inner product space if there exists a function $\langle \cdot, \cdot \rangle_V : V \times V \rightarrow \mathbb{R}$, which is:

1. *Symmetric:* $\langle v, v' \rangle_V = \langle v', v \rangle_V$ for all $v, v' \in V$
2. *Linear:* $\langle \alpha v_1 + \beta v_2, v' \rangle_V = \alpha \langle v_1, v' \rangle_V + \beta \langle v_2, v' \rangle_V$ for all $\alpha, \beta \in \mathbb{R}$ and all $v_1, v_2, v' \in V$
3. *Positive definite:* $\langle v, v \rangle_V \geq 0$ for all $v \in V$, and $\langle v, v \rangle_V = 0$ if and only if $v = 0$

Let $(V, \langle \cdot, \cdot \rangle_V)$ be an inner product space. Then we can define a *norm* on V via

$$\|v\|_V \triangleq \sqrt{\langle v, v \rangle_V}.$$

It is easy to check that this is, indeed, a norm —

1. It is homogeneous: for any $v \in V$ and any $\alpha \in \mathbb{R}$,

$$\|\alpha v\|_V = \sqrt{\langle \alpha v, \alpha v \rangle_V} = \sqrt{\alpha^2 \langle v, v \rangle_V} = |\alpha| \sqrt{\langle v, v \rangle_V} = |\alpha| \cdot \|v\|_V$$

2. It satisfies the triangle inequality: for any $v, v' \in V$,

$$\|v + v'\|_V \leq \|v\|_V + \|v'\|_V. \quad (22)$$

To prove this, we first need to establish another key property of $\|\cdot\|_V$: the *Cauchy–Schwarz inequality*, which generalizes its classical Euclidean counterpart and says that

$$|\langle v, v' \rangle_V| \leq \|v\|_V \|v'\|_V. \quad (23)$$

To prove (23), we start with the observation that $\|v - \lambda v'\|_V^2 = \langle v - \lambda v', v - \lambda v' \rangle_V \geq 0$ for any $\lambda \in \mathbb{R}$. Expanding this, we get

$$\langle v - \lambda v', v - \lambda v' \rangle_V = \lambda^2 \|v'\|_V^2 - 2\lambda \langle v, v' \rangle_V + \|v\|_V^2 \geq 0.$$

This is a quadratic function of λ , and from the above we see that its graph does not cross the horizontal axis. Therefore, we must have

$$4|\langle v, v' \rangle_V|^2 \leq 4\|v\|_V^2 \|v'\|_V^2 \quad \iff \quad |\langle v, v' \rangle_V| \leq \|v\|_V \|v'\|_V.$$

Now we can write

$$\begin{aligned} (\|v\|_V + \|v'\|_V)^2 &= \|v\|_V^2 + 2\|v\|_V \|v'\|_V + \|v'\|_V^2 \\ &\geq \|v\|_V^2 + 2\langle v, v' \rangle_V + \|v'\|_V^2 \\ &= \langle v, v \rangle_V + \langle v, v' \rangle_V + \langle v', v \rangle_V + \langle v', v' \rangle_V \\ &= \langle v + v', v + v' \rangle_V \\ &\equiv \|v + v'\|_V^2, \end{aligned}$$

where the first step uses the Cauchy–Schwarz inequality, the second step uses the definition of $\|\cdot\|_V$ and the symmetry of $\langle \cdot, \cdot \rangle_V$, the third step uses the linearity of $\langle \cdot, \cdot \rangle_V$, and the final step is, again, by definition. Since all norms are nonnegative, we can take square roots of both sides to get the triangle inequality.

3. Finally, $\|v\|_V \geq 0$, and $\|v\|_V = 0$ if and only if $v = 0$ – this is obvious from definitions.

Thus, an inner product space can be equipped with a norm that has certain special properties (mainly, the Cauchy–Schwarz inequality, since a lot of useful things follow from it alone). Now that we have a norm, we can talk about *convergence* of sequences in V :

Definition 4. Let $\{v_n\}_{n=1}^\infty$ be a sequence of elements of V . We say that it converges to $v \in V$ if

$$\lim_{n \rightarrow \infty} \|v_n - v\|_V = 0. \quad (24)$$

Remark 4. This definition is valid for any norm on V , not necessarily a norm induced by an inner product.

