

## Handout 3: Learnability, stability

Lecturer &amp; author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

**Aim.** To introduce concepts PAC, fitting vs stability trade off, stability, and their implementation in regularization problems and convex problems.

### Reading list & references:

- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
  - Ch. 13 Stable rules do not over-fit
- Bousquet, O., Boucheron, S., & Lugosi, G. (2003). Introduction to statistical learning theory. In Summer school on machine learning (pp. 169-207). Berlin, Heidelberg: Springer Berlin Heidelberg. (Suitable for PG students)

### 1. LEARNABLE PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

**Definition 1.** (Agnostic PAC Learnability for General Loss Functions) A hypothesis class  $\mathcal{H}$  is agnostic PAC learnable with respect to a set  $\mathcal{S}$  and a loss function  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ , if there exist a function  $m_{\mathcal{H}} : (0, 1) \rightarrow \mathbb{N}$  and a learning algorithm  $\mathfrak{A}(\cdot)$  with the following property: for every  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 1)$ , and distribution  $g$  over  $\mathcal{Z}$ , when running algorithm  $\mathfrak{A}(\cdot)$  given  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  i.i.d examples generated by  $g$  then  $\mathfrak{A}(\cdot)$  returns  $h \in \mathcal{H}$  such as

$$(1.1) \quad \Pr \left( \hat{R}_{\mathcal{S}}(h) - \min_{h^* \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}}(h^*) \right) \leq \epsilon \right) \geq 1 - \delta$$

*Note 2.* It may be easier to work with expectations: by using Markov inequality (1.1) becomes

$$(1.2) \quad \Pr \left( \hat{R}_{\mathcal{Z}}(h) - \min_{h^* \in \mathcal{H}} \left( \hat{R}_{\mathcal{Z}}(h^*) \right) \leq \epsilon \right) \geq 1 - \underbrace{\frac{1}{\epsilon} \mathbb{E} \left( \hat{R}_{\mathcal{Z}}(h) - \min_{h' \in \mathcal{H}} \left( \hat{R}_{\mathcal{Z}}(h') \right) \right)}_{\delta}$$

and hence we need to work with expectations and bounded above.

**Hint:** Markov inequality  $\Pr(X \geq a) \leq \frac{1}{a} \mathbb{E}(X)$  for  $X > 0$ .

### 2. OVER-FITTING

*Note 3.* Following we discuss the association of stability and over-fitting. We give an (arguable) definition of over-fitting and based on this we show that: “a learning algorithm does not over-fit if and only if it is on-average-replace-one-stable”.

### 3. ANALYSIS BASED ON THE FITTING-STABILITY TRADE-OFF

*Note 4.* Let  $R^* = \min_{\mathcal{H}} (R(h))$  be an ideal/optimal (hence minimum) Risk. The Risk of a learning algorithm  $\mathfrak{A}(\mathcal{S})$  can be decomposed as

$$(3.1) \quad R_g(\mathfrak{A}(\mathcal{S})) - R^* = \underbrace{\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) - R^*}_{(I)} + \underbrace{R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))}_{(II)}$$

*Note 5.* Over-fitting is reasonably quantified by  $R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))$ . However,  $\hat{R}_{\mathcal{S}}(\cdot)$  is a random variable and for our computational convenience we focus on expectation w.r.t.  $\mathcal{S}$ . Hence, we provide the following (arguable) definition of over-fitting (on which we base our analysis).

**Definition 6.** For a learning algorithm  $\mathfrak{A}$ , as a measure of over-fitting we consider the expected difference between true and empirical risk

$$(3.2) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)$$

and we say that  $\mathfrak{A}$  suffers from over-fitting when (3.2) is ‘too’ large.

*Note 7.* The expected Risk of a learning algorithm  $\mathfrak{A}(\mathcal{S})$  can be written as

$$(3.3) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = \underbrace{\mathbb{E}_{\mathcal{S} \sim g} \left( \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)}_{(I)} + \underbrace{\mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)}_{(II)}$$

(by applying expectations in (3.1) and ignoring  $R^*$ ), where (I) indicates how well  $\mathfrak{A}(\mathcal{S})$  fits the training set  $\mathcal{D}$ , and (II) indicates discrepancy between the true and empirical risks of  $\mathfrak{A}(\mathcal{S})$ . In Section 4, we argue that (II) measures the over-fitting and stability of  $\mathfrak{A}(\mathcal{S})$ .

