

## Handout 3: Learnability and stability in learning problems

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. 2, 3, 13
- Vapnik, V. (2000). The nature of statistical learning theory. Springer science & business media. (too advanced)

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

*Note 1.* We formally define the broad learning problem we will work on.

*Note 2.* Learning algorithms  $\mathcal{A}$  use training data sets  $\mathcal{S}$  which may miss characteristics of the unknown data-generating process  $g$ . Essentially, approximations (aka errors) are inevitable; some characteristics of data generating process  $g$  will be missed even if we use a very representative training set  $\mathcal{S}$ . We cannot hope that the learning algorithm will find a hypothesis whose error is smaller than the minimal possible error.

*Note 3.* The PAC learning problem requires no prior assumptions about the data-generating process  $g$ . It requires that the learning algorithm  $\mathcal{A}$  will find a predictor  $\mathcal{A}(\mathcal{S})$  whose error is not much larger than the best possible error of a predictor in some given benchmark hypothesis class. So essentially, in practice, the researcher's effort falls on the hypothesis class  $\mathcal{H}$ .

**Definition 4.** (Agnostic PAC Learnability for General Loss Functions) A hypothesis class  $\mathcal{H}$  is agnostic PAC learnable with respect to a domain  $\mathcal{Z}$  and a loss function  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ , if there exist a function  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  and a learning algorithm  $\mathcal{A}$  with the following property:

- for every  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 1)$ , and distribution  $g$  over  $\mathcal{Z}$ , when running algorithm  $\mathcal{A}$  given training set  $\mathcal{S}_m = \{z_1, \dots, z_m\}$  with  $z_i \stackrel{\text{iid}}{\sim} g$  for  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  then  $\mathcal{A}$  returns  $\mathcal{A}(\mathcal{S}) \in \mathcal{H}$  such as

$$(1.1) \quad \Pr_{\mathcal{S} \sim g} \left( R_g(\mathcal{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

*Note 5.* It may be easier to work with expectations: by using Markov inequality (1.1) becomes

$$(1.2) \quad \Pr_{\mathcal{S} \sim g} \left( R_g(\mathcal{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \frac{1}{\epsilon} \mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathcal{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \right)$$

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  $X \geq 0$ .

*Remark 6.* About 1.2: The accuracy parameter  $\epsilon$  determines how far the output rule  $\mathfrak{A}(\mathcal{S}_m)$  of  $\mathfrak{A}$  can be from the optimal one (this corresponds to the ‘approximately correct’). The confidence parameter  $\delta$  indicates how likely the classifier is to meet that accuracy requirement (corresponds to the “probably” part of “PAC”).

## 2. ANALYSIS OF THE RISK BASED ON THE TRADE-OFF FITTING VS STABILITY

*Note 7.* Let  $R^* = \min_{h \in \mathcal{H}} (R(h))$  be an ideal/optimal (hence minimum) Risk, and  $\mathfrak{A}(\mathcal{S})$  the learning rule from a learning algorithm  $\mathfrak{A}$  trained against dataset  $\mathcal{S}$ . The Risk of a learning algorithm  $\mathfrak{A}$  can be decomposed as

$$(2.1) \quad R_g(\mathfrak{A}(\mathcal{S})) - R^* = \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) - R^* + \underbrace{R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))}_{\text{over-fitting}}$$

*Note 8.* Over-fitting can be represented 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, here, for our computational convenience, we focus on its expectation w.r.t.  $\mathcal{S} \sim g$ . Hence, we provide the following (arguable) definition for over-fitting on which we base our analysis.

**Definition 9.** For a learning algorithm  $\mathfrak{A}$ , as a measure of over-fitting we consider the expected difference between true Risk and empirical Risk

$$(2.2) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))$$

**Definition 10.** We say that learning algorithm  $\mathfrak{A}$  suffers from over-fitting when (2.2) is ‘too’ large.

*Note 11.* The expected Risk of a learning algorithm  $\mathfrak{A}(\mathcal{S})$  can be decomposed as

$$(2.3) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = \underbrace{\mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{\text{(I)}} + \underbrace{\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{\text{(II)}}$$

by applying expectations in (2.1) and ignoring  $R^*$ . (I) indicates how well  $\mathfrak{A}(\mathcal{S})$  fits the training set  $\mathcal{S}$ , and (II) indicates the discrepancy between the true and empirical risks of  $\mathfrak{A}(\mathcal{S})$ . In Section 3, we argue that the over-fitting term (II) is directly related to a certain type of stability of  $\mathfrak{A}(\mathcal{S})$ .

*Note 12.* The following result connects the need for upper bounding (and minimizing this bound) the Expected Risk in (2.3) and PAC learning (Definition 4).