Any norm-convergent sequence has the property that, as n gets larger, its elements get closer and closer to one another. Specifically, suppose that $\{v_n\}$ converges to v . Then (24) implies that for any $\varepsilon > 0$ we can choose n large enough, so that $\|v_n - v\|_V < \varepsilon/2$ for all $m \geq n$. But the triangle inequality gives

$$\|v_n - v_m\|_V \leq \|v_n - v\|_V + \|v_m - v\|_V < \varepsilon, \quad \forall m \geq n.$$

In other words, we have

$$\lim_{m \rightarrow \infty} \|v_n - v_m\| = 0.$$

Since this holds for every n , we can write

$$\lim_{\min(m,n) \rightarrow \infty} \|v_n - v_m\| = 0. \quad (25)$$

Any sequence $\{v_n\}$ that has the property (25) is called a *Cauchy sequence*. We have just proved that any convergent sequence is Cauchy. However, the converse is not necessarily true: a Cauchy sequence does not have to be convergent. This motivates the following definition:

Definition 5. A normed space $(V, \|\cdot\|_V)$ is complete if any Cauchy sequence $\{v_n\}$ of its elements is convergent. If the norm $\|\cdot\|_V$ is induced by an inner product, then we say that V is a Hilbert space.

There is a standard procedure of starting with an inner product and the corresponding normed space and then *completing* it by adding the limits of all Cauchy sequences. We will not worry too much about this procedure. Here are a few standard examples of Hilbert spaces:

1. The Euclidean space $V = \mathbb{R}^d$ with the usual inner product

$$\langle v, v' \rangle = \sum_{j=1}^d v_j v'_j.$$

The corresponding norm is the familiar ℓ_2 norm, $\|v\| = \sqrt{\langle v, v \rangle}$.

2. More generally, if A is a positive definite $d \times d$ matrix, then the inner product

$$\langle v, v' \rangle_A \triangleq \langle v, Av' \rangle$$

induces the A -weighted norm $\|v\|_A \triangleq \sqrt{\langle v, v \rangle_A} = \sqrt{\langle v, Av \rangle}$, which makes \mathbb{R}^d into a Hilbert space. The preceding example is a special case with $A = I_d$, the $d \times d$ identity matrix.

3. The space $L^2(\mathbb{R}^d)$ of all *square-integrable* functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$, i.e.,

$$\int_{\mathbb{R}^d} f^2(x) dx < \infty,$$

is a Hilbert space with the inner product

$$\langle f, g \rangle_{L^2(\mathbb{R}^d)} \triangleq \int_{\mathbb{R}^d} f(x)g(x) dx$$

and the corresponding norm

$$\|f\|_{L^2(\mathbb{R}^d)} \triangleq \sqrt{\int_{\mathbb{R}^d} f^2(x) dx}.$$

4. Let (Ω, \mathcal{B}, P) be a probability space. Then the space $L^2(P)$ space of all real-valued random variables $X: \Omega \rightarrow \mathbb{R}$ with finite second moment, i.e.,

$$\mathbb{E}X^2 = \int_{\Omega} X^2(\omega)P(d\omega) < +\infty,$$

is a Hilbert space with the inner product

$$\langle X, X' \rangle_{L^2(P)} \triangleq \mathbb{E}[XX'] = \int_{\Omega} X(\omega)X'(\omega)P(d\omega)$$

and the corresponding norm

$$\|X\|_{L^2(P)} \triangleq \sqrt{\int_{\Omega} |X(\omega)|^2 P(d\omega)} \equiv \sqrt{\mathbb{E}X^2}.$$

From now on, we will denote a typical Hilbert space by $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$; the induced norm will be denoted by $\|\cdot\|_{\mathcal{H}}$.

