

Empirical Risk Minimization: Abstract risk bounds and Rademacher averages

Maxim Raginsky

September 24, 2013

In the last lecture, we have left off with a theorem that gave a sufficient condition for the *Empirical Risk Minimization* (ERM) algorithm

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} \hat{P}_{Z^n}(\ell_f) \quad (1)$$

$$= \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)) \quad (2)$$

to be PAC for a given learning problem with hypothesis space \mathcal{F} and loss function ℓ . This condition pertained to the behavior of the uniform deviation of empirical means from true means over the *induced class* $\mathcal{L}_{\mathcal{F}} = \{\ell_f : f \in \mathcal{F}\}$. Specifically, we proved that ERM is a PAC algorithm if

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P^n \left(\sup_{f \in \mathcal{F}} |\hat{P}_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon \right) = 0, \quad \forall \varepsilon > 0, \quad (3)$$

where \mathcal{P} is the class of probability distributions generating the training data.

1 An abstract framework for ERM

To study ERM in a general framework, we will adopt a simplified notation often used in the literature. We have a space Z and a class \mathcal{F} of functions $f : Z \rightarrow [0, 1]$. Let $\mathcal{P}(Z)$ denote the space of all probability distributions on Z . For each sample size n , the training data are in the form of an n -tuple $Z^n = (Z_1, \dots, Z_n)$ of Z -valued random variables drawn according to some unknown $P \in \mathcal{P}$. For each P , we can compute the *expected risk* of any $f \in \mathcal{F}$ by

$$P(f) = \mathbb{E}_P f(Z) = \int_Z f(z) P(dz). \quad (4)$$

The *minimum risk* over \mathcal{F} is

$$L_P^*(\mathcal{F}) \triangleq \inf_{f \in \mathcal{F}} P(f). \quad (5)$$

A learning algorithm is a sequence $\mathcal{A} = \{A_n\}_{n \geq 1}$ of mappings $A_n : Z^n \rightarrow \mathcal{F}$, and the objective is to ensure that

$$P(\hat{f}_n) \approx L_P^*(\mathcal{F}) \quad \text{eventually with high probability.} \quad (6)$$

The ERM algorithm works by taking

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} \hat{P}_{Z^n}(f) \quad (7)$$

$$= \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(Z_i). \quad (8)$$

This way of writing down our problem hides most of the ingredients that were specified in Haussler's framework of model-free learning, so it is important to keep in mind that Z is an input/output pair (X, Y) and the functions $f \in \mathcal{F}$ are really the induced losses for some loss function ℓ and hypothesis class \mathcal{G} .

We have already seen that the consistency of ERM hinges on the uniform deviation behavior of empirical means in \mathcal{F} . In order to have a clean way of keeping track of all the relevant quantities, let us introduce some additional notation. First of all, we need a way of comparing the behavior of any two probability distributions P and P' on the class \mathcal{F} . A convenient way of doing this is through the \mathcal{F} -seminorm

$$\|P - P'\|_{\mathcal{F}} \triangleq \sup_{f \in \mathcal{F}} |P(f) - P'(f)| \quad (9)$$

$$= \sup_{f \in \mathcal{F}} |\mathbb{E}_P f - \mathbb{E}_{P'} f| \quad (10)$$

$$= \sup_{f \in \mathcal{F}} \left| \int_Z f(z) P(dz) - \int_Z f(z) P'(dz) \right|. \quad (11)$$

We say that $\|\cdot\|_{\mathcal{F}}$ is a *seminorm* because it has all the properties of a norm (in particular, it satisfies the triangle inequality), but it may happen that $\|P - P'\|_{\mathcal{F}} = 0$ for $P \neq P'$. Next, given a random sample Z^n we define the *uniform deviation*

$$\Delta_n(Z^n) \triangleq \|\hat{P}_{Z^n} - P\|_{\mathcal{F}} \equiv \sup_{f \in \mathcal{F}} |\hat{P}_{Z^n}(f) - P(f)|. \quad (12)$$

