

Instructors: Purushottam Kar  
Authors: Gurpreet Singh  
Date: January 17, 2018

## Parametric Learning and Convergence

### Parametric Learning

Before we understand parametric learning theory, it is best to know the standard notation we will be using in this text.

#### Notations

1. **Instance/Feature Space** ( $\mathcal{X}$ ) — This is the space of the random variable that defines the data distribution.
2. **Output Space** ( $\mathcal{Y}$ ) — This defines the space of the label or value mapped to each data point. There can be different Output Spaces depending on the problem
  - Labelling:  $\{-1, 1\}, \mathbb{R}^p$
  - Permutation:  $S_n$
  - Hierarchy / Tree
  - Alternate Representation:  $\mathbb{R}^D \rightarrow \mathbb{R}^d$  (e.g. Dimensionality Reduction)
3. **Hypothesis/Model Space** ( $\mathcal{F}$ ) — This defines the set of all functions  $\{f : \mathcal{X} \rightarrow \mathcal{Y}\}$  belonging to a certain function class, for example, all linear functions or all neural networks. Generally equal to  $\mathcal{Y}^{\mathcal{X}}$ .
4. **Distribution** ( $\mathcal{D}$ ) — In case of noisy setting or agnostic setting, this is the probability distribution over  $\mathcal{X} \times \mathcal{Y}$ , typically unknown.  
However, if the setting is realizable (non-agnostic), then the distribution is only over  $\mathcal{X}$ , and we map  $f^* : x \mapsto y$ .
5. **Training Sample** ( $S$ ) — This is a tuple of  $n$  data points and their mapped label/parameter sampled from the distribution or randomly obtained.  $S = ((x^1, y^1), (x^2, y^2) \dots (x^n, y^n)) \in (\mathcal{X} \times \mathcal{Y})^n$   
If the setting is agnostic, then  $S \stackrel{iid}{\sim} \mathcal{D}$
6. **Loss Function** ( $l$ ) — Loss function defines the similarity / dissimilarity in the estimated label value and the actual / observed label value. We say  $l : (f(\mathcal{X}), \mathcal{Y}) \rightarrow \mathbb{R}$ . There are different types of loss functions, for example the squared loss function  $l^{\text{sq}}$  and the 0 – 1 loss function  $l^{0-1}$ .
7. **Algorithm** ( $A$ ) — We say that the algorithm maps the generated / sampled training sample to a function belonging in the hypothesis space *i.e.*

$$A : S \mapsto \hat{f}_s \in \mathcal{F}$$

$$A : \bigcup_{n=1}^{\infty} (\mathcal{X}, \mathcal{Y})^n \mapsto \mathcal{F}$$

8. **l-risk** ( $\text{er}_D^l[f]$ ) — This is defined as the expected loss/error that is obtained for a function  $f \in \mathcal{F}$  given a loss function  $l : (f(\mathcal{X}), \mathcal{Y}) \rightarrow \mathbb{R}$ .

$$\text{er}_D^l[f] \triangleq \mathbb{E}_{(x,y) \sim D} [l(f(x), y)]$$

**9. empirical risk** ( $\text{er}_D^l[f]$ ) — This is defined as the weighted loss/error on the training sample that is obtained for a function  $f \in \mathcal{F}$  given a loss function  $l : (f(\mathcal{X}), \mathcal{Y}) \rightarrow \mathbb{R}$

$$\text{er}_S^l[f] \triangleq \frac{1}{n} \sum_{i=1}^n l(f(x^i), y^i)$$

**Exercise 2.1.** Show that  $\text{er}_D^l[f] = \mathbb{E}[\text{er}_S^l[f]]$

**Exercise 2.2.** Show that  $\text{er}_D^{l^{0-1}}[f] = \text{er}_S^{0-1}[f]$  where

$$l^{0-1}(\hat{y}, y) = 1 - \mathbb{I}[y = \hat{y}]$$

**Exercise 2.3.** Find out  $\text{er}_D^{l_\alpha^{0-1}}[f]$  in terms of the empirical risk where

$$l_\alpha^{0-1} = \begin{cases} \alpha & \hat{y} \neq y, y = 1 \\ 1 - \alpha & \hat{y} \neq y, y = 0 \\ 0 & \hat{y} = y \end{cases}$$

## Toy Binary Classification Example

Assume we have a binary classification problem, with a finite hypothesis space  $\mathcal{F} = \{f_1, f_2 \dots f_m\}$ , where  $\forall f \in \mathcal{F}, f : \mathcal{X} \rightarrow \mathcal{Y}$

We first sample a training sample,  $S \stackrel{iid}{\sim} D^n$ . Then, we find the best model in  $\mathcal{F}$  *i.e.*

$$f^* = \arg \min_{f \in \mathcal{F}} \text{er}_D^{0-1}[f]$$

One possible solution is to estimate the l-risk using empirical risk. Hence, we now minimize empirical risk (Empirical Risk Minimization)

$$f^* \approx \hat{f} = \arg \min_{f \in \mathcal{F}} \text{er}_S[f]$$

## Pointwise and Uniform Convergence

We continue with our analysis of the toy binary classification example. As stated earlier, we had defined the best function  $f^*$  and  $\hat{f}$  as

$$\begin{aligned} f^* &= \arg \min_{f \in \mathcal{F}} \text{er}_D^l[f] \\ \hat{f} &= \arg \min_{f \in \mathcal{F}} \text{er}_S^l[f] \\ f^* &\approx \hat{f} \end{aligned}$$

Since our hypothesis space is finite, we can, without loss of generality, say that  $f^* = f_1$ . Then,

- (i) we do not want  $f_2, f_3 \dots f_m$  to perform well on  $S$  *i.e.* have high(er) empirical risk or training error.
- (ii) we do not want  $f_1$  to perform poorly on  $S$  *i.e.* have low(er) empirical risk.
- (iii) we try to ensure all  $f_i$  give faithful and honest performance on  $S$  *i.e.*  $\forall f \in \mathcal{F}, \text{er}_D^l[f] \approx \text{er}_S^l[f]$

We say that we need  $f_1$  to give good performance, however we need to define what is “good”. We say that  $S$  is good with respect to a function  $f \in \mathcal{F}$  ( $S \in \text{good}_f(\epsilon)$ ) if

$$|\text{er}_D^l[f] - \text{er}_S^l[f]| \leq \epsilon$$

This is known as pointwise convergence. We define it more formally below.