**Proposition 13.** Let  $\mathfrak{A}$  be a learning algorithm that guarantees the following:

- If  $m \geq m_{\mathcal{H}}(\epsilon)$  then for every distribution  $g$ , it is

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}_m))) \leq \min_{h' \in \mathcal{H}} (R_g(h')) + \epsilon$$

Then  $\mathfrak{A}$  satisfies the PAC guarantee in Definition 4:

- for every  $\delta \in (0, 1)$ , if  $m \geq m_{\mathcal{H}}(\epsilon\delta)$  then

$$\Pr_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

*Proof.* Let  $\xi = R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h'))$ . From Markov's inequality  $\Pr(\xi \geq \mathbb{E}(\xi)/\delta) \leq \frac{1}{\mathbb{E}(\xi)\delta} \mathbb{E}(\xi) = \delta$ . Namely,  $\Pr(\xi \leq \mathbb{E}(\xi)/\delta) = 1 - \delta$ . But it is given that if  $m \geq m_{\mathcal{H}}((\epsilon\delta))$  for every distribution  $g$  it is  $\mathbb{E}(\xi) \leq (\epsilon\delta)$ . So by substitution  $\Pr(\xi \leq (\epsilon\delta)/\delta) = 1 - \delta$  implies  $\Pr(\xi \leq \epsilon) = 1 - \delta$ . Now substitute back  $\xi$  and we conclude the proof.  $\square$

*Note 14.* Hence, we aim to design a learning algorithm  $\mathfrak{A}(\mathcal{S})$  that both fits the training set and is stable; i.e, keep both (I) and (II) in (2.3) small. As seen later, in certain learning problems, there may be a trade-off between empirical risk term (I) and (II) in (2.3). Consequently, we aim to upper bound (2.3) by upper bounding (I) and (II) individually and decide which term we should benefit against the other to achieve a certain total bound.

### 3. STABILITY AND ITS ASSOCIATION WITH OVER-FITTING

*Notation 15.* Consider a learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$ . 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 16.* It is reasonable to consider that a learning algorithm  $\mathfrak{A}$  can be stable if a small change of the algorithm input does not change the algorithm output much. Mathematically formalizing this, 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 being replaced by another independent example  $z' \sim g$ , then a good learning algorithm  $\mathfrak{A}$  would produce a small value of

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

**Definition 17.** A learning algorithm  $\mathfrak{A}$  is **on-average-replace-one-stable** with rate a decreasing function  $\epsilon : \mathbb{N} \rightarrow \mathbb{R}$  if for every distribution  $g$

$$(3.1) \quad \mathbb{E}_{\substack{\mathcal{S} \sim g \\ z' \sim g, i \sim U\{1, \dots, m\}}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right) \leq \epsilon(m)$$

*Note 18.* The following result demonstrates the direct association of stability and over-fitting as defined in Definitions 10 & 17.

**Theorem 19.** For any learning algorithm  $\mathfrak{A}$  it is

$$(3.2) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = \mathbb{E}_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} \left( \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right)$$

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  $i$

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

and

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

□

*Note 20.* Hence, we desire to design learning algorithms corresponding to as small as possible (3.2).

#### 4. IMPLEMENTATION IN REGULARIZED LOSS LEARNING PROBLEMS

**Definition 21.** A Regularized Loss Minimization (RLM) learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$  with a regularization function  $J : \mathcal{H} \rightarrow \mathbb{R}$  aims at a RLM learning rule that is the minimizer

$$(4.1) \quad h^* = \arg \min_{h \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}}(h) + J(h) \right)$$

where  $\hat{R}_{\mathcal{S}}(\cdot) = \frac{1}{m} \sum_{i=1}^m \ell(\cdot, z_i)$ , and  $\mathcal{S} = \{z_1, \dots, z_m\} \subset$  an IID training sample.

*Remark 22.* The motivation for considering the regularization function  $J$  in (4.1) is to: (1.) control complexity and (2.) improve stability; as we will see later.

*Note 23.* Here, we make our example more specific and narrow it to the Ridge RLM learning problem (could have been LASSO etc.).

**Definition 24.** The Ridge RLM learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$ , with  $\mathcal{H} = \mathcal{W} \subset \mathbb{R}^d$ , uses regularization function  $J(w; \lambda) = \lambda \|w\|_2^2$  with  $\lambda > 0$ ,  $w \in \mathcal{W}$  and produces learning rule

$$(4.2) \quad \mathfrak{A}(\mathcal{S}) = \arg \min_{w \in \mathcal{W}} \left( \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right)$$

*Note 25.* Recall (Term 1) how the regularization function in Ridge RLM learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$  penalizes complexity. Essentially, it implies a sequence of hypothesis  $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots$  with  $\mathcal{H}_i = \{w \in \mathbb{R}^d : \|w\|_2 < i\}$  due to duality of the corresponding minimization problem.