To keep things simple, we do not indicate the underlying distribution P or the function class \mathcal{F} explicitly. We will do this from now on, unless some confusion is possible, in which case we will use appropriate indices. Thus, we will write $L(f)$, $L^*(\mathcal{F})$, etc., and you should always keep in mind that all expectations are computed w.r.t. the (unknown) data-generating distribution $P \in \mathcal{P}(Z)$. In the same spirit, we will denote by $P_n(f)$ the *empirical risk* of f on the sample Z^n :

$$P_n(f) = \hat{P}_{Z^n}(f) = \frac{1}{n} \sum_{i=1}^n f(Z_i). \quad (13)$$

The following result is key to understanding the role of the uniform deviations $\Delta_n(Z^n)$ in controlling the performance of the ERM algorithm.

Proposition 1. *The ERM algorithm satisfies the following inequalities:*

$$P(\hat{f}_n) \leq L^*(\mathcal{F}) + 2\Delta_n(Z^n) \quad (14)$$

$$P(\hat{f}_n) \leq P_n(\hat{f}_n) + \Delta_n(Z^n). \quad (15)$$

Proof. We have already proved the two inequalities of the proposition in the last lecture, except now they are written in our new abstract notation. Let us give the proof again in order to get comfortable with the notation. Let f^* be any minimizer of $P(f)$ over \mathcal{F} . Then

$$P(\hat{f}_n) - L^*(\mathcal{F}) = P(\hat{f}_n) - P(f^*) \quad (16)$$

$$= P(\hat{f}_n) - P_n(\hat{f}_n) + P_n(\hat{f}_n) - P_n(f^*) + P_n(f^*) - P(f^*), \quad (17)$$

where $P_n(\hat{f}_n) - P_n(f^*) \leq 0$ by definition of ERM,

$$P(\hat{f}_n) - P_n(\hat{f}_n) \leq \sup_{f \in \mathcal{F}} [P_n(f) - P(f)] \leq \|P_n - P\|_{\mathcal{F}} = \Delta_n(Z^n), \quad (18)$$

and the same holds for $P_n(f^*) - P(f^*)$. This proves both (14) and (15). \square

The bound (14) says that, if the uniform deviation $\Delta_n(Z^n)$ is small, then the expected risk of the ERM hypothesis will be close to the minimum risk $L^*(\mathcal{F})$; in addition, the bound (15) says that the empirical estimate $P_n(\hat{f}_n)$ is an accurate estimate of the generalization performance of \hat{f}_n . Both bounds suggest that the success of ERM depends on how small the uniform deviation $\Delta_n(Z^n)$ can be. Thus, we need to develop tools for analyzing the behavior of $\Delta_n(Z^n)$.

2 Bounding the uniform deviation: Rademacher averages

It turns out that the behavior of the uniform deviation $\Delta_n(Z^n)$ is closely connected to how the values of the functions $f \in \mathcal{F}$ on randomly selected n -tuples Z^n correlate with random signs. Intuitively, this can be motivated as follows. In order for ERM to succeed, the function class \mathcal{F} has to be “discriminating:” we should be able to clearly separate all near-minimizers of the empirical risk from functions whose empirical risks (and hence expected risks) are high, but only if the sample is representative of the true data-generating distribution. If the class \mathcal{F} is discriminating not only on the actual sample, but also on its random perturbations, then we cannot expect the empirical risks to truly reflect the generalization ability of the functions in \mathcal{F} . As we will soon see, the degree of correlation of the “projections” of \mathcal{F} onto random samples with random signs is captured by the quantities known as the *Rademacher averages*.

First, we need some preparatory results. Let Y be a real-valued random variable. We say that it is *symmetric* if $-Y$ has the same distribution as Y . This is equivalent to saying that

$$\mathbb{P}(Y \geq a) = \mathbb{P}(Y \leq -a), \quad \forall a \in \mathbb{R}. \quad (19)$$

A random variable σ taking values -1 or $+1$ with probability $1/2$ is called a *Rademacher random variable*.