An enormous advantage of working with Hilbert spaces is the availability of the notion of *orthogonality* and *orthogonal projection*. Two elements h, g of a Hilbert space \mathcal{H} are said to be *orthogonal* if $\langle h, g \rangle_{\mathcal{H}} = 0$.

Now consider a closed linear subspace \mathcal{H}_1 of \mathcal{H} , where “closed” means that the limit of any convergent sequence $\{h_n\}$ of elements of \mathcal{H}_1 is also contained in \mathcal{H}_1 . Then we have the following basic facts:

Theorem 5. Let \mathcal{H}_1^{\perp} be the set of all $h^{\perp} \in \mathcal{H}$, such that $\langle g, h^{\perp} \rangle_{\mathcal{H}} = 0$ for all $g \in \mathcal{H}_1$. Then:

1. \mathcal{H}_1^{\perp} is also a closed linear subspace of \mathcal{H} .
2. Any element g of \mathcal{H} can be uniquely decomposed as $g = h + h^{\perp}$, where $h \in \mathcal{H}_1$ and $h^{\perp} \in \mathcal{H}_1^{\perp}$.
3. Define the orthogonal projection $\Pi: \mathcal{H} \rightarrow \mathcal{H}_1$ onto \mathcal{H}_1 through

$$\Pi g \triangleq h \quad \text{if } g = h + h^{\perp} \text{ with } h \in \mathcal{H}_1, h^{\perp} \in \mathcal{H}_1^{\perp}.$$

Then Π has the following properties:

- (a) It is a linear operator.
- (b) $\Pi^2 = \Pi$, i.e., $\Pi(\Pi g) = \Pi g$ for any $g \in \mathcal{H}$.
- (c) If $g \in \mathcal{H}_1$, then $\Pi g = g$.
- (d) For any $g \in \mathcal{H}$ and any $h \in \mathcal{H}_1$,

$$\langle \Pi g, h \rangle_{\mathcal{H}} = \langle g, h \rangle_{\mathcal{H}}.$$

- (e) For any $g \in \mathcal{H}$, $h = \Pi g \in \mathcal{H}_1$ is the unique solution of the optimization problem

$$\text{minimize } \|h - g\| \text{ subject to } h \in \mathcal{H}_1.$$

Remark 5. It is important for \mathcal{H}_1 to be a *closed* linear subspace of \mathcal{H} for the above results to hold.

4.2 Reproducing kernel Hilbert spaces

Now let us return to our original goal. Suppose we have a fixed kernel K on our feature space X (which we assume to be a closed subset of \mathbb{R}^d). Let $\mathcal{L}_K(X)$ be the *linear span* of the set $\{K(x', \cdot) : x' \in X\}$, i.e., the set of all functions $f : X \rightarrow \mathbb{R}$ of the form

$$f(x) = \sum_{j=1}^N c_j K(x_j, x) \quad (26)$$

for all possible choices of $N \in \mathbb{N}$, $c_1, \dots, c_N \in \mathbb{R}$, and $x_1, \dots, x_N \in X$. It is easy to see that $\mathcal{L}_K(X)$ is a *vector space*: for any two functions f, f' of the form (26), their sum is also of that form; if we multiply any $f \in \mathcal{L}_K(X)$ by a scalar $c \in \mathbb{R}$, we will get another element of $\mathcal{L}_K(X)$; and the zero function is clearly in $\mathcal{L}_K(X)$. It turns out that, for any (Mercer) kernel K , we can *complete* $\mathcal{L}_K(X)$ into a *Hilbert space* of functions that can potentially represent *any* continuous function from X into \mathbb{R} , provided K is chosen appropriately.

