

# Statistical Learning Theory Notes

Nong Minh Hieu<sup>1</sup>

<sup>1</sup> School of Physical and Mathematical Sciences, Nanyang Technological University (NTU - Singapore)

## Contents

<b>1</b>	<b>Probability settings</b>	<b>3</b>
1.1	Classification problem . . . . .	3
1.2	Goal of classification . . . . .	5
<b>2</b>	<b>Bayes classifier</b>	<b>6</b>
2.1	Properties of Bayes Risk . . . . .	6
2.2	Likelihood Ratio Test . . . . .	8
2.3	Plug-in classifier . . . . .	9
2.4	End of chapter exercises . . . . .	11
<b>3</b>	<b>Hoeffding's inequality</b>	<b>14</b>
3.1	Markov's Inequality . . . . .	14
3.2	Hoeffding's Inequality . . . . .	15
3.3	Convergence of Empirical Risk . . . . .	16
3.4	KL-divergence & Hypothesis Testing . . . . .	17
3.5	End of chapter exercises . . . . .	20
<b>4</b>	<b>Empirical Risk Minimization</b>	<b>22</b>
4.1	Uniform Deviation Bounds . . . . .	22
4.2	PAC Learning & Sample Complexity . . . . .	25
4.3	Zero-error case . . . . .	25
4.4	End of chapter exercises . . . . .	28
<b>5</b>	<b>Vapnik-Chevronenkis Theory</b>	<b>31</b>
5.1	VC Dimension . . . . .	31
5.2	Sauer's Lemma . . . . .	32
5.3	VC Theorem for classifiers . . . . .	34
5.4	VC Classes . . . . .	36
5.5	VC Theorem for sets . . . . .	38
5.6	End of chapter exercises . . . . .	41
<b>6</b>	<b>Rademacher Complexity</b>	<b>43</b>
6.1	Bounded Difference Inequality . . . . .	43
6.2	Rademacher Complexity . . . . .	43
6.3	Bounds for binary classification . . . . .	43
6.4	Proof of VC Inequality . . . . .	43
<b>A</b>	<b>Related topics</b>	<b>44</b>
A.1	Neyman-Pearson Lemma . . . . .	44
A.1.1	Type I & Type II errors . . . . .	44
A.1.2	Neyman-Pearson Lemma . . . . .	45

<b>B</b>	<b>List of Definitions</b>	<b>47</b>
<b>C</b>	<b>Important Theorems</b>	<b>47</b>
<b>D</b>	<b>Important Corollaries</b>	<b>47</b>
<b>E</b>	<b>Important Propositions</b>	<b>48</b>
<b>F</b>	<b>References</b>	<b>49</b>

# 1 Probability settings

## 1.1 Classification problem

**Definition 1.1** (Classifier ( $h$ )).

In **classification problems**, we consider pairs  $(x, y)$  where  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Where:

- $\mathcal{X}$  is the space of **feature vectors**.
- $\mathcal{Y}$  is the space of **labels**.

A classifier is a function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  which aims to assign correct labels to given feature vectors.

**Remark :** The key assumptions of classification problems are:

- There exists a joint distribution  $P_{XY}$  on  $\mathcal{X} \times \mathcal{Y}$ .
- The pairs  $(x, y)$  (observed data) are random samples of the random variables pair  $(X, Y)$  which has the distribution  $P_{XY}$ .

**Definition 1.2** (Decomposition of  $P_{XY}$ ).

We can decompose  $P_{XY}$  in either of the following two ways:

$$\begin{aligned} P_{XY} &= P_{X|Y} P_Y \\ P_{XY} &= P_{Y|X} P_X \end{aligned}$$

Which can be understood as two possible ways to generate the pairs  $(x, y)$  from the joint distribution  $P_{XY}$ .

- The first way is to generate a random label  $y \sim P_Y$ . Then, generate the feature vector corresponding to that label  $x \sim P_{X|Y=y}$ .
- The second way is to generate a random vector  $x \sim P_X$ . Then, generate the label corresponding to that feature vector  $y \sim P_{Y|X=x}$ .