Lemma 1. *Let U and U' be two i.i.d. real-valued random variables. Then $Y = U - U'$ is symmetric.*

Proof. Let $F(u) = \mathbb{P}(U \geq u)$. Then

$$\mathbb{P}(Y \geq a) = \mathbb{P}(U - U' \geq a) \quad (20)$$

$$= \mathbb{E} \left[\mathbb{E} \left[\mathbf{1}_{\{U - U' \geq a\}} \middle| U' \right] \right] \quad (21)$$

$$= \mathbb{E} F(a + U') \quad (22)$$

$$= \mathbb{E} F(a + U), \quad (23)$$

where the last line is because U and U' are i.i.d. An analogous calculation for $\mathbb{P}(Y \leq -a)$ gives the same result. \square

Lemma 2. *Let Y be a symmetric random variable, and let σ be a Rademacher random variable independent of Y . Then $W = \sigma Y$ has the same distribution as Y .*

Proof. Direct calculation:

$$\mathbb{P}(W \geq a) = \frac{1}{2} \mathbb{P}(Y \geq a) + \frac{1}{2} \mathbb{P}(Y \leq -a) = \mathbb{P}(Y \geq a), \quad (24)$$

where the second step is due to the symmetry of Y . Since this holds for an arbitrary $a \in \mathbb{R}$, we conclude that W has the same distribution as Y . \square

Corollary 1. *Let Y_1, \dots, Y_n be n independent symmetric random variables, and let $\sigma_1, \dots, \sigma_n$ be independent Rademacher random variables that are also independent of the Y_i 's. Then the sum $Y_1 + \dots + Y_n$ has the same distribution as $\sigma_1 Y_1 + \dots + \sigma_n Y_n$.*

Now we are ready to define Rademacher averages.

Definition 1. *The Rademacher average of a bounded set $\mathcal{A} \subset \mathbb{R}^n$ is*

$$R_n(\mathcal{A}) \triangleq \mathbb{E} \sup_{a \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right|, \quad (25)$$

where the expectation is over n i.i.d. Rademacher random variables $\sigma_1, \dots, \sigma_n$.

Now consider a class \mathcal{F} of functions $f : Z \rightarrow [0, 1]$ from our formulation of the ERM problem. The key result, which we will now prove, is that the uniform deviations $\Delta_n(Z^n)$ are controlled by the Rademacher averages of the *random sets*

$$\mathcal{F}(Z^n) \triangleq \{(f(Z_1), \dots, f(Z_n)) : f \in \mathcal{F}\}. \quad (26)$$

A useful way to think about $\mathcal{F}(Z^n)$ is as a “projection” of \mathcal{F} onto the random sample Z^n .

Theorem 1. *Fix a space Z and let \mathcal{F} be a class of functions $f : Z \rightarrow [0, 1]$. Then for any $P \in \mathcal{P}(Z)$*

$$\mathbb{E} \Delta_n(Z^n) \leq 2 \mathbb{E} R_n(\mathcal{F}(Z^n)). \quad (27)$$

Proof. The proof uses a clever technique known as “symmetrization,” which goes back to the seminal work of Vapnik and Chervonenkis [VC71], but in its modern form is due to Giné and Zinn [GZ84]. The main idea is as follows. Consider a random i.i.d. sample Z^n from P and introduce an independent “ghost” sample $\bar{Z}^n = (\bar{Z}_1, \dots, \bar{Z}_n)$ from the same P . We will denote expectations w.r.t. \bar{Z}^n by $\bar{\mathbb{E}}$. Let \bar{P}_n denote the empirical distribution of \bar{Z}^n . Then for any bounded function $g : Z \rightarrow \mathbb{R}$ we can write

$$P(g) = \bar{\mathbb{E}} \left[\frac{1}{n} \sum_{i=1}^n g(\bar{Z}_i) \right] = \bar{\mathbb{E}} \bar{P}_n(g). \quad (28)$$

