

# MATH/COMP 562 LECTURE NOTES

## STABILITY THEORY

ADAM M. OBERMAN

### CONTENTS

1. Convex analysis	1
1.1. Lipschitz and strong convexity (for functions)	1
2. Stability Theory Setup	2
3. Stability Theory	2
3.1. Strongly convex losses	3
3.2. Define stability	3
3.3. Stability Theorem	3
4. Learning bounds for stable algorithms	5
4.1. Bounded differences	5
References	7

Notes adapted from [MRT18, Chapter 14] and [SSBD14, Ch 13]

### 1. CONVEX ANALYSIS

#### 1.1. Lipschitz and strong convexity (for functions).

**Definition 1.1.** The function  $f : W \rightarrow \mathbb{R}$  is Lipschitz continuous with constant  $C_f$ , if

$$|f(w_1) - f(w_2)| \leq C_f \|w_1 - w_2\|, \quad \forall w_1, w_2 \in W$$

**Definition 1.2.** The function  $f : W \subset \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\lambda$ -strongly convex (on  $W$ ), if  $f(w) - \lambda\|w\|^2$  is convex on  $W$ .

*Exercise 1.1.* Show that  $f$   $\lambda$  convex on  $\mathbb{R}^d$  means  $f(w) - \lambda\|w - w_0\|^2$  is convex for any  $w_0 \in \mathbb{R}^d$

**Lemma 1.3.** Suppose  $f$  is  $\lambda$ -strongly convex. Let  $w^*$  be a minimizer of  $f$ . Then

$$(SC) \quad f(w) - f(w^*) \geq \lambda \|w - w^*\|^2, \quad \forall w$$

*Remark 1.4.* Intuition of this: consider  $f : \mathbb{R} \rightarrow \mathbb{R}$  smooth. Then Taylor expansion around  $w^*$ :

$$f(w^* + v) = f(w^*) + \nabla f(w^*) \cdot v + v^2 f''(w^*)/2 + O(v^3)$$

Set  $\lambda = f''$  for a local version of the inequality. So the global version of this idea given by strong convexity.

## 2. STABILITY THEORY SETUP

We are considering a supervised learning problem (classification or regression). For now, consider the case of binary classification, or one dimensional regression. We may assume normalized vector data  $\mathcal{X} \subset [-1, 1]^d \subset \mathbb{R}^d$ . We assume labels are either  $\mathcal{Y} \subset \mathbb{R}$  or  $\mathcal{Y} = \{-1, 1\}$

**Definition 2.1.** Given the ML setup  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , and a class of functions  $\mathcal{H}$ ,  $h : \mathcal{X} \rightarrow \mathcal{Y}$ . A learning algorithm is a operator  $A$  which takes a finite subset,  $Z \subset \mathcal{Z}$  and returns a function  $h = A(Z)$ .

Note, in the parametric setting  $\mathcal{H} = \{h(x, w) \mid w \in W\}$ , so we have  $h_w(x) = A(Z)$ . We can simply write  $w = A(Z)$ .

Write  $z = (x, y)$  and for the dataset,  $S^m = \{(x_i, y_i)\}_{i=1}^m$  write

$$Z = (z_1, \dots, z_m) = ((x_1, y_1), \dots, (x_m, y_m))$$

**Definition 2.2** (bounded data). We say the learning problem  $(\mathcal{X}, \mathcal{Y})$  is bounded if there exist constants  $r_x, r_y$  such that

$$\|x\| \leq r_x, \quad \forall x \in \mathcal{X},$$

and

$$|y| \leq r_y, \quad \forall y \in \mathcal{Y}$$

We are given a loss function

$$\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}^+$$

*Example 2.3.* In the case of regression,

$$\ell(f, y) = (f - y)^2$$

*Example 2.4.* For classification, with,  $\mathcal{Y} = \{-1, +1\}$ , consider,

$$\ell(f, y) = \begin{cases} -\log(\sigma(f)), & y = +1 \\ -\log(1 - \sigma(f)), & y = -1 \end{cases}$$

## 3. STABILITY THEORY

**Definition 3.1.** Define the bounded linear hypotheses

$$\mathcal{H}_{lin, W} = \{h(x, w) = w \cdot x \mid x \in \mathcal{X}, w \in W\}$$

Let  $W \subset \mathbb{R}^d$  be a convex and bounded, with radius  $C_w$ .