*Note 26.* Below, we will try to analyze the behavior of Ridge RLM learning rule (4.2) w.r.t. the Risk decomposition (2.3). In particular, to upper bounded w.r.t. the shrinkage term  $\lambda$ , training sample size  $m$ , and other characteristics.

##### 4.1. Bounding the empirical risk (I) in (2.3).

*Note 27.* From (4.2), 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

$$(4.3) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w) + \lambda \|w\|_2^2; \quad \forall w \in \mathcal{W}$$

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 28.* From (4.3), we observe that part (I) in the expected risk decomposition (2.3), (aka the upper bound of the expected empirical risk) increases with the regularization term  $\lambda > 0$ . (!!!)  
–Although anticipated, it did not start well.

#### 4.2. Bounding the empirical risk (II) in (2.3).

*Note 29.* As we will see to bound (II) in (2.3) we will have to impose an additional condition on the behavior of loss function apart from convexity; e.g. Lipschitzness.

**Assumption 30.** *The loss function  $\ell(\cdot, z)$  in (4.2) is convex for any  $z \in \mathcal{Z}$ .*

*Note 31.* Let  $\tilde{R}_S(w) = \hat{R}_S(w) + \lambda \|w\|_2^2$ .  $\tilde{R}_S(\cdot)$  is  $2\lambda$ -strongly convex as the sum of a convex function  $\hat{R}_S(\cdot)$  (Assumption 30) and a  $2\lambda$ -strongly convex function  $J(\cdot; \lambda) = \lambda \|\cdot\|_2^2$  (results directly from Definition 38 in Handout 2).

**Fact 32.** *(To be used in Note 33) If  $f$  is  $\lambda$ -strongly convex and  $u$  is a minimizer of  $f$  then for any  $w$*

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

*Note 33.* Let  $\mathfrak{A}(\mathcal{S})$  be the Ridge learning algorithm output minimizing (4.2). Because  $\tilde{R}_S(\cdot)$  is  $2\lambda$ -strongly convex and  $\mathfrak{A}(\mathcal{S})$  is its minimizer, according to Fact 32, for  $\mathfrak{A}(\mathcal{S})$  and any  $w \in \mathcal{W}$ , it is

$$(4.4) \quad \tilde{R}_S(w) - \tilde{R}_S(\mathfrak{A}(\mathcal{S})) \geq \lambda \|w - \mathfrak{A}(\mathcal{S})\|^2, \quad \forall w \in \mathcal{W}$$

*Note 34.* Also, for any  $w, u \in \mathcal{W}$ , it is

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

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

$$(4.5) \quad \tilde{R}_S(\mathfrak{A}(\mathcal{S}^{(i)})) - \tilde{R}_S(\mathfrak{A}(\mathcal{S})) \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\textcolor{red}{\mathcal{S}}), z') - \ell(\mathfrak{A}(\textcolor{red}{\mathcal{S}^{(i)}}), z')}{m}$$

*Note 35.* Then (4.4) and (4.5) imply

$$(4.6) \quad \lambda \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\|^2 \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\textcolor{red}{\mathcal{S}}), z') - \ell(\mathfrak{A}(\textcolor{red}{\mathcal{S}^{(i)}}), z')}{m}$$

*Note 36.* Now that we brought it in form (4.6), we can impose/use an additional assumption on the loss to bound it. In this example we use Lipschitzness.

**Assumption 37.** *The loss function  $\ell(\cdot, z)$  in (4.2) is convex and  $\rho$ -Lipschitz for any  $z \in \mathcal{Z}$ .*

*Note 38.* Given  $\rho$ -Lipschitzness in Assumption 37, it is

$$(4.7) \quad \ell \left( \mathfrak{A} \left( \mathcal{S}^{(i)} \right), z_i \right) - \ell \left( \mathfrak{A} (\mathcal{S}), z_i \right) \leq \rho \left\| \mathfrak{A} \left( \mathcal{S}^{(i)} \right) - \mathfrak{A} (\mathcal{S}) \right\|$$

$$(4.8) \quad \ell \left( \mathfrak{A} (\mathcal{S}), z' \right) - \ell \left( \mathcal{S}^{(i)}, z' \right) \leq \rho \left\| \mathfrak{A} \left( \mathcal{S}^{(i)} \right) - \mathfrak{A} (\mathcal{S}) \right\|$$

and hence plugging 4.7 and (4.8) in (4.6) yields

$$(4.9) \quad \left\| \mathfrak{A} \left( \mathcal{S}^{(i)} \right) - \mathfrak{A} (\mathcal{S}) \right\| \leq 2 \frac{\rho}{\lambda m}$$

*Note 39.* Plugging (4.9) in (4.7) yields

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

*Note 40.* Using Theorem 19, we get an upper bound for the stability / over-fitting

