

# References

## 1. Pattern Classification

Duda, Hart, Stork

## 2. FOUNDATIONS OF MACHINE LEARNING [FOML]

Mohri et al.

[select portions]

$$\mathcal{D} = \left\{ (x^{(i)}, y^{(i)}) \mid x^{(i)} \in \mathcal{X}, y^{(i)} \in \{-1, 1\} \right. \\ \left. i \in [N] \right\}$$

$$\mathcal{D} \sim \mathcal{P}^{(N)}$$


---

Assume that there exist  
 $\omega^*$  such that

$$\text{sign}(\omega^*{}^T x^{(i)}) = y^{(i)}$$

$$i \in [N].$$

$$\|x^{(i)}\| \leq R$$

$$\omega^{(i)} \neq 0$$

$$i = 1, \dots, N$$

$$\Rightarrow y^{(i)} (\omega^*{}^T x^{(i)}) > 0$$

$$\gamma = \min_i \frac{y^{(i)} (\omega^*{}^T x^{(i)})}{\|\omega^*\|}$$

$\omega^{(n)}$  be the current estimate.

Let  $(x^{(n)}, y^{(n)}) \in \mathcal{D}$  such that

$$\text{sign}((\omega^{(n)})^T x^{(n)}) \neq y^{(n)}$$

$$\Rightarrow y^{(n)} (\omega^{(n)T} x^{(n)}) < 0$$

$$\omega^{(n+1)} = \begin{cases} \omega^{(n)} + y^{(n)} x^{(n)} & [\text{update}] \\ \omega^{(n)} & \text{otherwise} \end{cases}$$

On a linearly separable dataset  
the perceptron algorithm  
terminates after making  
at most  $\frac{R^2}{\gamma^2}$  updates

$$\frac{R^2}{\gamma^2} \leq \frac{R^2 \|\omega^*\|^2}{\min_i \{y_i (\omega^{*T} x^{(i)})\}^2}$$

Generalization error of  
the classifier returned by  
Perceptron?

$$\bar{\omega} = \sum_{k=1}^M y^{(i_k)} x^{(i_k)}$$

$$\begin{array}{l} i \in \{1, \dots, N\} \\ \text{sign}(\bar{\omega}^T x^{(i)}) \\ = y^{(i)} \end{array}$$

$V = \{i_1, \dots, i_M\} \rightarrow$  Indices which  
are updated

$\mathcal{D}$  be a Sample of size  $N$ .

Linearly Separable [There exists  $v > 0$ ].

Let Perceptron Algorithm return  
a Classifier  $h_{\mathcal{D}}^{(P)}$ .

$$P(h_{\mathcal{D}}^{(P)}(x) \neq y) = R(h_{\mathcal{D}}^{(P)})$$

$$E_{\mathcal{D} \sim P^{[N]}} R(h_{\mathcal{D}}^{(P)})$$

Expected generalization  
error by a classifier  
returned by Perceptron algorithm  
acting on a Linearly Separable  
Sample of  $N$  i.i.d. draws from  $P$ .

$$E_{\mathcal{D} \sim P^{[N]}} R(h_{\mathcal{D}}^{(P)})$$

Hack:

- ①  $\bar{\mathcal{D}} \sim P^{(N+1)}$  Sample of size  $N+1$ .
- ② Let it be linearly separable
- ③ Create  $N+1$  Datasets from  $\bar{\mathcal{D}}$  by removing  $i$ th datapoint.  
 $\mathcal{D}^{(i)} = \bar{\mathcal{D}} - \{x^{(i)}, y^{(i)}\} \quad i=1, \dots, N+1$

- Size of  $\mathcal{D}^{(i)} = N$
- $\mathcal{D}$  is linearly separable
- $\Rightarrow \mathcal{D}^{(i)}$  is linearly separable.

An algorithm  $A$  acting on a sample  $D$  of size  $m$  return a classifier  $h_D^{(A)}$ .