### Proposition 1.1: Law of total expectation

Given  $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ . The **law of total expectation** states that:

$$\begin{aligned} \mathbb{E}_{XY} [\phi(X, Y)] &= \mathbb{E}_Y [\mathbb{E}_{X|Y} [\phi(X, Y)]] \\ &= \mathbb{E}_X [\mathbb{E}_{Y|X} [\phi(X, Y)]] \end{aligned}$$

Similar to how  $P_{XY}$  is decomposed, law of total expectation describes two way of taking the average value:

- Loop through the labels and take average over the feature vectors corresponding to each label.
- Loop through the feature vectors and take average over the labels corresponding to each vector.

**Proof** (Proposition 1.1). 

---

We have:

$$\begin{aligned}
\mathbb{E}_{XY}[\phi(X, Y)] &= \int_{\mathcal{X}} \int_{\mathcal{Y}} \phi(x, y) P_{XY}(x, y) dy dx \\
&= \int_{\mathcal{X}} \int_{\mathcal{Y}} \phi(x, y) P_X(x) P_{Y|X}(y|x) dy dx \\
&= \int_{\mathcal{X}} P_X(x) \int_{\mathcal{Y}} \phi(x, y) P_{Y|X}(y|x) dy dx \\
&= \int_{\mathcal{X}} P_X(x) \mathbb{E}_{Y|X=x}[\phi(X, Y)] dx \\
&= \mathbb{E}_X[\mathbb{E}_{Y|X}[\phi(X, Y)]]
\end{aligned}$$

Applying the same technique, we have  $\mathbb{E}_{XY}[\phi(X, Y)] = \mathbb{E}_Y[\mathbb{E}_{X|Y}[\phi(X, Y)]]$ .  $\square$ .

**Remark :** Usually, the label space is discrete and finite, meaning  $\mathcal{Y} = \{0, 1, 2, \dots, m\}$  for some  $m < \infty$ . Hence, the expectations over  $Y$  can be written as discrete sums:

$$\begin{aligned}
\mathbb{E}_{XY}[\phi(X, Y)] &= \mathbb{E}_Y[\mathbb{E}_{X|Y}[\phi(X, Y)]] = \sum_{y \in \mathcal{Y}} \mathbb{E}_{X|Y=y}[\phi(X, Y)] \\
&= \mathbb{E}_X[\mathbb{E}_{Y|X}[\phi(X, Y)]] = \mathbb{E}_X\left[\sum_{y \in \mathcal{Y}} \mathbb{E}_{Y=y|X}[\phi(X, Y)]\right]
\end{aligned}$$

**Definition 1.3** (Hypothesis space ( $\mathcal{H}$ )). 

---

The hypothesis space is a collection (family) of classifiers  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that have some common properties:

$$\mathcal{H} = \left\{ h : \mathcal{X} \rightarrow \mathcal{Y} \mid \text{some common properties} \right\}$$

For example, let  $\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = (0, 1)$ . In logistic regression, we assume the classifiers to be logit functions:

$$\mathcal{H}_{\text{logit}} = \left\{ h : \mathbb{R}^d \rightarrow (0, 1) \mid h(x) = \text{logit}(\beta x) = \frac{1}{1 + e^{-\beta x}}, \beta \in \mathbb{R}^{1 \times d} \right\}$$

**Definition 1.4** (Learning algorithm ( $\mathcal{L}_n$ )). 

---

To learn a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , suppose that we have access to a training dataset of  $n$  data pairs  $\{(X_k, Y_k)\}_{k=1}^n$  which are assumed to be **i.i.d sampled from**  $P_{XY}$ . The domain of the training data is then  $(\mathcal{X} \times \mathcal{Y})^n$ . A **learning algorithm**, denoted as  $\mathcal{L}_n$  is a function/procedure that derives a classifier  $\hat{h}_n : \mathcal{X} \rightarrow \mathcal{Y}$  from the training data.

$$\begin{aligned}
\mathcal{L}_n &: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H} \\
\hat{h}_n &= \mathcal{L}_n((X_1, Y_1), \dots, (X_n, Y_n))
\end{aligned}$$

## 1.2 Goal of classification

**Definition 1.5** (Risk ( $R(h)$ )).

The **risk** of a classifier is defined as followed:

$$R(h) = P(h(X) \neq Y) = \mathbb{E}[\mathbf{1}_{\{h(X) \neq Y\}}]$$

Where  $(X, Y)$  are independent of the training data.