The following result is essential (for the proof, see Section 2.4 of Cucker and Zhou [CZ07]):

Theorem 6. *Let X be a closed subset of \mathbb{R}^d , and let $K : X \times X \rightarrow \mathbb{R}$ be a Mercer kernel. Then there exists a unique Hilbert space $(\mathcal{H}_K, \langle \cdot, \cdot \rangle_K)$ of real-valued functions on X with the following properties:*

1. *For all $x \in X$, the function $K_x(\cdot) \triangleq K(x, \cdot)$ is an element of \mathcal{H}_K , and $\langle K_x, K_{x'} \rangle_K = K(x, x')$ for all $x, x' \in X$.*
2. *The linear space $\mathcal{L}_K(X)$ is dense in \mathcal{H}_K , i.e., for any $f \in \mathcal{H}_K$ and any $\varepsilon > 0$ there exist some $N \in \mathbb{N}$, $c_1, \dots, c_N \in \mathbb{R}$, and $x_1, \dots, x_N \in X$, such that*

$$\left\| f - \sum_{j=1}^N c_j K_{x_j} \right\|_K < \varepsilon.$$

3. *For all $f \in \mathcal{H}_K$ and all $x \in X$,*

$$f(x) = \langle K_x, f \rangle_K. \quad (27)$$

Moreover, the functions in \mathcal{H}_K are continuous. The Hilbert space \mathcal{H}_K is called the *Reproducing Kernel Hilbert Space (RKHS) associated with K* ; the property (27) is referred to as the *reproducing kernel property*.

Remark 6. The reproducing kernel property essentially states that the value of any function $f \in \mathcal{H}_K$ at any point $x \in X$ can be “extracted” by projecting f onto the function $K_x(\cdot) = K(x, \cdot)$, i.e., a copy of the kernel K “centered” at the point x . It is easy to prove when $f \in \mathcal{L}_K(X)$. Indeed, if f has the form (26), then

$$\begin{aligned} \langle f, K_x \rangle_K &= \left\langle \sum_{j=1}^N c_j K_{x_j}, K_x \right\rangle_K \\ &= \sum_{j=1}^N c_j \langle K_{x_j}, K_x \rangle_K \\ &= \sum_{j=1}^N c_j K(x_j, x) \\ &= f(x). \end{aligned}$$

Since any $f \in \mathcal{H}_K$ can be expressed as a limit of functions from $\mathcal{L}_K(X)$, the proof of (27) for a general f follows by continuity.

Now we pick a kernel K on our feature space and consider classifiers of the form (9) with the underlying f taken from a suitable subset of the RKHS \mathcal{H}_K . One choice, which underlies such things as the Support Vector Machine, is to take a ball in \mathcal{H}_K : given some $\lambda > 0$, let

$$\mathcal{F}_\lambda \triangleq \{f \in \mathcal{H}_K : \|f\|_K \leq \lambda\}.$$

This set is the closure (in the $\|\cdot\|_K$ norm) of the convex set

$$\left\{ \sum_{j=1}^N c_j K_{x_j} : N \in \mathbb{N}; c_1, \dots, c_N \in \mathbb{R}; x_1, \dots, x_N \in X; \sum_{i,j=1}^N c_i c_j K(x_i, x_j) \leq \lambda^2 \right\} \subset \mathcal{L}_K(X),$$

and is itself convex. Now, as we already know, the performance of any learning algorithm that chooses an element $\hat{f}_n \in \mathcal{F}_\lambda$ in a data-dependent way is controlled by the Rademacher average $R_n(\mathcal{F}_\lambda(X^n))$. It turns out that this Rademacher average is fairly easy to estimate. Indeed, using the reproducing kernel property (27) and then the linearity of the inner product $\langle \cdot, \cdot \rangle_K$, we can write

$$\begin{aligned} R_n(\mathcal{F}_\lambda(X^n)) &= \sup_{f: \|f\|_K \leq \lambda} \frac{1}{n} \mathbb{E}_{\sigma^n} \left| \sum_{i=1}^n \sigma_i f(X_i) \right| \\ &= \sup_{f: \|f\|_K \leq \lambda} \frac{1}{n} \mathbb{E}_{\sigma^n} \left| \sum_{i=1}^n \sigma_i \langle f, K_{X_i} \rangle_K \right| \\ &= \sup_{f: \|f\|_K \leq \lambda} \frac{1}{n} \mathbb{E}_{\sigma^n} \left| \left\langle f, \sum_{i=1}^n \sigma_i K_{X_i} \right\rangle_K \right| \end{aligned}$$