Given a loss function  $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$ , define  $\ell : W \times \mathcal{Z} \rightarrow \mathbb{R}$  by

$$\ell(w, z) = \ell(w \cdot x, y)$$

Then, in this overloaded notation,

$$\ell(h_w(x), y) = \ell(w \cdot x, y) = \ell(w, z), \quad \forall h \in \mathcal{H}_{lin}$$

**3.1. Strongly convex losses.** Given a loss function  $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$ , define  $\ell_w : \mathcal{Z} \times \mathbb{R}^+ \rightarrow \mathbb{R}$  by

$$\ell_w(z, \lambda) = \ell(w, z, \lambda) = \ell(w \cdot x, y) + \lambda \|w\|^2$$

**Lemma 3.2.** Suppose  $\ell(s, y)$  is a convex function  $s$ , (for all  $z$ ), and  $\lambda \geq 0$ . Then  $\ell(w, z, \lambda)$  is  $\lambda$ -strongly convex function of  $w$ , for all  $z = (x, y)$ .

*Proof.* Class Notes / HW □

We say a loss is Lipschitz continuous if it is Lipschitz continuous as a function of  $w$ , independent of  $z$ .

**Definition 3.3.** The loss  $\ell : W \times Z \rightarrow \mathbb{R}$  is Lipschitz continuous with constant  $C_\ell$ , if

(Lip) 
$$|\ell(w_1, z) - \ell(w_2, z)| \leq C_\ell \|w_1 - w_2\|, \quad \forall w_1, w_2 \in W, z \in Z$$

**Definition 3.4.** Given a dataset  $Z = \{z_1, \dots, z_m\}$ , where each  $z_i \in \mathcal{Z}$ . Given a loss function  $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}^+$  Define for  $w \in W$ ,

$$L(w, Z, \lambda) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i) + \lambda \|w\|^2$$

**3.2. Define stability.**

*Remark 3.5.* In what follows, we will usually have  $\beta(m) = C/m$

**Definition 3.6** (Replace one Stability). Given  $\ell, \rho(z)$ . Let  $Z_1, Z_2$  be two datasets of size  $m$ , which differ in exactly one element.

The operator  $A$  is replace one stable in  $w$  if there exists  $C_w > 0$  such that

(S1) 
$$\|A(Z_1) - A(Z_2)\| = \|w_1 - w_2\| \leq \frac{C_w}{m}$$

The operator  $A$  is uniformly replace one stable in  $\ell$  with rate  $\beta(m)$ , if there exists a function  $\beta = \beta(m)$

(1) 
$$|\ell(w_1, z) - \ell(w_2, z)| \leq \beta(m), \quad \forall z \in \mathcal{Z}$$

The operator  $A$  is replace one stable in the expected loss, with rate  $\beta(m)$ , if there exists a function  $\beta = \beta(m)$  such that

(S2) 
$$L(A(Z_1)) - L(A(Z_2)) \leq \beta(m)$$

for all datasets  $Z_1, Z_2$  of size  $m$  which differ by only one element.

**3.3. Stability Theorem.** Define

$$w_i = A(Z_i) = \arg \min_{w \in W} L(w, Z_i, \lambda), \quad i = 1, 2$$

**Lemma 3.7.** Suppose the loss is  $C_\ell$  Lipschitz continuous. Then (S1) implies (1) and (S2) with  $C_L = C_\ell C_w$

*Proof.*

$$\begin{aligned} \ell(w_1, z) - \ell(w_2, z) &\leq C_\ell \|w_1 - w_2\| && \text{by (Lip)} \\ &\leq C_\ell C_w \frac{1}{m} && \text{by (S1)} \end{aligned}$$

The second result comes from taking expectations of the first inequality.

□

**Definition 3.8.** Given the loss,  $\ell$ , the dataset,  $Z$ , and  $\lambda > 0$ , define the regularized empirical loss

$$(REL) \quad L(w, Z, \lambda) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i) + \lambda \|w\|^2$$

The regularized loss minimization problem is to set

$$(RLM) \quad w = A(Z, \lambda) = \arg \min_w L(w, Z, \lambda)$$

**Theorem 3.9.** Given the loss  $\ell : W \times Z$  which is

- (1) convex in  $w$
- (2) Lipschitz continuous in  $w$ , as defined by (Lip)

Then the regularized loss minimization problem, (RLM), is replace one stable in  $w$  with constant  $C_w = C_\ell/\lambda$  and replace on stable in the expected loss, with constant  $C_L = C_\ell^2/\lambda$ .