With this, we have

$$\Delta_n(Z^n) = \|P_n - P\|_{\mathcal{F}} \quad (29)$$

$$= \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \quad (30)$$

$$= \sup_{f \in \mathcal{F}} |P_n(f) - \bar{\mathbb{E}} \bar{P}_n(f)| \quad (31)$$

$$\leq \sup_{f \in \mathcal{F}} \bar{\mathbb{E}} |P_n(f) - \bar{P}_n(f)| \quad (32)$$

$$\leq \bar{\mathbb{E}} \sup_{f \in \mathcal{F}} |P_n(f) - \bar{P}_n(f)|, \quad (33)$$

where the first inequality uses convexity of the absolute value function, while the second is because $\sup \bar{\mathbb{E}}[\cdot] \leq \bar{\mathbb{E}} \sup[\cdot]$. Now let us take expectations of both sides w.r.t. Z^n to get

$$\mathbb{E} \Delta_n(Z^n) \leq \mathbb{E} \sup_{f \in \mathcal{F}} |P_n(f) - \bar{P}_n(f)|, \quad (34)$$

where now the expectation on the right is w.r.t. both Z^n and \bar{Z}^n , which are independent of each other. Let us inspect the difference $P_n(f) - \bar{P}_n(f)$:

$$P_n(f) - \bar{P}_n(f) = \frac{1}{n} \sum_{i=1}^n [f(Z_i) - f(\bar{Z}_i)]. \quad (35)$$

For each i Z_i and \bar{Z}_i are i.i.d., so the differences $f(Z_i) - f(\bar{Z}_i)$ are symmetric by Lemma 1. Introducing n i.i.d. Rademacher random variables $\sigma_1, \dots, \sigma_n$ independent of Z^n and \bar{Z}^n , from Corollary 1 we know that

$$\frac{1}{n} \sum_{i=1}^n [f(Z_i) - f(\bar{Z}_i)] \stackrel{d}{=} \frac{1}{n} \sum_{i=1}^n \sigma_i [f(Z_i) - f(\bar{Z}_i)], \quad (36)$$

where $\stackrel{d}{=}$ means “equality in distribution.” The same holds if we take the supremum of both sides over $f \in \mathcal{F}$. Hence,

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_n(f) - \bar{P}_n(f)| = \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n [f(Z_i) - f(\bar{Z}_i)] \right| \quad (37)$$

$$= \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i [f(Z_i) - f(\bar{Z}_i)] \right|, \quad (38)$$

where in the last line the expectation is over Z^n, \bar{Z}^n , and $\sigma^n = (\sigma_1, \dots, \sigma_n)$. Now note that

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left[f(Z_i) - f(\bar{Z}_i) \right] \right| \leq \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right| + \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(\bar{Z}_i) \right| \quad (39)$$

$$= 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right|, \quad (40)$$

where the first line is by the triangle inequality and the second line uses the fact that Z^n has the same distribution as \bar{Z}^n . Now, since Z^n and σ^n are independent,

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right| = \mathbb{E}_{Z^n} \mathbb{E}_{\sigma^n} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right| = \mathbb{E} R_n(\mathcal{F}(Z^n)). \quad (41)$$

This completes the proof. \square

The above theorem implies the following key result on ERM:

Corollary 2. *For any $P \in \mathcal{P}(Z)$ and any n , the ERM hypothesis \hat{f}_n satisfies the bound*

$$P(\hat{f}_n) \leq L^*(\mathcal{F}) + 4\mathbb{E} R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{2 \log(\frac{1}{\delta})}{n}} \quad (42)$$

with probability at least $1 - \delta$.