Now, using the Cauchy–Schwarz inequality (23), it is not hard to show that

$$\sup_{f: \|f\|_K \leq \lambda} |\langle f, g \rangle_K| = \lambda \|g\|_K$$

for any $g \in \mathcal{H}_K$. Therefore,

$$R_n(\mathcal{F}_\lambda(X^n)) = \frac{\lambda}{n} \mathbb{E}_{\sigma^n} \left\| \sum_{i=1}^n \sigma_i K_{X_i} \right\|_K.$$

Now we exploit the following easily proved fact: for any n functions $g_1, \dots, g_n \in \mathcal{H}_K$,

$$\mathbb{E}_{\sigma^n} \left\| \sum_{i=1}^n \sigma_i g_i \right\|_K \leq \sqrt{\sum_{i=1}^n \|g_i\|_K^2}. \quad (28)$$

The proof of this is in two steps: First, we use the concavity of the square root to write

$$\mathbb{E}_{\sigma^n} \sqrt{\left\| \sum_{i=1}^n \sigma_i g_i \right\|_K^2} \leq \sqrt{\mathbb{E} \left\| \sum_{i=1}^n \sigma_i g_i \right\|_K^2}.$$

Then we expand the squared norm:

$$\left\| \sum_{i=1}^n \sigma_i g_i \right\|_K^2 = \left\langle \sum_{i=1}^n \sigma_i g_i, \sum_{i=1}^n \sigma_i g_i \right\rangle_K = \sum_{i,j=1}^n \sigma_i \sigma_j \langle g_i, g_j \rangle_K.$$

And finally we take the expectation over σ^n and use the fact that $\mathbb{E}[\sigma_i \sigma_j] = 1$ if $i = j$ and 0 otherwise to get

$$\mathbb{E} \left\| \sum_{i=1}^n \sigma_i g_i \right\|_K^2 = \sum_{i=1}^n \langle g_i, g_i \rangle_K = \sum_{i=1}^n \|g_i\|_K^2.$$

Hence, we obtain

$$R_n(\mathcal{F}(X^n)) \leq \frac{1}{n} \sqrt{\sum_{i=1}^n \langle K_{X_i}, K_{X_i} \rangle_K} = \frac{1}{n} \sqrt{\sum_{i=1}^n K(X_i, X_i)}.$$

Finally, taking the expectation w.r.t. X^n and once more using concavity of the square root, we have

$$\mathbb{E} R_n(\mathcal{F}(X^n)) \leq \frac{\lambda \sqrt{\mathbb{E} K(X, X)}}{\sqrt{n}}.$$

4.3 Empirical risk minimization in an RKHS

Another advantage of working with kernels is that, in many cases, a minimizer of empirical risk over a sufficiently regular subset of an RKHS will have the form of a linear combination of kernels centered at the training feature points. The results ensuring this are often referred to in the literature as *representer theorems*. Here is one such result (due, in a slightly different form, to Schölkopf, Herbrich, and Smola [SHS01]), sufficiently general for our purposes:

Theorem 7 (The generalized representer theorem). *Let X be a closed subset of \mathbb{R}^d and let Y be a subset of the reals. Consider a nonnegative loss function $\ell : Y \times Y \rightarrow \mathbb{R}^+$. Let K be a Mercer kernel on X , and let \mathcal{H}_K be the corresponding RKHS.*

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be an i.i.d. sample from some distribution $P = P_{XY}$ on $X \times Y$, let $\mathcal{H}_K^{X^n}$ be the closed linear subspace of \mathcal{H}_K spanned by $\{K_{X_i} : 1 \leq i \leq n\}$, and let Π_n denote the orthogonal projection onto $\mathcal{H}_K^{X^n}$. Let \mathcal{F} be a subset of \mathcal{H}_K , such that $\Pi_n(\mathcal{F}) \subseteq \mathcal{F}$. Then

$$\inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)) = \inf_{f \in \Pi_n(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)), \quad (29)$$

and if a minimizer of the left-hand side of (29) exists, then it can be taken to have the form

$$\hat{f}_n = \sum_{i=1}^n c_i K_{X_i} \quad (30)$$

for some $c_1, \dots, c_n \in \mathbb{R}$.