*Proof.* Let  $S_1, S_2$  differ by one point,

$$(RO) \quad S_1 = \{z_1, \dots, z_{m-1}, z'_1\}, \quad S_2 = \{z_1, \dots, z_{m-1}, z'_2\}$$

Define  $w_1, w_2$  by (RLM),

$$w_i = w_i(S_i, \lambda) = \arg \min_w L(w, S_i, \lambda), \quad i = 1, 2$$

Since  $\ell$  is convex in  $w$ ,  $L(w, S, \lambda)$  is  $\lambda$ -strongly convex in  $w$ . Thus, applying strong convexity of (REL), and the definition of  $w_1, w_2$ , we obtain

$$\lambda \|w_1 - w_2\|^2 \leq L(w_2, S_1, \lambda) - L(w_1, S_1, \lambda) \quad \text{by (SC) and (RLM)}$$

$$\lambda \|w_1 - w_2\|^2 \leq L(w_1, S_2, \lambda) - L(w_2, S_2, \lambda) \quad \text{by (SC) and (RLM)}$$

Adding the two inequalities above,

$$2\lambda \|w_1 - w_2\|^2 \leq L(w_2, S_1, \lambda) - L(w_1, S_1, \lambda) + L(w_1, S_2, \lambda) - L(w_2, S_2, \lambda)$$

Using the fact that the datasets differ by only one point, we have

$$L(w_2, S_1, \lambda) - L(w_2, S_2, \lambda) = \frac{1}{m} (\ell(w_2, z'_1) - \ell(w_2, z'_2)) \quad \text{by (RO)}$$

$$L(w_1, S_2, \lambda) - L(w_1, S_1, \lambda) = \frac{1}{m} (\ell(w_1, z'_2) - \ell(w_1, z'_1)) \quad \text{by (RO)}$$

Combining the last three lines,

$$2\lambda \|w_1 - w_2\|^2 \leq \frac{1}{m} (\ell(w_2, z'_1) - \ell(w_2, z'_2) + \ell(w_1, z'_2) - \ell(w_1, z'_1))$$

Now apply the Lipschitz condition to obtain

$$\ell(w_2, z'_1) - \ell(w_1, z'_1) \leq C_\ell \|w_1 - w_2\| \quad \text{by (Lip) at } z'_1$$

$$-\ell(w_2, z'_2) + \ell(w_1, z'_2) \leq C_\ell \|w_1 - w_2\| \quad \text{by (Lip) at } z'_2$$

Combining, gives

$$2\lambda \|w_1 - w_2\|^2 \leq \frac{2C_\ell}{m} \|w_1 - w_2\|$$

Simplify the last inequality to obtain

$$\|w_1 - w_2\| \leq \frac{C_\ell}{\lambda m}$$

as desired. The second result follows directly from Lemma 3.7.  $\square$

*Remark 3.10.* This proof used a symmetric differences technique. We started by fixing the dataset and changing  $w$  to get the first inequality. Later we fixed the  $w$  and changed the dataset to go from  $L$  to  $\ell$ . Finally, we fixed the  $z$  and changed the  $w$  again to get another inequality.

#### 4. LEARNING BOUNDS FOR STABLE ALGORITHMS

**Definition 4.1.** Define

$$(2) \quad \Phi(S) = L(A(S)) - L_S(A(S))$$

to be the gap between the expected loss and the *training* loss of an algorithm.

*Remark 4.2.* The learning bounds for stable algorithms come from applying McDiarmid's inequality to  $\Phi$ . This is similar to what was done for Rademacher complexity bounds. We need to

- (1) Show that  $\Phi$  satisfies the bounded differences inequality
- (2) Apply McDiarmid's inequality to  $\Phi$
- (3) Bound the expected value of  $\Phi$  and relate this to the quantity of interest (in this case, stability).

##### 4.1. Bounded differences.

**Lemma 4.3.** *Let  $\Phi$  be defined by (2). Let  $A$  be a uniformly  $\beta$ -stable algorithm (in expectation). Suppose, in addition, that the loss is bounded by  $M$ . Then  $\Phi(S)$  satisfies the bounded difference inequality*

$$(3) \quad |\Phi(S_1) - \Phi(S_2)| \leq 2\beta(m) + \frac{M}{m}$$

(where  $S_1, S_2$  differ by one point).

*Proof of Lemma 4.3.* Using local notation, since copied from [MRT18].

Let  $\Phi$  be defined for all samples  $S$  by

$$\Phi(S) = R(h_S) - \widehat{R}_S(h_S)$$

Let  $S'$  be another sample of size  $m$  with points drawn i.i.d. according to  $\mathcal{D}$  that differs from  $S$  by exactly one point. We denote that point by  $z_m$  in  $S$ ,  $z'_m$  in  $S'$ , i.e.,

$$S = (z_1, \dots, z_{m-1}, z_m) \quad \text{and} \quad S' = (z_1, \dots, z_{m-1}, z'_m)$$

By definition of  $\Phi$ , the following inequality holds:

$$|\Phi(S') - \Phi(S)| \leq |R(h_{S'}) - R(h_S)| + \left| \widehat{R}_{S'}(h_{S'}) - \widehat{R}_S(h_S) \right|$$

We bound each of these two terms separately.