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

*Note 41.* After this saga, the researcher could come to the conclusion that: A Ridge RLM learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$  with the loss function which is convex and  $\rho$ -Lipschitz, regularizer  $J(\cdot; \lambda) = \lambda \|\cdot\|^2$ ,  $\lambda > 0$ , and learning rule trained against an iid sample  $\mathcal{S} = \{z_i\}_{i=1}^m$  from  $g$  is on-average-replace-one-stable with rate  $\epsilon(m) = 2 \frac{\rho^2}{\lambda m}$ ; i.e.

$$(4.10) \quad \mathbb{E}_{\mathcal{S} \sim g} \left( R_g (\mathfrak{A} (\mathcal{S})) - \hat{R}_{\mathcal{S}} (\mathfrak{A} (\mathcal{S})) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

*Note 42.* From (4.10), we see that stability improves (and over-fitting decreases) as the shrinkage parameter  $\lambda$  increases.

### 4.3. Bounding the Risk (2.3).

*Note 43.* Given the bounds (4.3) and (4.10), the decomposition of the expected Risk in (2.3) yields that: A Ridge RLM learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$  with the loss function which is convex and  $\rho$ -Lipschitz, regularizer  $J(\cdot; \lambda) = \lambda \|\cdot\|^2$ ,  $\lambda > 0$ , and learning rule trained against an iid sample  $\mathcal{S} = \{z_i\}_{i=1}^m$  from  $g$  has Expected Risk bound

$$(4.11) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \underbrace{R_g (w) + \lambda \|w\|_2^2}_{(I)} + \underbrace{2 \frac{\rho^2}{\lambda m}}_{(II)}; \quad \forall w \in \mathcal{W}$$

*Note 44.* From 4.11, we see that there is a trade-off between Empirical Risk (I) and stability/overfitting (II) with regards the regularization parameter  $\lambda$ . It is desirable to use the optimal  $\lambda > 0$  corresponding to the smallest bound in (4.11); 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)!

**Assumption 45.** Assume that the learning problem is additionally  $B$ -bounded by  $B > 0$ ; i.e.  $\mathcal{H} = \left\{ w \in \mathbb{R}^d : \|w\|_2^2 \leq B \right\}$ .

*Note 46.* If additionally the learning problem is  $B$ -bounded then the upper bound in (4.11) becomes

$$(4.12) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \underbrace{\min_{w \in \mathcal{H}} R_g (w)}_{(I)} + \underbrace{\lambda B + 2 \frac{\rho^2}{\lambda m}}_{(II)}$$

*Note 47.* Furthermore, if we choose a shrinkage parameter  $\lambda_m = \sqrt{\frac{2}{m}} \frac{\rho}{B}$  in (4.2) the upper bound in (4.11) becomes

$$(4.13) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g (w) + \lambda_m B \sqrt{\frac{8}{m}}$$

#### 4.4. Deriving PAC guarantees.

*Note 48.* In particular, if the size  $m$  of the training sample  $\mathcal{S}$  is

$$m \geq \frac{\lambda_m^2 B^2 8}{\epsilon^2} = \frac{8 \rho^2 B^2}{\epsilon^2}$$

for some  $\epsilon > 0$  then for every distribution  $g$ , the upper bound in (4.11) becomes

$$(4.14) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g (w) + \epsilon$$

and hence we can have a PAC guarantee that is summarized in the following.

*Summary 49.* Let  $(\mathcal{H}, \mathcal{Z}, \ell)$  be a learning problem convex-Lipschitz-bounded with parameters  $\rho$  and  $B$ , and let  $m$  be the training data size. The associated Ridge Regularized Loss Minimization rule with regularization function  $\lambda_m = \sqrt{\frac{2}{m}} \frac{\rho}{B}$  satisfies

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g (w) + \lambda_m B \sqrt{\frac{8}{m}}$$

Moreover, if  $m \geq \frac{8 \rho^2 B^2}{\epsilon^2}$ ,  $\epsilon > 0$  then for every distribution  $g$  we can get a PAC -guarantee

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g (\mathfrak{A} (\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g (w) + \epsilon$$

## APPENDIX A. RECALL STRONG CONVEX FUNCTIONS FROM HANDOUT 2.

- A function  $f$  is  $\lambda$ -strongly convex function is for all  $w, u$ , and  $\alpha \in (0, 1)$  we have

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

- The function  $f(w) = \lambda\|w\|_2^2$  is  $2\lambda$ -strongly convex
- If  $f$  is  $\lambda$ -strongly convex and  $g$  is convex then  $f + g$  is  $\lambda$ -strongly convex
- If  $f$  is  $\lambda$ -strongly convex and  $u$  is a minimizer of  $f$  then for any  $w$

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