Proof. From Theorem 1 it follows that, for any $t > 0$,

$$\mathbb{P}\left(\Delta_n(Z^n) \geq 2\mathbb{E} R_n(\mathcal{F}(Z^n)) + t\right) \leq \mathbb{P}\left(\Delta_n(Z^n) \geq \mathbb{E} \Delta_n(Z^n) + t\right). \quad (43)$$

The uniform deviation $\Delta(Z^n)$ has the bounded differences property with $c_1 = \dots = c_n = 1/n$. Hence, by McDiarmid's inequality

$$\mathbb{P}\left(\Delta_n(Z^n) \geq \mathbb{E} \Delta_n(Z^n) + t\right) \leq e^{-2nt^2}.$$

Letting $t = \sqrt{\frac{\log(1/\delta)}{2n}}$, we see that

$$\Delta_n(Z^n) \leq \mathbb{E} \Delta_n(Z^n) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

with probability at least $1 - \delta$. Together with (43), this implies that

$$\Delta_n(Z^n) \leq 2\mathbb{E} R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{\log(1/\delta)}{2n}} \quad (44)$$

with probability at least $1 - \delta$. Combining this with the first bound of Proposition 1, we conclude that

$$P(\hat{f}_n) \leq L^*(\mathcal{F}) + 4\mathbb{E} R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{2 \log(1/\delta)}{n}} \quad (45)$$

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

3 Structural results for Rademacher averages

The results developed above highlight the fundamental role played by Rademacher averages in bounding the generalization error of the ERM algorithm. In order to use these bounds, we need to get a better handle on the behavior of Rademacher averages.

Lemma 3 (Basic properties of Rademacher averages). *Let \mathcal{A} and \mathcal{B} be bounded subsets of \mathbb{R}^n , and let $c \in \mathbb{R}$ be a constant. Then*

$$R_n(\mathcal{A} \cup \mathcal{B}) \leq R_n(\mathcal{A}) + R_n(\mathcal{B}) \quad (46)$$

$$R_n(c\mathcal{A}) = |c|R_n(\mathcal{A}) \quad (47)$$

$$R_n(\mathcal{A} + \mathcal{B}) \leq R_n(\mathcal{A}) + R_n(\mathcal{B}), \quad (48)$$

where $c\mathcal{A} \triangleq \{ca : a \in \mathcal{A}\}$ and $\mathcal{A} + \mathcal{B} \triangleq \{a + b : a \in \mathcal{A}, b \in \mathcal{B}\}$. Moreover, let

$$\text{conv } \mathcal{A} \triangleq \left\{ \sum_{m=1}^N c_m a_m : N \in \mathbb{N}; a_m \in \mathcal{A}; c_m \geq 0, \forall m; \sum_{m=1}^N c_m = 1 \right\} \quad (49)$$

be the convex hull of \mathcal{A} and

$$\text{absconv } \mathcal{A} \triangleq \left\{ \sum_{m=1}^N c_m a_m : N \in \mathbb{N}; a_m \in \mathcal{A}; \sum_{m=1}^N |c_m| \leq 1 \right\} \quad (50)$$

be the absolute convex hull of \mathcal{A} . Then

$$R_n(\mathcal{A}) = R_n(\text{conv } \mathcal{A}) = R_n(\text{absconv } \mathcal{A}). \quad (51)$$

Proof. The proof is by direct calculation. First of all, by double-counting,

$$R_n(\mathcal{A} \cup \mathcal{B}) = \mathbb{E} \sup_{v \in \mathcal{A} \cup \mathcal{B}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i v_i \right| \quad (52)$$

$$\leq \mathbb{E} \sup_{a \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right| + \mathbb{E} \sup_{b \in \mathcal{B}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right| \quad (53)$$

$$= R_n(\mathcal{A}) + R_n(\mathcal{B}). \quad (54)$$

The case of $c\mathcal{A}$ is obvious. For $\mathcal{A} + \mathcal{B}$,

$$R_n(\mathcal{A} + \mathcal{B}) = \mathbb{E} \sup_{v \in \mathcal{A} + \mathcal{B}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i v_i \right| \quad (55)$$

$$= \mathbb{E} \sup_{a \in \mathcal{A}, b \in \mathcal{B}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i (a_i + b_i) \right| \quad (56)$$

$$\leq \mathbb{E} \sup_{a \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right| + \mathbb{E} \sup_{b \in \mathcal{B}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right| \quad (57)$$

$$= R_n(\mathcal{A}) + R_n(\mathcal{B}), \quad (58)$$

where the third step uses the triangle inequality.