**Definition 1.6** (Bayes Risk ( $R^*$ )).

The **Bayes risk** is the infimum of the risk taken over all  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , not just for  $h \in \mathcal{H}$ :

$$R^* = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y}} R(h)$$

**Definition 1.7** (Consistency of learning algorithms).

A learning algorithm  $\mathcal{L}_n$  is called:

- **Weakly consistent** if  $R(\hat{h}_n) \xrightarrow{P} R^*$ :

$$\lim_{n \rightarrow \infty} P(R(\hat{h}_n) \leq r) = P(R^* \leq r), \forall r \geq 0$$

- **Strongly consistent** if  $R(\hat{h}_n) \xrightarrow{a.s.} R^*$ :

$$P\left(\lim_{n \rightarrow \infty} |R(\hat{h}_n) - R^*| \geq \epsilon\right) = 0, \forall \epsilon > 0$$

- **Universally weakly/strongly consistent** if  $\mathcal{L}_n$  is weakly/strongly consistent for all  $P_{XY}$ .  
Meaning, consistency holds without any assumption about  $P_{XY}$ .

## 2 Bayes classifier

### 2.1 Properties of Bayes Risk

**Overview :** Recall that the Bayes classifier is the one with minimum risk and the corresponding risk is called the Bayes Risk. For  $\mathcal{Y} = \{0, 1\}$  and defined:

$$\eta(x) = P(Y = 1|X = x)$$

Define the following classifier:

$$h^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

#### Theorem 2.1: Properties of Bayes classifier

The following properties hold for the Bayes classifier with  $\mathcal{Y} = \{0, 1\}$  (Binary classification):

- (i)  $R(h^*) = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y}} \{R(h)\} = R^*$ .
- (ii)  $\underbrace{R(h) - R^*}_{\text{Excess risk}} = 2\mathbb{E}_X \left[ \left| \eta(x) - \frac{1}{2} \right| \mathbf{1}_{\{h(X) \neq h^*(X)\}} \right]$ .
- (iii)  $R^* = \mathbb{E} \left[ \min(\eta(X), 1 - \eta(X)) \right]$ .

**Proof** (Theorem 2.1). \_\_\_\_\_

*Proving each point:*

$$(i) \ R(h^*) = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y}} \{R(h)\} = R^*.$$

For all  $h: \mathcal{X} \rightarrow \mathcal{Y}$ , we have:

$$\begin{aligned} R(h) &= \mathbb{E}_{XY} [\mathbf{1}_{\{h(X) \neq Y\}}] \\ &= \mathbb{E}_{x \sim X} \left[ \mathbb{E}_{Y|X=x} [\mathbf{1}_{\{Y \neq h(x)\}}] \right] \\ &= \mathbb{E}_{x \sim X} \left[ \sum_{y \in \{0,1\}} \mathbf{1}_{\{y \neq h(x)\}} \right] \\ &= \mathbb{E}_{x \sim X} [\eta(x) \mathbf{1}_{\{h(x)=0\}} + (1 - \eta(x)) \mathbf{1}_{\{h(x)=1\}}] \end{aligned}$$

Since the two events  $\{h(x) = 1\}$  and  $\{h(x) = 0\}$  are mutually exclusive,  $R(h)$  is the smallest when we set  $h(x) = 1$  when  $\eta(x) \geq 1 - \eta(x) \implies \eta(x) \geq \frac{1}{2}$ . Therefore, we have:

$$h^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$(ii) \ \underbrace{R(h) - R^*}_{\text{Excess risk}} = 2\mathbb{E}_X \left[ \left| \eta(x) - \frac{1}{2} \right| \mathbf{1}_{\{h(X) \neq h^*(X)\}} \right].$$

We have:

$$\begin{aligned}
R(h) - R^* &= \mathbb{E}_{x \sim X} \left[ \mathbb{E}_{Y|X=x} \left[ \mathbf{1}_{\{Y \neq h(x)\}} \right] \right] - \mathbb{E}_{x \sim X} \left[ \mathbb{E}_{Y|X=x} \left[ \mathbf{1}_{\{Y \neq h^*(x)\}} \right] \right] \\
&= \mathbb{E}_{x \sim X} \left[ \sum_{y \in \{0,1\}} \mathbf{1}_{\{y \neq h(x)\}} P(Y = y|X = x) \right] - \mathbb{E}_{x \sim X} \left[ \sum_{y \in \{0,1\}} \mathbf{1}_{\{y \neq h^*(x)\}} P(Y = y|X = x) \right] \\
&= \mathbb{E}_{x \sim X} \left[ \eta(x) \left( \mathbf{1}_{\{h(x)=0\}} - \mathbf{1}_{\{h^*(x)=0\}} \right) + (1 - \eta(x)) \left( \mathbf{1}_{\{h(x)=1\}} - \mathbf{1}_{\{h^*(x)=1\}} \right) \right] \\
&= \mathbb{E}_{x \sim X} \left[ \eta(x) \left( \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=0\}} - \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=1\}} \right) \right. \\
&\quad \left. + (1 - \eta(x)) \left( \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=1\}} - \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=0\}} \right) \right] \\
&= \mathbb{E}_{x \sim X} \left[ (2\eta(x) - 1) \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=0\}} + (1 - 2\eta(x)) \mathbf{1}_{\{h(x) \neq h^*(x), h(x)=1\}} \right] \\
&= \mathbb{E}_{x \sim X} \left[ \left| 2\eta(x) - 1 \right| \mathbf{1}_{\{h(x) \neq h^*(x)\}} \right] \\
&= 2\mathbb{E}_X \left[ \left| \eta(X) - \frac{1}{2} \right| \mathbf{1}_{\{h(X) \neq h^*(X)\}} \right]
\end{aligned}$$

$$(iii) \ R^* = \mathbb{E} \left[ \min(\eta(X), 1 - \eta(x)) \right].$$

From (i) we have:

$$\begin{aligned}
R(h^*) &= \mathbb{E}_{x \sim X} \left[ \eta(x) \mathbf{1}_{\{h^*(x)=0\}} + (1 - \eta(x)) \mathbf{1}_{\{h^*(x)=1\}} \right] \\
&= \mathbb{E}_X \left[ \min(\eta(X), 1 - \eta(x)) \right]
\end{aligned}$$

□.

### Theorem 2.2: Properties of Bayes classifier (Multi-class)

For multi-class classification with more than two labels :  $\mathcal{Y} = \{1, 2, \dots, M\}$ , the Bayes classifier is defined as followed:

$$h^*(x) = \arg \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\}$$

$$\text{Where : } \eta_y(x) = P(Y = y | X = x)$$

The following properties hold for the Bayes classifier with  $\mathcal{Y} = \{1, 2, \dots, M\}$  (Multi-class classification):

- (i) **Bayes Risk**  $R^*$  :

$$R^* = \mathbb{E}_{x \sim X} \left[ 1 - \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\} \right] = \mathbb{E}_{x \sim X} \left[ \min_{y \in \mathcal{Y}} \overline{\eta}_y(x) \right]$$

- (ii) **Excess Risk**  $R(h) - R^*$  :

$$R(h) - R^* = \mathbb{E}_X \left[ \left( \eta_{y_x^*}(x) - \eta_{y_x}(x) \right) \mathbf{1}_{\{h(x) \neq h^*(x)\}} \right]$$

Where  $y_x = h(x)$  is the prediction made by an arbitrary classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$  and  $y_x^* = h^*(x)$  is the prediction made by the Bayes classifier.

**Proof** (Theorem 2.2).

(The proof of this theorem has been included in the solution of Exercise 2.1).

□.

## 2.2 Likelihood Ratio Test

**Overview** : Define  $\pi_1 = P(Y = 1)$  and  $\pi_0 = P(Y = 0)$  be the prior probabilities. Let  $p_1(x) = P(X = x | Y = 1)$  and  $p_0(x) = P(X = x | Y = 0)$  be the class-conditional densities. Note that we have:

$$\begin{aligned} \eta(x) &= P(Y = 1 | X = x) \\ &= \frac{P(X = x | Y = 1)P(Y = 1)}{P(X = x | Y = 1)P(Y = 1) + P(X = x | Y = 0)P(Y = 0)} \\ &= \frac{\pi_1 p_1(x)}{\pi_1 p_1(x) + \pi_0 p_0(x)} \\ &= \frac{1}{1 + \frac{\pi_0 p_0(x)}{\pi_1 p_1(x)}} \end{aligned}$$

Hence, we have:

$$\begin{aligned} \eta(x) \geq \frac{1}{2} &\iff \frac{\pi_0 p_0(x)}{\pi_1 p_1(x)} \\ &\iff \frac{p_1(x)}{p_0(x)} \geq \frac{\pi_0}{\pi_1} \end{aligned}$$



### Proposition 2.1: Likelihood ratio test

The Bayes classifier  $h^*$  can be re-defined as followed:

$$h^*(x) = \begin{cases} 1 & \text{if } \frac{p_1(x)}{p_0(x)} \geq \frac{\pi_0}{\pi_1} \\ 0 & \text{otherwise} \end{cases}$$

The fraction  $\frac{p_1(x)}{p_0(x)}$  is called the **likelihood ratio**.