First, by the  $\beta$ -stability of  $\mathcal{A}$ , (S2), we have

$$\begin{aligned} |R(h_S) - R(h_{S'})| &= \left| \mathbb{E}_z [L_z(h_S)] - \mathbb{E}_z [L_z(h_{S'})] \right| \leq \mathbb{E}_z [|L_z(h_S) - L_z(h_{S'})|] \leq \beta \\ \left| \widehat{R}_S(h_S) - \widehat{R}_{S'}(h_{S'}) \right| &= \frac{1}{m} \left| \left( \sum_{i=1}^{m-1} L_{z_i}(h_S) - L_{z_i}(h_{S'}) \right) + L_{z_m}(h_S) - L_{z'_m}(h_{S'}) \right| \\ &\leq \frac{1}{m} \left[ \left( \sum_{i=1}^{m-1} |L_{z_i}(h_S) - L_{z_i}(h_{S'})| \right) + |L_{z_m}(h_S) - L_{z'_m}(h_{S'})| \right] \end{aligned}$$

Using the uniform  $\beta$ -stability of  $\mathcal{A}$ , (1), for the first terms, along with boundedness of  $L$ , for the last term, we have

$$\left| \widehat{R}_S(h_S) - \widehat{R}_{S'}(h_{S'}) \right| \leq \frac{m-1}{m} \beta(m) + \frac{M}{m} \leq \beta(m) + \frac{M}{m}$$

Thus,  $\Phi$  satisfies (3).  $\square$

In the next result, we take expectations of  $\Phi$ .

**Lemma 4.4.** *Let  $\Phi$  be defined by (2). Then*

$$\mathbb{E}_{S \sim \mathcal{D}^m} \Phi(S) \leq \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [|L_z(h_S) - L_z(h_{S'})|]$$

*Proof.* Local notation scope: we are copying from Mohri, so using his notation.

Rewrite

$$\Phi(S) = R(h_S) - \widehat{R}_S(h_S)$$

We now bound the expectation term, first noting that by linearity of expectation

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S[R(h_S)] - \mathbb{E}_S[\widehat{R}_S(h_S)]$$

By definition of the generalization error,

$$\mathbb{E}_{S \sim \mathcal{D}^m} [R(h_S)] = \mathbb{E}_{S \sim \mathcal{D}^m} \left[ \mathbb{E}_{z \sim \mathcal{D}} [L_z(h_S)] \right] = \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [L_z(h_S)]$$

By the linearity of expectation,

$$\mathbb{E}_{S \sim \mathcal{D}^m} [\widehat{R}_S(h_S)] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S \sim \mathcal{D}^m} [L_{z_i}(h_S)] = \mathbb{E}_{S \sim \mathcal{D}^m} [L_{z_1}(h_S)]$$

where the second equality follows from the fact that the  $z_i$  are drawn i.i.d. and thus the expectations  $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{z_i}(h_S)]$ ,  $i \in [m]$ , are all equal. The last expression in the equation above is the expected loss of a hypothesis on one of its training points. We can rewrite it as

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{z_1}(h_S)] = \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [L_z(h_{S'})]$$

where  $S'$  is a sample of  $m$  points containing  $z$  extracted from the  $m+1$  points formed by  $S$  and  $z$ . Thus,

$$\begin{aligned} \left| \mathbb{E}_{S \sim \mathcal{D}^m} [\Phi(S)] \right| &= \left| \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [L_z(h_S)] - \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [L_z(h_{S'})] \right| \\ &\leq \mathbb{E}_{S, z \sim \mathcal{D}^{m+1}} [|L_z(h_S) - L_z(h_{S'})|] \end{aligned}$$

as desired.  $\square$

Combining the lemma and the theorem, it means we can apply McDiarmid's inequality to  $\Phi(S)$ . Thus we have

**Theorem 4.5.** *Assume that the loss function  $L$  is bounded by  $M \geq 0$ . Let  $\mathcal{A}$  be a uniformly  $\beta$ -stable learning algorithm. Let  $S$  be a sample of  $m$  points drawn i.i.d. according to distribution  $\mathcal{D}$ . Then, with probability at least  $1 - \delta$  over the sample  $S$  drawn, the following holds:*

$$R(h_S) \leq \widehat{R}_S(h_S) + \beta + (2m\beta + M) \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

In particular, in the case  $\beta = C/m$ , we have

$$R(h_S) \leq \widehat{R}_S(h_S) + \frac{C}{m} + (2C + M) \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

*Proof.* Use the two previous results, to apply McD inequality.

Also apply the uniform stable definition. □

#### REFERENCES

- [MRT18] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [SSBD14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.