Finally, consider the absolute convex hull of \mathcal{A} . Since $\mathcal{A} \subset \text{absconv } \mathcal{A}$, $R_n(\mathcal{A}) \leq R_n(\text{absconv } \mathcal{A})$. On the other hand, fix some $N \in \mathbb{N}$ and N real numbers c_1, \dots, c_N such that $\sum_{m=1}^N |c_m| = 1$, and consider the set

$$c_1 \mathcal{A} + \dots + c_N \mathcal{A} \equiv \{c_1 a_1 + \dots + c_N a_N : a_1, \dots, a_N \in \mathcal{A}\}. \quad (59)$$

Then

$$R_n(c_1 \mathcal{A} + \dots + c_N \mathcal{A}) \leq \sum_{i=1}^N |c_i| R_n(\mathcal{A}) \leq R_n(\mathcal{A}). \quad (60)$$

Since $\text{absconv } \mathcal{A}$ is the union of all sets of the form (59) for all choices of N and $\{c_m\}_{m=1}^N$, we see that $R_n(\text{absconv } \mathcal{A}) \leq R_n(\mathcal{A})$. Therefore, $R_n(\mathcal{A}) = R_n(\text{absconv } \mathcal{A})$. Since $\mathcal{A} \subset \text{conv } \mathcal{A} \subset \text{absconv } \mathcal{A}$, the same equality holds for the convex hull of \mathcal{A} . \square

The properties listed in the lemma show what happens to Rademacher averages when we form combinations of sets. This will be useful to us later, when we talk about hypothesis classes made up of simpler classes by means of operations like set-theoretic unions, intersections, complements or differences, logical connectives, or convex and linear combinations.

The next result, often referred to as the Finite Class Lemma, is very important:

Lemma 4 (Finite class lemma). *If $\mathcal{A} = \{a^{(1)}, \dots, a^{(N)}\} \subset \mathbb{R}^n$ is a finite set with $\|a^{(j)}\| \leq L$ for all $j = 1, \dots, N$ and $N \geq 2$, then*

$$R_n(\mathcal{A}) \leq \frac{2L\sqrt{\log N}}{n}. \quad (61)$$

Proof. Let σ^n be n i.i.d. Rademacher variables, and for every $j \in 1, \dots, N$ let

$$Y_j \triangleq \frac{1}{n} \sum_{i=1}^n \sigma_i a_i^{(j)}. \quad (62)$$

Then for any $s > 0$

$$\mathbb{E} e^{s Y_j} = \mathbb{E} \exp\left(\frac{s}{n} \sum_{i=1}^n \sigma_i a_i^{(j)}\right) = \prod_{i=1}^n \mathbb{E} e^{s \sigma_i a_i^{(j)}/n}, \quad (63)$$

where the second step uses the fact that σ^n are i.i.d. For each i , the random variable $\sigma_i a_i^{(j)}$ has zero mean and is bounded between $-a_i^{(j)}$ and $a_i^{(j)}$, so by the Hoeffding bound we have

$$\mathbb{E} e^{s \sigma_i a_i^{(j)}/n} \leq \exp\left(\frac{s^2 |a_i^{(j)}|^2}{2n^2}\right). \quad (64)$$

Therefore,

$$\mathbb{E}e^{sY_j} \leq \prod_{i=1}^n \exp\left(\frac{s^2 |a_i^{(j)}|^2}{2n^2}\right) = \exp\left(\frac{s^2}{2n^2} \sum_{i=1}^n |a_i^{(j)}|^2\right) = \exp\left(\frac{s^2 \|a^{(j)}\|^2}{2n^2}\right) \leq \exp\left(\frac{s^2 L^2}{2n^2}\right). \quad (65)$$