## 2.3 Plug-in classifier

**Definition 2.1** (Plug-in classifier).

A **plug-in classifier** is based on an estimate of  $\eta(x)$ . This estimate is then plugged into the definition of the Bayes classifier. Suppose that  $\widehat{\eta}_n$  is an estimate of  $\eta$  based on  $n$  training samples  $\{(X_i, Y_i)\}_{i=1}^n$ . We define  $\widehat{h}_n$  as:

$$\widehat{h}_n = \begin{cases} 1 & \text{if } \widehat{\eta}_n(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

### Corollary 2.1: Excess risk of plug-in classifier

We have the following upper-bound for the excess risk of the plug-in classifier:

$$R(\widehat{h}_n) - R^* \leq 2\mathbb{E}_X \left[ \left| \eta(X) - \widehat{\eta}_n(X) \right| \right]$$

**Proof** (Corollary 2.1).

From theorem 2.1, we have:

$$R(\widehat{h}_n) - R^* = 2\mathbb{E}_X \left[ \left| \eta(X) - \frac{1}{2} \right| \mathbf{1}_{\{\widehat{h}_n(X) \neq h^*(X)\}} \right]$$

The indicator term will be non-zero in the above equality if one of the following cases occurs:

$$\begin{cases} \widehat{h}_n(X) = 1, h^*(X) = 0 \\ \widehat{h}_n(X) = 0, h^*(X) = 1 \end{cases} \implies \begin{cases} \widehat{\eta}_n(X) \geq \frac{1}{2}, \eta(X) < \frac{1}{2} \\ \widehat{\eta}_n(X) < \frac{1}{2}, \eta(X) \geq \frac{1}{2} \end{cases}$$

**Case 1 :**  $\widehat{\eta}_n(X) \geq \frac{1}{2}, \eta(X) < \frac{1}{2}$

We have:

$$\begin{aligned} \eta(X) - \widehat{\eta}_n(X) &\leq \eta(X) - \frac{1}{2} \quad (\text{Both sides negative}) \\ \implies \left| \eta(X) - \widehat{\eta}_n(X) \right| &\geq \left| \eta(X) - \frac{1}{2} \right| \end{aligned}$$

**Case 2 :**  $\widehat{\eta}_n(X) < \frac{1}{2}, \eta(X) \geq \frac{1}{2}$

We have:

$$\widehat{\eta}_n(X) - \eta(X) \geq \widehat{\eta}_n(X) - \frac{1}{2} \geq \eta(X) - \frac{1}{2} \quad (\text{All positive})$$

Therefore, we have:

$$\left| \eta(X) - \widehat{\eta}_n(X) \right| \geq \left| \eta(X) - \frac{1}{2} \right|$$

For both cases, we have the same  $\left| \eta(X) - \widehat{\eta}_n(X) \right| \geq \left| \eta(X) - \frac{1}{2} \right|$  inequality. Therefore, we have:

$$R(\widehat{h}_n) - R^* \leq 2\mathbb{E}_X \left[ \left| \eta(X) - \widehat{\eta}_n(X) \right| \right]$$

□.

## 2.4 End of chapter exercises

### Exercise 2.1

Extend theorem 2.1 to the multi-class classification case where  $\mathcal{Y} = \{1, 2, \dots, M\}$ . In other words, prove theorem 2.2.