**Definition 2.1.** We say that a hypothesis class,  $\mathcal{F}$  has pointwise convergence for a given  $\epsilon$  and a loss function  $l$ , if for all functions  $f \in \mathcal{F}$ , and a sample  $S$  of size  $n$ , the property

$$|\text{er}_D^l[f] - \text{er}_S^l[f]| \leq \epsilon$$

holds true.

**Theorem 2.1.** If  $S \stackrel{iid}{\sim} D^n$ , then  $\forall f \in \mathcal{F}$ ,

$$\mathbb{P}[S \in \text{good}_f(\epsilon)] \geq 1 - 2 \exp\left(\frac{-n\epsilon^2}{3}\right)$$

or equivalently

$$\mathbb{P}[|\text{er}_D^l[f] - \text{er}_S^l[f]| > \epsilon] \leq 2 \exp\left(\frac{-n\epsilon^2}{3}\right)$$

**Exercise 2.4.** Prove Theorem 2.1

*Hint: Take a Bernoulli random variable and use Chernoff's Bound*

Pointwise convergence is a good property for individual functions, however, we want such a property to hold for all functions collectively. Hence, we define a much stronger convergence property, called uniform convergence.

In order to define Uniform Convergence, we must first define a desired situation for the defined hypothesis class.

$$\mathbb{P}[\text{er}_D^l[\hat{f}] > \text{er}_D^l[f^*]] < \delta$$

First, let us define the “goodness” of a sample in terms of uniform convergence. We say a sample  $S$  is good *i.e.*  $S \in \text{good}(\epsilon)$  if for all functions  $f \in \mathcal{F}$ ,  $S \in \text{good}_f(\epsilon)$

Since we define  $\hat{f} = \arg \min_{f \in \mathcal{F}} \text{er}_S^l[f]$ , we can say

$$\text{er}_S^l[\hat{f}] \leq \text{er}_S^l[f^*]$$

Suppose  $S \in \text{good}(\epsilon)$ , then

$$\begin{aligned} \text{er}_D^l[\hat{f}] &\leq \text{er}_S^l[\hat{f}] + \epsilon \\ &\leq \text{er}_S^l[f^*] + \epsilon \\ &\leq \text{er}_D^l[f^*] + 2\epsilon \end{aligned}$$

Hence, we can say that if

$$\begin{aligned} S \in \text{good}(\epsilon) &\implies \text{er}_D^l[\hat{f}] \leq \text{er}_D^l[f^*] + \frac{\epsilon}{2} \\ \therefore \text{er}_D^l[\hat{f}] > \text{er}_D^l[f^*] + \frac{\epsilon}{2} &\implies S \notin \text{good}(\epsilon) \end{aligned}$$

Here, we use an identity, the proof of which is left as an exercise.

$$\text{If } A \implies B \text{ then } \mathbb{P}[A] \leq \mathbb{P}[B]$$

Therefore, we can say

$$\mathbb{P}\left[\text{er}_D^l[\hat{f}] > \text{er}_D^l[f^*]\right] \leq \mathbb{P}\left[S \notin \text{good}\left(\frac{\epsilon}{2}\right)\right]$$

We can further reduce the RHS of the above inequality in terms of the size of the sample  $S$ ,  $\epsilon$  and the size of the hypothesis set, and hence state the condition for uniform convergence.

**Theorem 2.2.** A hypothesis function,  $\mathcal{F}$  is said to have the uniform convergence property if for some sample  $S$  of size  $n$ , and some  $\epsilon$

$$\mathbb{P}\left[\left|\text{er}_D^l[\hat{f}] - \text{er}_D^l[f^*]\right| > 2\epsilon\right] \leq \mathbb{P}[S \notin \text{good}(\epsilon)] \leq 2 \exp\left(\frac{-n\epsilon^2}{3}\right) \cdot m$$

**Proof.** From earlier, we know that if pointwise convergence holds, then uniform convergence also holds

$$\begin{aligned} \forall f \in \mathcal{F}, S \in \text{good}_f(\epsilon) &\implies S \in \text{good}(\epsilon) \\ \therefore \mathbb{P}[S \in \text{good}(\epsilon)] &\geq \mathbb{P}[\forall f \in \mathcal{F}, S \in \text{good}_f(\epsilon)] \\ \implies \mathbb{P}[S \notin \text{good}(\epsilon)] &\leq \mathbb{P}[\exists f \in \mathcal{F}, S \notin \text{good}_f(\epsilon)] \end{aligned}$$

Since on the RHS, we just have a union, using basic probability theory inequalities, we can write

$$\begin{aligned} \mathbb{P}[S \notin \text{good}(\epsilon)] &\leq \sum_{k=1}^m \mathbb{P}[S \notin \text{good}_{f_k}(\epsilon)] \\ &\leq \sum_{k=1}^m 2 \exp\left(\frac{-n\epsilon^2}{3}\right) \\ &= 2 \exp\left(\frac{-n\epsilon^2}{3}\right) \cdot m \end{aligned}$$

Hence, we can say that the theorem is valid. □

<++>