Repeating the same argument for each $-Y_j$, we see that

$$\mathbb{E}e^{-sY_j} \leq \exp\left(\frac{s^2 L^2}{2n^2}\right). \quad (66)$$

Now we recall the following statement, which was given as a homework problem in Spring 2011 at Duke University¹: Let U_1, \dots, U_K be K random variables (not necessarily independent) that are *subgaussian* with parameter $\nu > 0$, i.e.,

$$\mathbb{E}[e^{sU_k}] \leq e^{s^2 \nu^2 / 2}, \quad \forall s > 0. \quad (67)$$

Then

$$\mathbb{E}\left[\max_{1 \leq k \leq K} U_k\right] \leq \nu \sqrt{2 \log K}. \quad (68)$$

Consider now the $2N$ random variables $Y_1, -Y_1, \dots, Y_N, -Y_N$. According to (65) and (66), they are subgaussian with parameter $\nu = L/n$. Hence,

$$\mathbb{E}\left[\max_{1 \leq j \leq N} |Y_j|\right] = \mathbb{E}[\max(Y_1, -Y_1, \dots, Y_N, -Y_N)] \leq \frac{L\sqrt{2 \log(2N)}}{n} \leq \frac{2L\sqrt{\log N}}{n}, \quad (69)$$

where the last step uses the fact that, since $N \geq 2$, $2N \leq N^2$. Finally,

$$R_n(\mathcal{A}) = \mathbb{E}\left[\max_{1 \leq j \leq N} \left|\frac{1}{n} \sum_{i=1}^n \sigma_i a_i^{(j)}\right|\right] = \mathbb{E}\left[\max_{1 \leq j \leq N} |Y_j|\right] \leq \frac{2L\sqrt{\log N}}{n}, \quad (70)$$

which is what we wanted to prove. \square

We will start exploring the implications of the Finite Class Lemma more fully in the next lecture, but we can give a brief preview here. Consider a learning problem of the type described in Section 1 in the special case when \mathcal{F} consists of *binary-valued* functions on Z , i.e., $\mathcal{F} = \{f : Z \rightarrow \{0, 1\}\}$. From Theorem 1, we know that

$$\mathbb{E}\Delta_n(Z^n) \leq 2\mathbb{E}R_n(\mathcal{F}(Z^n)), \quad (71)$$

where

$$\mathcal{F}(Z^n) \triangleq \{(f(Z_1), \dots, f(Z_n)) : f \in \mathcal{F}\}. \quad (72)$$

¹See Problem 2 in http://maxim.ece.illinois.edu/teaching/spring11/homework/homework_1.pdf

Note that because each f can take values 0 or 1, $\mathcal{F}(Z^n) \subseteq \{0, 1\}^n$. Moreover, since for any $Z^n \in Z^n$ and any $f \in \mathcal{F}$ we have

$$\sqrt{\sum_{i=1}^n |f(Z_i)|^2} \leq \sqrt{n}, \quad (73)$$

the set $\mathcal{F}(Z^n)$ for a fixed Z^n satisfies the conditions of the Finite Class Lemma with $N = |\mathcal{F}(Z^n)| \leq 2^n$ and $L = \sqrt{n}$. Hence,

$$R_n(\mathcal{F}(Z^n)) \leq 2\sqrt{\frac{\log |\mathcal{F}(Z^n)|}{n}}. \quad (74)$$