*Note 8.* The ultimate goal in a learning problem is to design a learning algorithm  $\mathfrak{A}(\mathcal{S})$  that both fit the training set and be stable; i.e. minimizing both terms in (3.3).

*Note 9.* As seen later, there may be a trade-off between (I) and (II); expected empirical risk term and stability term.

### 4. STABILITY AND OVER-FITTING

*Notation 10.* Let  $\mathcal{S} = \{z_1, \dots, z_m\}$  be a training sample, and  $\mathfrak{A}$  be a learning algorithm with output  $\mathfrak{A}(\mathcal{S})$ .

*Note 11.* Reasonably a learning algorithm can be stable if a small change of the input to the algorithm does not change the output of the algorithm much. Formalizing this in maths, we can say that if  $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$  is another training dataset equal to  $\mathcal{S}$  but the  $i$ th element which is replaced by another  $z' \sim g$ , then a good learning algorithm  $\mathfrak{A}$  would produce a small value of

$$\ell \left( \mathfrak{A}(\mathcal{S}^{(i)}), z_i \right) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \geq 0$$

**Definition 12.** We say that a learning algorithm  $\mathfrak{A}$  is **on-average-replace-one-stable** with rate  $\epsilon : \mathbb{N} \rightarrow \mathbb{R}$  if for every distribution  $g$

$$\mathbb{E} \left( \ell \left( \mathfrak{A}(\mathcal{S}^{(i)}), z_i \right) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right) \leq \epsilon(m)$$

where  $\epsilon(\cdot)$  has to be a decreasing function.

*Note 13.* Following we discuss the association of stability and over-fitting, based on the over-fitting Definition

**Definition 14.**<sup>1</sup> A learning algorithm  $\mathfrak{A}$  suffers from over-fitting if the expected difference between true and empirical risk

$$E_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)$$

is large.

**Theorem 16.** For any learning algorithm  $\mathfrak{A}$

$$E_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = E_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right)$$

where  $g$  is a distribution,  $\mathcal{S} = \{z_1, \dots, z_m\}$  and  $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$  are training datasets with  $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$ .

*Proof.* As  $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$ , then for every  $i$

$$E_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = E_{\substack{\mathcal{S} \sim g \\ z' \sim g}} (\ell(\mathfrak{A}(\mathcal{S}), z')) = E_{\substack{\mathcal{S} \sim g \\ z' \sim g}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) \right)$$

and

$$E_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) = E_{\substack{\mathcal{S} \sim g \\ i \sim U\{1, \dots, m\}}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) \right)$$

□

## 5. IMPLEMENTATION IN REGULARIZED LOSS LEARNING PROBLEMS

**Definition 17.** Assume  $(\mathcal{H}, \mathcal{Z}, \ell)$ . Regularized Loss Minimization (RLM) learning rule is the one that results as the output of jointly minimizing the empirical risk  $\hat{R}_{\mathcal{Z}}(h)$  and a regularization function  $J : \mathcal{H} \rightarrow \mathbb{R}$  that is

$$(5.1) \quad h^* = \underbrace{\arg \min}_{h \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}}(h) + J(h) \right)$$

*Remark 18.* The motivation for considering the regularization function  $J$  in (5.1) is to: (1.) control complexity and (2.) improve stability; as we will see later.

*Note 19.* We make our example more specific and narrow it to the Ridge RLM learning problem (could be LASSO, Elastic Net, etc.).

---

1

*Note 15.* Over-fitting is reasonably quantified by  $R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))$ . In Definition 14, we use the expectation for our mathematical convenience; other summaries/moments could be used instead.

**Definition 20.** The Ridge RLM learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$ , here  $\mathcal{H} = \mathcal{W} \subset \mathbb{R}^d$ , uses regularization function  $J(w; \lambda) = \lambda \|w\|_2^2$  with  $\lambda > 0$ ,  $w \in \mathcal{W}$  and produces learning rule

$$(5.2) \quad \mathfrak{A}(\mathcal{S}) = \arg \min_{w \in \mathcal{W}} \left( \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right)$$

*Note 21.* Recall (Term 1) that the regularization function in Ridge RLM learning problem penalizes complexity. Essentially, implies a sequence of hypothesis  $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots$  with  $\mathcal{H}_i = \{w \in \mathbb{R}^d : \|w\|_2 < i\}$ .