**Solution** (Exercise 2.1).

We re-define the Bayes classifier  $h^*$  as followed:

$$h^*(x) = \arg \max_{y \in \mathcal{Y}} \{\eta_y(x)\},$$

$$\eta_y(x) = P(Y = y | X = x)$$

We have:

$$\sum_{y \in \mathcal{Y}} \eta_y(x) = 1, \forall x \in \mathcal{X}$$

(i) **Calculate Bayes risk  $R^*$**

For any classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , we have:

$$R(h) = \mathbb{E}_{x \sim X} \left[ \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) \right]$$

Letting  $\hat{y}_x = h(x)$  being  $h$ 's prediction for a given feature vector  $x \in \mathcal{X}$ , we have:

$$R(h) = \mathbb{E}_{x \sim X} \left[ \sum_{y \in \mathcal{Y}; y \neq \hat{y}_x} \eta_y(x) \right] = \mathbb{E}_{x \sim X} \left[ 1 - \eta_{\hat{y}_x}(x) \right]$$

In order to minimize  $R(h)$ , we need  $\eta_{\hat{y}_x}(x)$  to be maximized for all  $x \in \mathcal{X}$ . Hence, we have:

$$R^* = \mathbb{E}_{x \sim X} \left[ 1 - \max_{y \in \mathcal{Y}} \{\eta_y(x)\} \right]$$

Therefore, we have  $h^*(x) = \arg \max_{y \in \mathcal{Y}} \{\eta_y(x)\}$  is the Bayes classifier and the Bayes risk  $R^* =$

$$\mathbb{E}_{x \sim X} \left[ 1 - \max_{y \in \mathcal{Y}} \{\eta_y(x)\} \right].$$

(ii) **Calculate excess risk  $R(h) - R^*$**

For any  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , we have:

$$\begin{aligned} R(h) - R^* &= \mathbb{E}_{x \sim X} \left[ \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) \right] - \mathbb{E}_{x \sim X} \left[ 1 - \max_{y \in \mathcal{Y}} \{\eta_y(x)\} \right] \\ &= \mathbb{E}_{x \sim X} \left[ \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \{\eta_y(x)\} - 1 \right] \end{aligned}$$

Denote  $h^*(x) = y_x^*$  and  $h(x) = y_x$ . When  $h(x) = h^*(x) = y_x^*$ , we have:

$$\begin{aligned}
\sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \{\eta_y(x)\} &= \sum_{y \in \mathcal{Y}; y \neq y_x} \eta_y(x) + \eta_{y_x^*}(x) \\
&= \sum_{y \in \mathcal{Y}; y \neq y_x^*} \eta_y(x) + \eta_{y_x^*}(x) \\
&= \sum_{y \in \mathcal{Y}} \eta_y(x) = 1 \\
\Rightarrow \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \{\eta_y(x)\} - 1 &= 0
\end{aligned}$$

When  $h(x) \neq h^*(x)$ , we have:

$$\begin{aligned}
\sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \{\eta_y(x)\} - 1 &= \sum_{y \in \mathcal{Y}; y \neq y_x} \eta_y(x) + \eta_{y_x^*}(x) - 1 \\
&= 2\eta_{y_x^*}(x) - 1 + \sum_{y \in \mathcal{Y} \setminus \{y_x, y_x^*\}} \eta_y(x) \\
&= 2\eta_{y_x^*}(x) - (\eta_{y_x}(x) + \eta_{y_x^*}(x)) \\
&= \eta_{y_x^*}(x) - \eta_{y_x}(x).
\end{aligned}$$

Therefore, we can re-write the excess risk by multiplying the entire integrand with the indicator function  $\mathbf{1}_{\{h(x) \neq h^*(x)\}}$  as followed:

$$R(h) - R^* = \mathbb{E}_{x \sim X} \left[ \left( \eta_{y_x^*}(x) - \eta_{y_x}(x) \right) \mathbf{1}_{\{h(x) \neq h^*(x)\}} \right]$$

### (iii) Simpler form of Bayes risk

From (i) we have:

$$R^* = \mathbb{E}_X \left[ 1 - \max_{y \in \mathcal{Y}} \{\eta_y(x)\} \right] = \mathbb{E}_X \left[ \min_{y \in \mathcal{Y}} \{\bar{\eta}_y(x)\} \right]$$

Where  $\bar{\eta}_y(x) = P(Y \neq y | X = x)$ .

□.

### Exercise 2.2

Define the  $\alpha$ -cost-sensitive risk of a classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$  as followed:

$$R_\alpha(h) = \mathbb{E}_{XY} \left[ (1 - \alpha) \mathbf{1}_{\{Y=1, h(X)=0\}} + \alpha \mathbf{1}_{\{Y=0, h(X)=1\}} \right]$$

Define the Bayes classifier and prove an analogue of theorem 2.1.