Remark 7. Note that both the subspace $\mathcal{H}_K^{X^n}$ and the corresponding orthogonal projection Π_n are *random objects*, since they depend on the random features X^n .

Proof. Since $K_{X_i} \in \mathcal{H}_K^{X^n}$ for every i , by Theorem 5 we have

$$\langle f, K_{X_i} \rangle_K = \langle \Pi_n f, K_{X_i} \rangle_K, \quad \forall f \in \mathcal{H}_K.$$

Moreover, from the reproducing kernel property (27) we deduce that

$$f(X_i) = \langle f, K_{X_i} \rangle_K = \langle \Pi_n f, K_{X_i} \rangle_K = \Pi_n f(X_i).$$

Therefore, for every $f \in \mathcal{F}$ we can write

$$\frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \Pi_n f(X_i)).$$

This implies that

$$\inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)) = \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \Pi_n f(X_i)) = \inf_{g \in \Pi_n(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, g(X_i)). \quad (31)$$

Now suppose that $\hat{f}_n \in \mathcal{F}$ achieves the infimum on the left-hand side of (31). Then its projection $\Pi_n \hat{f}_n$ onto $\mathcal{H}_K^{X^n}$ achieves the infimum on the right-hand side. Moreover, since $\Pi_n(\mathcal{F}) \subseteq \mathcal{F}$ by hypothesis, we may conclude that $f = \Pi_n(\hat{f}_n)$, i.e., $f \in \mathcal{H}_K^{X^n}$. Since every element of $\mathcal{H}_K^{X^n}$ has the form (30), the theorem is proved. \square

In the classification setting, we may take $Y = \{-1, +1\}$ and consider the problem of minimizing the empirical surrogate loss

$$A_{\varphi, n}(f) = \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i))$$

over the ball \mathcal{F}_λ in a suitable RKHS \mathcal{H}_K . By the above theorem, we may write this problem in the following form:

$$\min_{c_1, \dots, c_n \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \varphi \left(-Y_i \sum_{j=1}^n c_j K(X_i, X_j) \right) \quad (32a)$$

$$\text{subject to } \sum_{i,j=1}^n c_i c_j K(X_i, X_j) \leq \lambda^2 \quad (32b)$$

Suppose the surrogate loss function φ is convex. Then the objective function in (32) is convex as well, and the decision variables $c_1, \dots, c_n \in \mathbb{R}$ are subject to a quadratic constraint. Thus, (32) is an instance of a *quadratically constrained convex program* (QCCP). Moreover, when φ is such that the objective is *quadratic* in c_1, \dots, c_n , then we have a *quadratically constrained quadratic problem* (QCQP), which can be solved very efficiently using interior point methods. For detailed background see the text of Boyd and Vandenberghe [BV04]. Many popular machine learning algorithms can be cast in the form (32). For instance, if we let φ be the hinge loss $\varphi(u) = (u+1)_+$, then (32) corresponds to the *Support Vector Machine* (SVM) algorithm — more precisely, the SVM is the *scalarized* version of (32), i.e., it has the form

$$\min_{c_1, \dots, c_n \in \mathbb{R}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(1 - Y_i \sum_{j=1}^n c_j K(X_i, X_j) \right)_+ + \tau \sum_{i,j=1}^n c_i c_j K(X_i, X_j) \right\}$$

for some regularization parameter $\tau > 0$.