*Note 22.* Below, we will try to analyze the behavior of Ridge RLM learning rule (5.2) w.r.t. the Risk decomposition (3.3). In particular, to upper bounded w.r.t. the shrinkage term  $\lambda$ , training sample size  $m$ , and other characteristics.

### 5.1. Bounding the empirical risk (I) in (3.3).

*Note 23.* From (5.2), we have

$$\begin{aligned} \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) &\leq \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) + \lambda \|\mathfrak{A}(\mathcal{S})\|_2^2 \\ &\leq \hat{R}_{\mathcal{S}}(w') + \lambda \|w'\|_2^2; \quad \forall w' \in \mathcal{W} \end{aligned}$$

and by taking expectations w.r.t.  $\mathcal{S}$ , it is

$$(5.3) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w') + \lambda \|w'\|_2^2; \quad \forall w' \in \mathcal{W}$$

because  $\mathbb{E}_{\mathcal{S} \sim g} \left( \hat{R}_{\mathcal{S}}(\cdot) \right) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathcal{S} \sim g} (\ell(\cdot, z_i)) = R_g(\cdot)$ .

*Note 24.* We observe that part (I) in the expected risk decomposition (3.3), (aka the upper bound of the expected empirical risk) increases with the regularization term  $\lambda > 0$ . (!!!) –Although expected, it did not start well.

### 5.2. Bounding the empirical risk (II) in (3.3).

*Note 25.* We do that by constraining the loss function to be convex and Lipschitz.

**Assumption 26.** The loss function  $\ell(\cdot, z)$  in (5.2) is convex for any  $z \in \mathcal{Z}$ .

*Note 27.* Let  $\tilde{R}_{\mathcal{S}}(w) = \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2$ .  $\tilde{R}_{\mathcal{S}}(\cdot)$  is  $2\lambda$ -strongly convex as the sum of a convex function  $\hat{R}_{\mathcal{S}}(\cdot)$  (Assumption 26) and a  $2\lambda$ -strongly convex function  $J(\cdot; \lambda) = \lambda \|\cdot\|_2^2$  (results directly from Definition 40).

*Note 28.* Because  $\tilde{R}_{\mathcal{S}}(\cdot)$  is  $2\lambda$ -strongly convex and  $\mathfrak{A}_{\text{Ridge}}(\mathcal{S})$  is its minimizer, according to Lemma 42 in Handout (...), for  $\mathfrak{A}(\mathcal{S})$  and any  $w \in \mathcal{W}$ , it is

$$(5.4) \quad \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \geq \lambda \|w - \mathfrak{A}(\mathcal{S})\|_2^2, \quad \forall w \in \mathcal{W}$$

Note 29. Also, for any  $w, u \in \mathcal{W}$ , it is

$$\begin{aligned}\tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(u) &= \left( \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right) - \left( \hat{R}_{\mathcal{S}}(u) + \lambda \|u\|_2^2 \right) \\ &= \left( \hat{R}_{\mathcal{S}^{(i)}}(w) + \lambda \|w\|_2^2 \right) - \left( \hat{R}_{\mathcal{S}^{(i)}}(u) + \lambda \|u\|_2^2 \right) \\ &\quad + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(w, z') - \ell(u, z')}{m} \\ &= \tilde{R}_{\mathcal{S}^{(i)}}(w) - \tilde{R}_{\mathcal{S}^{(i)}}(u) + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(w, z') - \ell(u, z')}{m}\end{aligned}$$

Choosing  $w = \mathfrak{A}(\mathcal{S}^{(i)})$  and  $u = \mathfrak{A}(\mathcal{S})$ , and the fact that  $\tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}^{(i)})) \leq \tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}))$  it is

$$(5.5) \quad \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(u) \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z')}{m}$$

Note 30. Then (5.5) and (5.4) imply

$$(5.6) \quad \lambda \left\| \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}^{(i)})) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right\| \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z')}{m}$$

Note 31. Now that we brought it in that form, we can use an additional assumption on the loss to bound it.