**Solution** (Exercise 2.2).

Using the law of total expectation, we have:

$$\begin{aligned}
R_\alpha(h) &= \mathbb{E}_{x \sim X} \left[ \sum_{y \in \{0,1\}} \left[ (1 - \alpha) \mathbf{1}_{\{y=1, h(x)=0\}} + \alpha \mathbf{1}_{\{y=0, h(x)=1\}} \right] P(Y = y | X = x) \right] \\
&= \mathbb{E}_{x \sim X} \left[ (1 - \alpha) \eta(x) \mathbf{1}_{\{h(x)=0\}} + \alpha (1 - \eta(x)) \mathbf{1}_{\{h(x)=1\}} \right]
\end{aligned}$$

Since  $\mathbf{1}_{\{h(x)=0\}}$  and  $\mathbf{1}_{\{h(x)=1\}}$  are mutually exclusive, in order for  $R_\alpha(h)$  to be minimize, we define the following Bayes classifier:

$$h^*(x) = \begin{cases} 1 & \text{if } \alpha(1 - \eta(x)) \leq (1 - \alpha)\eta(x) \\ 0 & \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } \eta(x) \geq \alpha \\ 0 & \text{otherwise} \end{cases}$$

We can also derive a likelihood-ratio test version of the Bayes classifier, we have:

$$\begin{aligned} \eta(x) \geq \alpha &\implies \frac{1}{1 + \frac{\pi_0 p_0(x)}{\pi_1 p_1(x)}} \geq \alpha \\ &\implies 1 + \frac{\pi_0 \cdot p_0(x)}{\pi_1 \cdot p_1(x)} \leq \frac{1}{\alpha} \\ &\implies \frac{p_1(x)}{p_0(x)} \geq \frac{\alpha}{1 - \alpha} \cdot \frac{\pi_0}{\pi_1} \end{aligned}$$

Hence, we can rewrite the Bayes classifier as followed:

$$h^*(x) = \begin{cases} 1 & \text{if } \frac{p_1(x)}{p_0(x)} \geq \frac{\alpha}{1 - \alpha} \cdot \frac{\pi_0}{\pi_1} \\ 0 & \text{otherwise} \end{cases}$$

(i) **Bayes Risk**  $R_\alpha^*$

We have:

$$\begin{aligned} R_\alpha^* &= R_\alpha(h^*) \\ &= \mathbb{E}_{x \sim X} \left[ (1 - \alpha)\eta(x)\mathbf{1}_{\{h^*(x)=0\}} + \alpha(1 - \eta(x))\mathbf{1}_{\{h^*(x)=1\}} \right] \\ &= \mathbb{E}_X \left[ \min(\alpha(1 - \eta(X)), (1 - \alpha)\eta(X)) \right] \end{aligned}$$

(ii) **Excess Risk**  $R_\alpha(h) - R_\alpha^*$

For an arbitrary  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , we have:

$$\begin{aligned} R_\alpha(h) - R_\alpha^* &= \mathbb{E}_{x \sim X} \left[ (1 - \alpha)\eta(x) \left( \mathbf{1}_{\{h(x)=0\}} - \mathbf{1}_{\{h^*(x)=0\}} \right) + \alpha(1 - \eta(x)) \left( \mathbf{1}_{\{h(x)=1\}} - \mathbf{1}_{\{h^*(x)=1\}} \right) \right] \\ &= \mathbb{E}_{x \sim X} \left[ (1 - \alpha)\eta(x) \left( \mathbf{1}_{\{h(x)=0, h^*(x)=1\}} - \mathbf{1}_{\{h(x)=1, h^*(x)=0\}} \right) \right. \\ &\quad \left. + \alpha(1 - \eta(x)) \left( \mathbf{1}_{\{h(x)=1, h^*(x)=0\}} - \mathbf{1}_{\{h(x)=0, h^*(x)=1\}} \right) \right] \\ &= \mathbb{E}_{x \sim X} \left[ \mathbf{1}_{\{h(x)=0, h^*(x)=1\}} (\eta(x) - \alpha) + \mathbf{1}_{\{h(x)=1, h^*(x)=0\}} (\alpha - \eta(x)) \right] \\ &= \mathbb{E}_X \left[ \left| \eta(X) - \alpha \right| \mathbf{1}_{\{h(X) \neq h^*(X)\}} \right] \end{aligned}$$

□.