5 Convex risk minimization

Choosing a convex surrogate loss function φ has many advantages in general. First of all, we may arrange things in such a way that minimizing the surrogate loss $A_\varphi(f)$ over all measurable $f : X \rightarrow \mathbb{R}$ is equivalent to determining the Bayes classifier (1):

Theorem 8. Let $P = P_{XY}$ be the joint distribution of the feature $X \in \mathbb{R}^d$ and the binary label $Y \in \{-1, +1\}$, and let $\eta(x) = \mathbb{P}[Y = 1|X = x]$ be the corresponding regression function. Consider a surrogate loss function φ , which is strictly convex and differentiable. Then the unique minimizer of the surrogate loss $A_\varphi(f) = \mathbb{E}[\varphi(-Yf(X))]$ over all (measurable) functions $f : X \rightarrow \mathbb{R}$ has the form

$$f^*(x) = \operatorname{argmin}_{u \in \mathbb{R}} h_{\eta(x)}(u),$$

where for each $\eta \in [0, 1]$ we have $h_\eta(u) \triangleq \eta\varphi(-u) + (1 - \eta)\varphi(u)$. Moreover, $f^*(x)$ is positive if and only if $\eta(x) > 1/2$, i.e., the induced sign classifier $g_{f^*}(x) = \operatorname{sgn}(f^*(x))$ is the Bayes classifier (1).

Proof. By the law of iterated expectation,

$$A_\varphi(f) = \mathbb{E}[\varphi(-Yf(X))] = \mathbb{E}[\mathbb{E}[\varphi(-Yf(X))|X]].$$

Hence,

$$\begin{aligned} \inf_f A_\varphi(f) &= \inf_f \mathbb{E}[\mathbb{E}[\varphi(-Yf(X))|X]] \\ &= \mathbb{E}\left[\inf_{u \in \mathbb{R}} \mathbb{E}[\varphi(-Yu)|X = x]\right]. \end{aligned}$$

For every $x \in X$, we have

$$\begin{aligned} \mathbb{E}[\varphi(-Yu)|X = x] &= \mathbb{P}[Y = 1|X = x]\varphi(-u) + \mathbb{P}[Y = -1|X = x]\varphi(u) \\ &= \eta(x)\varphi(-u) + (1 - \eta)\varphi(u) \\ &\equiv h_{\eta(x)}(u). \end{aligned}$$

Since φ is strictly convex and differentiable, so is h_η for every $\eta \in [0, 1]$. Therefore, $\inf_{u \in \mathbb{R}} h_\eta(u)$ exists, and is achieved by a unique u^* ; in particular,

$$f^*(x) = \operatorname{argmin}_{u \in \mathbb{R}} h_{\eta(x)}(u).$$

To find the u^* minimizing h_η , we differentiate h_η w.r.t. u and set the derivative to zero. Since

$$h'_\eta(u) = -(1 - \eta)\varphi'(-u) + \eta\varphi'(u),$$

the point of minimum u^* is the solution to the equation

$$\frac{\varphi'(u)}{\varphi'(-u)} = \frac{\eta}{1 - \eta}.$$

Suppose $\eta > 1/2$; then

$$\frac{\varphi'(u)}{\varphi'(-u)} > 1.$$

Since φ is strictly convex, its derivative φ' is strictly increasing. Hence, $u^* > -u^*$ which implies that $u^* > 0$. Conversely, if $u^* \leq 0$, then $u^* \leq -u^*$, so $\varphi'(u^*) \leq \varphi'(-u^*)$, which means that $\eta/(1 - \eta) \leq 1$, i.e., $\eta \leq 1/2$. Thus, we conclude that $f^*(x)$, which is the minimizer of $h_{\eta(x)}$, is positive if and only if $\eta(x) > 1/2$, i.e., $\operatorname{sgn}(f^*(x))$ is the Bayes classifier. \square

Secondly, under some additional regularity conditions it is possible to relate the minimum surrogate loss

$$A_\varphi^* \triangleq \inf_f A_\varphi(f)$$

to the Bayes rate

$$L^* = \inf_f \mathbb{P}(Y \neq f(X)) :$$

Theorem 9. Assume that the surrogate loss function φ satisfies the conditions of Theorem 3, and that there exist positive constants $s \geq 1$ and c , such that the inequality

$$L(f) - L^* \leq c \left(A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad (33)$$

holds for any measurable function $f : X \rightarrow \mathbb{R}$. Consider the learning algorithm that minimizes empirical surrogate loss over some class \mathcal{F} :