**Assumption 32.** The loss function  $\ell(\cdot, z)$  in (5.2) is convex for any  $z \in \mathcal{Z}$  and  $\rho$ -Lipschitz.

Note 33. Given  $\rho$ -Lipschitzness in Assumption 32, it is

$$(5.7) \quad \begin{aligned}\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) &\leq \rho \left\| \mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S}) \right\| \\ \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z') &\leq \rho \left\| \mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S}) \right\|\end{aligned}$$

and hence (5.6) yields

$$(5.8) \quad \left\| \mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S}) \right\| \leq 2 \frac{\rho}{\lambda m}$$

Note 34. Plugging (5.8) in (5.7) yields

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 35. Using Theorem 16, we get an upper bound for the stability / over-fitting

$$\mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = \mathbb{E}_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 36. After this saga, the researcher could come to the conclusion that: In a Ridge regularization learning problem with loss function which is convex and  $\rho$ -Lipschitz, and the regularizer is  $J(\cdot; \lambda) = \lambda \|\cdot\|^2$  with  $\lambda > 0$ , the learning rule trained against iid sample  $\mathcal{S} = \{z_i\}_{i=1}^m$  is on-average-replace-one-stable with rate  $\epsilon(m) = 2 \frac{\rho^2}{\lambda m}$ ; i.e.

$$(5.9) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

*Note 37.* From (5.9), we see that stability improves (and over-fitting decreases) as the shrinkage parameter  $\lambda$  increases.

### 5.3. Bounding the Risk (3.3).

*Note 38.* Given the bounds (5.3) and (5.9), the decomposition of the expected Risk in (3.3) yields that: In a Ridge regularization learning problem with loss function which is convex and  $\rho$ -Lipschitz, and the regularizer is  $J(\cdot; \lambda) = \lambda \|\cdot\|^2$  with  $\lambda > 0$ , the learning rule trained against iid sample  $\mathcal{S} = \{z_i\}_{i=1}^m$  has expected Risk bound

$$(5.10) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \underbrace{R_g(w') + \lambda \|w'\|_2^2}_{(I)} + \underbrace{2 \frac{\rho^2}{\lambda m}}_{(II)}; \quad \forall w' \in \mathcal{W}$$

*Note 39.* From 5.10, we see that there is a trade-off between Empirical Risk (I) and stability (II) with regards the regularization parameter  $\lambda$ . We wish to use the optimal  $\lambda > 0$  corresponding to the smallest bound in (5.10); it has to both fit the training data well (but perhaps not too well) and be very stable to different training data from the same  $g$  (but perhaps not too stable)!

**Definition 40.** (Strongly Convex functions) A function  $f$  is  $\lambda$ -strongly convex function is for all  $w, u$ , and  $\alpha \in (0, 1)$  we have

$$(5.11) \quad f(\alpha w + (1 - \alpha)u) \leq \alpha f(w) + (1 - \alpha)f(u) - \frac{\lambda}{2} \alpha(1 - \alpha) \|w - u\|^2$$

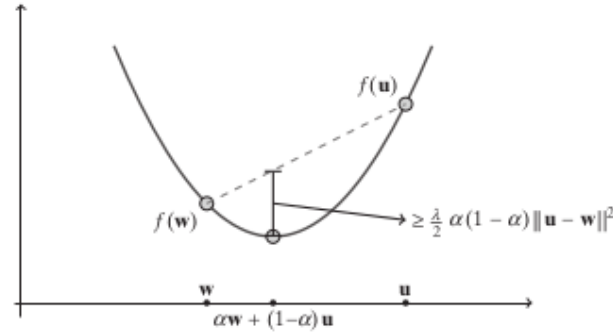


FIGURE 5.1. Strongly convex function

**Proposition 41.**

- (1) The function  $f(w) = \lambda \|w\|_2^2$  is  $2\lambda$ -strongly convex
- (2) If  $f$  is  $\lambda$ -strongly convex and  $g$  is convex then  $f + g$  is  $\lambda$ -strongly convex

**Lemma 42.** If  $f$  is  $\lambda$ -strongly convex and  $u$  is a minimizer of  $f$  then for any  $w$

$$f(w) - f(u) \geq \frac{\lambda}{2} \|w - u\|^2$$

*Proof.* Exercise 8 in the Exercise sheet. □