In general, since $\log |\mathcal{F}(Z^n)| \leq n$, the bound just says that $R_n(\mathcal{F}(Z^n)) \leq 2$, which is not that useful. However, as we will see in the next few lectures, for a broad range of binary function classes \mathcal{F} it will not be possible to pick out every single element in $\{0, 1\}^n$ by taking the random “slices” $\mathcal{F}(Z^n)$, provided n is sufficiently large. To make these notions precise, let us define the quantity

$$\mathbb{S}_n(\mathcal{F}) \triangleq \sup_{z^n \in Z^n} |\mathcal{F}(z^n)|, \quad (75)$$

which is called the n th shatter coefficient of \mathcal{F} . Then we have the bound

$$R_n(\mathcal{F}(Z^n)) \leq 2\sqrt{\frac{\log \mathbb{S}_n(\mathcal{F})}{n}}. \quad (76)$$

Next, let

$$V(\mathcal{F}) \triangleq \max\{n \in \mathbb{N} : \mathbb{S}_n(\mathcal{F}) = 2^n\}. \quad (77)$$

This number is the famous *Vapnik–Chervonenkis* (or *VC dimension*) of \mathcal{F} , which has originated in their work [VC71]. It is clear that if $\mathbb{S}_n(\mathcal{F}) < 2^n$ for some n , then $\mathbb{S}_m(\mathcal{F}) < 2^m$ for all $m > n$. Hence, $V(\mathcal{F})$ is always well-defined (though it may be infinite). When it is finite, we say that \mathcal{F} is a *VC class*. What this means is that, for n large enough, a certain structure emerges in the sets $\mathcal{F}(z^n)$, which prevents us from being able to form any combination of binary labels by sweeping through the entire \mathcal{F} . A fundamental result, which was independently derived by Sauer [Sau72] and Shelah [She72] in different contexts (combinatorics and mathematical logic respectively) and also appeared in a weaker form in the original work of Vapnik and Chervonenkis [VC71], says the following:

Lemma 5 (Sauer–Shelah). *If \mathcal{F} is a VC class, i.e., $V(\mathcal{F}) < \infty$, then*

$$\mathbb{S}_n(\mathcal{F}) \leq \sum_{i=1}^{V(\mathcal{F})} \binom{n}{i} \leq (n+1)^{V(\mathcal{F})}. \quad (78)$$

Thus, we arrive at the following important result, which we will revisit in the next lecture:

Theorem 2. *If \mathcal{F} is a VC class of binary functions, then*

$$\mathbb{E}R_n(\mathcal{F}(Z^n)) \leq 2\sqrt{\frac{V(\mathcal{F})\log(n+1)}{n}}. \quad (79)$$

Consequently, for a VC class \mathcal{F} , the risk of ERM computed on an i.i.d. sample of size n from an arbitrary distribution $P \in \mathcal{P}(Z)$ is bounded by

$$P(\hat{f}_n) \leq L^*(\mathcal{F}) + 8\sqrt{\frac{V(\mathcal{F})\log(n+1)}{n}} + \sqrt{\frac{2\log(\frac{1}{\delta})}{n}} \quad (80)$$

with probability at least $1 - \delta$. In fact, using a much more refined technique called *chaining* originating in the work of Dudley [Dud78], it is possible to remove the logarithm in (79) to obtain the bound

$$\mathbb{E}R_n(\mathcal{F}(Z^n)) \leq C\sqrt{\frac{V(\mathcal{F})}{n}}, \quad (81)$$

where $C > 0$ is some universal constant independent of n and \mathcal{F} . We will not cover chaining in this class, but we will use the above formula.

References

- [Dud78] R. M. Dudley. Central limit theorems for empirical measures. *Annals of Probability*, 6:899–929, 1978.
- [GZ84] E. Giné and J. Zinn. Some limit theorems for empirical processes. *The Annals of Probability*, 12:929–989, 1984.
- [Sau72] N. Sauer. On the density of families of sets. *Journal of Combinatorial Theory, Series A*, 13:145–147, 1972.
- [She72] S. Shelah. A combinatorial problem: stability and order for models and theories in infinity languages. *Pacific Journal of Mathematics*, 41:247–261, 1972.
- [VC71] V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and Its Applications*, 16:264–280, 1971.