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} A_{\varphi,n}(f) = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i)). \quad (34)$$

Then

$$L(\hat{f}_n) - L^* \leq 2^{1/s} c \left(4M_\varphi \mathbb{E} R_n(\mathcal{F}(X^n)) + B \sqrt{\frac{\log(1/\delta)}{2n}} \right)^{1/s} + c \left(\inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad (35)$$

with probability at least $1 - \delta$.

Proof. We have the following:

$$L(\hat{f}_n) - L^* \leq c \left(A_\varphi(\hat{f}_n) - A_\varphi^* \right)^{1/s} \quad (36)$$

$$= c \left(A_\varphi(\hat{f}_n) - \inf_{f \in \mathcal{F}} A_\varphi(f) + \inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad (37)$$

$$\leq c \left(A_\varphi(\hat{f}_n) - \inf_{f \in \mathcal{F}} A_\varphi(f) \right)^{1/s} + c \left(\inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad (38)$$

$$\leq 2^{1/s} c \left(\sup_{f \in \mathcal{F}} |A_{\varphi,n}(f) - A_\varphi(f)| \right)^{1/s} + c \left(\inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad (39)$$

$$\leq 2^{1/s} c \left(4M_\varphi \mathbb{E} R_n(\mathcal{F}(X^n)) + B \sqrt{\frac{\log(1/\delta)}{2n}} \right)^{1/s} + c \left(\inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^* \right)^{1/s} \quad \text{w.p. } \geq 1 - \delta, \quad (40)$$

where:

- (36) follows from (33);
- (38) follows from the inequality $(a + b)^{1/s} \leq a^{1/s} + b^{1/s}$ that holds for all $a, b \geq 0$ and all $s \geq 1$

- (39) and (40) follow from the same argument as the one used in the proof of Theorem 3.

This completes the proof. □

Remark 8. Condition (33) is often easy to check. For instance, Zhang [Zha04] proved that it is satisfied, provided the inequality

$$\left| \frac{1}{2} - \eta \right|^s \leq (2c)^s \left(1 - \inf_u h_\eta(u) \right) \quad (41)$$

holds for all $\eta \in [0, 1]$. For instance, (41) holds for the exponential loss $\varphi(u) = e^u$ and the logit loss $\varphi(u) = \log_2(1 + e^u)$ with $s = 2$ and $c = 2\sqrt{2}$; for the hinge loss $\varphi(u) = (u + 1)_+$, (41) holds with $s = 1$ and $c = 4$.

What Theorem 9 says is that, assuming the expected Rademacher average $\mathbb{E}R_n(\mathcal{F}(X^n)) = O(1/\sqrt{n})$, the difference between the generalization error of the Convex Risk Minimization algorithm (34) and the Bayes rate L^* is, with high probability, bounded by the combination of two terms: the $O(n^{-1/2s})$ “estimation error” term and the $(\inf_{f \in \mathcal{F}} A_\varphi(f) - A_\varphi^*)^{1/s}$ “approximation error” term. If the hypothesis space \mathcal{F} is rich enough, so that $\inf_{f \in \mathcal{F}} A_\varphi(f) = A_\varphi^*$, then the difference between $L(\hat{f}_n)$ and L^* is, with high probability, bounded as $O(1/n^{-2s})$, *independently* of the dimension d of the feature space.

References

- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [CZ07] F. Cucker and D. X. Zhou. *Learning Theory: An Approximation Theory Viewpoint*. Cambridge University Press, 2007.
- [SHS01] B. Schölkopf, R. Herbrich, and A. Smola. A generalized representer theorem. In D. Helmbold and B. Williamson, editors, *Computational Learning Theory*, volume 2111 of *Lecture Notes in Computer Science*, pages 416–426. Springer, 2001.
- [Zha04] T. Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32(1):56–134, 2004.