The leave one out error of  $A$  on  $D$  is

$$\bar{R}_D^{\text{LOO}}(A) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h_{D^{(i)}}^{(A)}(x^{(i)}) \neq y^{(i)}\}}$$

$$D^{(i)} : D - \{(x^{(i)}, y^{(i)})\}$$

$$E_{D \sim \mathcal{P}^m} \bar{R}_D^{\text{LOO}}(A)$$

Expected LOO error. on a random sample of size  $m$ .

$$E_{\mathcal{D} \sim p(m)} \bar{R}_{\mathcal{D}}^{\text{Loo}}(A)$$

$$= \frac{1}{m} E_{\mathcal{D} \sim p(m)} \sum_{i=1}^m \mathbb{1}_{\{h_{\mathcal{D}^{(i)}}^{(A)}(x^{(i)}) \neq y^{(i)}\}}$$

$$E_{\mathcal{D} \sim p(m)} = E_{\mathcal{D}^{(i)} \sim p(m-1)} E_{x^{(i)}, y^{(i)} \sim p}$$

IID



$$E_{X^{(i)}, Y^{(i)} \sim \mathcal{P}^{(A)} \{h_{\mathcal{Q}^{(i)}}^{(A)}(x^{(i)}) \neq y^{(i)}\}}$$

$$= P(h_{\mathcal{Q}^{(i)}}^{(A)}(x^{(i)}) \neq y^{(i)})$$

$$= R(h_{\mathcal{Q}^{(i)}}^{(A)})$$

$$E_{\mathcal{Q} \sim \mathcal{P}^{(m)}} \bar{R}_{\mathcal{Q}}^{(A)} \stackrel{\text{LOO}}{=}$$

$$= \frac{1}{m} \sum_{i=1}^m E_{\mathcal{Q}^{(i)} \sim \mathcal{P}^{(m-1)}} R(h_{\mathcal{Q}^{(i)}}^{(A)})$$

$$= E_{\mathcal{Q} \sim \mathcal{P}^{(m-1)}} R(h_{\mathcal{Q}}^{(A)})$$

Average  
 Generalization error  
 of a Classifier derived  
 by an Algorithm  $A$   
 acting on a sample of  
 size  $m-1$

$$E_{\mathcal{D} \sim \mathcal{P}^{(m-1)}} R(h_{\mathcal{D}}^{(A)})$$

= Expected LOO error  
 of  $A$  on sample  $m$ .

$$E_{\mathcal{Q} \sim P(m)} \bar{R}_{\mathcal{Q}}^{LOO}(A)$$

$$= E_{\mathcal{Q} \sim P(m-1)} R(h_{\mathcal{Q}}^{(A)})$$

Reading : Sec 5.2.4.  
 Lemma 5.3.  
 FOML

# Perceptron

$$E_{\mathcal{D} \sim P(N)} R(h_{\mathcal{D}}^{(P)})$$

$$= E_{\mathcal{D} \sim P(N+1)} \bar{R}_{\mathcal{D}}^{\text{LOO}}(P)$$

Consider linearly separable  
sample  $\mathcal{D}$  of size  $N+1$

Let the perceptron make  
 $M$  updates using the indices from

$$U = \{i_1, i_2, \dots, i_M\}$$

$$h_{\mathcal{D}}^{(p)}(x) = \text{sign}(\bar{\omega}_D^T x)$$

$$\bar{\omega}_D = \sum_{k=1}^M y^{(i_k)} x^{(i_k)}$$

$$\mathcal{D}^{(j)} = \mathcal{D} - \{x^{(j)}, y^{(j)}\}$$

forall  $j \in [N]$

$$\text{sign}(\bar{\omega}_D^T x^{(j)}) = y^{(j)}$$

$$\uparrow \{h_{\mathcal{D}}^{(p)}(x^{(j)}) \neq y^{(j)}\} = 0$$

(I)

if  $j \notin U$

$$h_{\mathcal{Q}(j)}^{(P)}(x) = h_{\mathcal{Q}}^{(P)}(x)$$

$$\Rightarrow \mathbb{1}_{\{h_{\mathcal{Q}(j)}^{(P)}(x^{(j)}) \neq y^{(j)}\}} = 0$$

(II)

Because of (I)

If  $j \in U$

$$\mathbb{1}_{\{h_{\mathcal{Q}(j)}^{(P)}(x^{(j)}) \neq y^{(j)}\}} \leq 1$$

(III)

$$\overline{R}_Q^{LOO}(P) =$$

$$\frac{1}{N+1} \sum_{j=1}^{N+1} \mathbb{1}_{\{h_Q^{(P)}(x^{(j)}) \neq y^{(j)}\}}$$

$$\leq \frac{1}{N+1} \left( \underbrace{\sum_{j \notin U} 0}_{\text{II}} + \underbrace{\sum_{j \in U} 1}_{\text{III}} \right) = \frac{M}{N+1}$$

$$E_{\mathcal{D}} \sim p(N) R(h_{\mathcal{D}}^{(P)})$$

$$\leq E_{\mathcal{D}} \sim p(N+1) \frac{\min(M(\mathcal{D}), \frac{R^2(\mathcal{D})}{\gamma^2(\mathcal{D})})}{N+1}$$

[Theorem 8.9]  
FOML]