

## Lecture 9: Deriving stability of a classification

Lecturer: Abir De

Scribe: Mahadevan Subramanian

## 9.1 Defining Stability

Let us define sets of points  $S$  to be consisting of points from  $\mathbb{R}^n \times \{0, 1\}$ .  $S$  is a data set where each point has some binary classification to it. Let us define the set of all these data sets as  $\mathcal{S}$ .

We define an algorithm  $A : \mathcal{S} \rightarrow \mathbb{R}^d$  to be stable if  $\mathbf{Stab}_1(A)$  is a tautology.

$$\mathbf{Stab}_1(A) = \forall S \in \mathcal{S}, \forall S' \in \mathcal{S} \left[ (|S \setminus S'| = |S' \setminus S|) \wedge (|S| = |S'|) \rightarrow \left( \|A(S) - A(S')\| = \mathcal{O}\left(\frac{1}{|S|}\right) \right) \right]$$

The condition for this stability is for sets  $S$  and  $S'$  such that they differ only in element hence  $S' = e' \cup (S \setminus e)$  where  $e \neq e'$ .

Similar to this condition let us define a new condition  $\mathbf{Stab}_2(A)$ .

$$\mathbf{Stab}_2(A) = \forall S \in \mathcal{S}, \forall e \in S \left( \|A(S) - A(S \setminus e)\| = \mathcal{O}\left(\frac{1}{|S|}\right) \right)$$

**The overarching question:** For some  $A$ , do we have  $\mathbf{Stab}_1(A) \rightarrow \mathbf{Stab}_2(A)$ ?

**Lemma 9.1.** For all algorithms  $A : \mathcal{S} \rightarrow \mathbb{R}^d$ ,  $\mathbf{Stab}_2(A) \rightarrow \mathbf{Stab}_1(A)$

*Proof.* If we have  $\mathbf{Stab}_2(A)$  then  $\|A(S) - A(S \setminus e)\| = \mathcal{O}(1/|S|)$  and  $\|A(S') - A(S' \setminus e')\| = \mathcal{O}(1/|S'|)$ .

$$\begin{aligned} \|A(S) - A(S')\| &\leq \|A(S) - A(S \setminus e)\| + \|A(S' \setminus e') - A(S')\| \\ &= \mathcal{O}(1/|S|) + \mathcal{O}(1/|S'|) \\ &= \mathcal{O}\left(\frac{1}{|S|}\right) \end{aligned}$$

Hence we will have  $\mathbf{Stab}_1(A)$  hold given  $\mathbf{Stab}_2(A)$ . □

## 9.2 Applying stability to classification

Let us say we have a dataset  $D = \{(x_i, y_i)\}$ . Let us say we have some convex loss function  $l(w^T x, y)$  which is Lipschitz continuous. Let us define the following function over  $S \subset D$  which has regularization

$$F_w(S) = \sum_S (l(w^T x_i, y_i) + \lambda \|w\|^2)$$

Using this function we can define the following vector which minimizes the sum of the loss as

$$w^*(S) = \operatorname{argmin}_w F_w(S)$$

**Proposition 9.2.** For the defined  $F_w(S)$  with a convex and Lipschitz  $l(w^T x, y)$ ,  $\mathbf{Stab}_1(w^*)$  is true.

*Proof.* Let us define the notation  $l(w^*(S), e) = l(w^*(S)^T x, y)$ . Now we take a close look at the value  $F_{w^*(S')}(S) - F_{w^*(S)}(S)$ . We must have the following hold

$$F_{w^*(S')}(S) - F_{w^*(S)}(S) = F_{w^*(S')}(S') - F_{w^*(S)}(S') + l(w^*(S'), e) - l(w^*(S), e) + l(w^*(S), e') - l(w^*(S'), e')$$

Since  $w^*(S') = \operatorname{argmin}_w F_w(S')$  we have  $F_{w^*(S')}(S') - F_{w^*(S)}(S') \leq 0$  hence

$$F_{w^*(S')}(S) - F_{w^*(S)}(S) \leq l(w^*(S'), e) - l(w^*(S), e) + l(w^*(S), e') - l(w^*(S'), e') \leq 2L \|w^*(S) - w^*(S')\|$$

The last part of the inequality comes by combining the triangle inequality with the Lipschitz condition of  $l(w^*(S'), e) - l(w^*(S), e) \leq L \|w^*(S) - w^*(S')\|$ .

We can also expand  $F_{w^*(S')}(S) - F_{w^*(S)}(S)$  as a Taylor expansion about the point  $w^*(S)$ .

$$F_{w^*(S')}(S) - F_{w^*(S)}(S) = \left. \frac{\partial F_w(S)}{\partial w} \right|_{w=w^*(S)} (w - w^*(S)) + \frac{1}{2} (w - w^*(S))^T H(w - w^*(S)) + \dots$$

Here  $H(F_w(S))$  is the Hessian for the function  $F_w(S)$  with respect to  $w$ . We know that  $w^*(S)$  minimizes  $F_w(S)$  hence the first term vanishes and we are left with the inequality

$$F_{w^*(S')}(S) - F_{w^*(S)}(S) \geq \frac{1}{2} (w^*(S') - w^*(S))^T H(F_{w^*(S')}(S)) (w^*(S') - w^*(S))$$

We know that  $l(w, e)$  is a convex function hence the Hessian  $H(l(w, e))$  is positive semi-definite. Hence we can surely conclude that the Hessian of the sum of all  $l(w, e)$  terms is also positive semi-definite.

Now we can look at the regularization term, this will have to add a  $2\lambda|S|I$  to the Hessian by definition and so we can conclude that  $H(F_w(S)) \geq 2\lambda|S|I$  since the loss terms Hessian will anyways be positive semi-definite. Hence we have

$$F_{w^*(S')}(S) - F_{w^*(S)}(S) \geq \frac{2\lambda|S|}{2} (w^*(S') - w^*(S))^T (w^*(S') - w^*(S)) \geq \lambda|S| \|w^*(S') - w^*(S)\|^2$$

By combining the two inequalities we obtain by using first the Lipschitz condition and then that of convexity we obtain

$$\lambda|S| \|w^*(S') - w^*(S)\|^2 \leq F_{w^*(S')}(S) - F_{w^*(S)}(S) \leq 2L \|w^*(S) - w^*(S')\|$$

This subsequently reduces to

$$\|w^*(S') - w^*(S)\| \leq \frac{2L}{\lambda|S|} = \mathcal{O}\left(\frac{1}{|S|}\right)$$

Hence we have proven that with a convex and Lipschitz  $l(w^T x, y)$ ,  $\mathbf{Stab}_1(w^*)$  is true.

### 9.3 Group Details and Individual Contribution

Group 1 for the scribe for lecture 9. All work in this copy done by Mahadevan Subramanian, Roll no. 190260027.  $\square$