### 3 Hoeffding's inequality

#### 3.1 Markov's Inequality

##### Proposition 3.1: Markov's Inequality

Let  $U$  be a non-negative random variable on  $\mathbb{R}$ , then for all  $t > 0$ , we have:

$$P(U \geq t) \leq \frac{1}{t} \mathbb{E}[U]$$

**Proof** (Proposition 3.1). \_\_\_\_\_

We have:

$$\begin{aligned} tP(U \geq t) &= t\mathbb{E}[\mathbf{1}_{\{U \geq t\}}] \\ &= t \int_0^\infty \mathbf{1}_{\{x \geq t\}} f_U(x) dx \\ &= t \int_t^\infty f_U(x) dx \\ &\leq \int_t^\infty x f_U(x) dx \\ &\leq \int_0^\infty x f_U(x) dx = \mathbb{E}[U] \\ \implies P(U \geq t) &\leq \frac{1}{t} \mathbb{E}[U] \end{aligned}$$

□.

##### Corollary 3.1: Chebyshev's Inequality

Let  $Z$  be a random variable on  $\mathbb{R}$  with mean  $\mu$  and variance  $\sigma^2$ , we have:

$$P(|Z - \mu| \geq t) \leq \frac{\sigma^2}{t^2}$$

**Proof** (Corollary 3.1). \_\_\_\_\_

Using Markov's inequality, we have:

$$\begin{aligned} P(|Z - \mu| \geq t) &= P(|Z - \mu|^2 \geq t^2) \\ &\leq \frac{\mathbb{E}[|Z - \mu|^2]}{t^2} = \frac{\sigma^2}{t^2} \end{aligned}$$

□.

##### Corollary 3.2: Chernoff's bounding method

Let  $Z$  be a random variable on  $\mathbb{E}$ , for any  $t > 0$ , we have:

$$P(Z \geq t) \leq \inf_{s>0} e^{-st} M_Z(s)$$

**Proof** (Corollary 3.2). \_\_\_\_\_

We have:

$$\begin{aligned}
P(Z \geq t) &= P(sZ \geq st), \quad (t > 0) \\
&= P(e^{sZ} \geq e^{st}) \\
&\leq \frac{\mathbb{E}[e^{sZ}]}{e^{st}} = e^{-st} M_Z(s) \quad (\text{Markov's inequality})
\end{aligned}$$

Since the above inequality holds for all  $s > 0$ , we can just take the infimum to obtain the tightest bound. Hence, we have:

$$P(Z \geq t) \leq \inf_{s>0} e^{-st} M_Z(s)$$

□.

### 3.2 Hoeffding's Inequality

Before diving into Hoeffding's inequality, we need to go through the following lemma (whose proof will not be included) that will help us prove the Hoeffding's inequality:

#### Lemma 3.1: Hoeffding's lemma

Let  $V$  be a random variable on  $\mathbb{R}$  with  $\mathbb{E}[V] = 0$  and suppose that  $a \leq V \leq b$  with probability one. We have:

$$\mathbb{E}[e^{sV}] \leq \exp\left(\frac{s^2(b-a)^2}{8}\right)$$

**Proof** (Lemma 3.1). \_\_\_\_\_  
*(The proof for this lemma can be found here [6]).*

□.

#### Theorem 3.1: Hoeffding's Inequality

Let  $Z_1, Z_2, \dots, Z_n$  be independent random variables on  $\mathbb{R}$  such that  $a_i \leq Z_i \leq b_i$  with probability one for all  $1 \leq i \leq n$ . Let  $S_n = \sum_{i=1}^n Z_i$ . We have:

$$P\left(|S_n - \mathbb{E}[S_n]| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right), \quad \forall t > 0$$

**Proof** (Theorem 3.1). \_\_\_\_\_

Using the Chernoff's bounds, we have:

$$\begin{aligned}
P\left(\left|S_n - \mathbb{E}[S_n]\right| \geq t\right) &\leq \inf_{s>0} e^{-st} M_{S_n - \mathbb{E}[S_n]}(s) \\
&= \inf_{s>0} e^{-st} \mathbb{E}\left[e^{s(S_n - \mathbb{E}[S_n])}\right] \\
&= \inf_{s>0} e^{-st} \mathbb{E}\left[\exp\left(s \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i])\right)\right] \\
&= \inf_{s>0} e^{-st} \mathbb{E}\left[\prod_{i=1}^n \exp\left(s(Z_i - \mathbb{E}[Z_i])\right)\right] \\
&= \inf_{s>0} e^{-st} \prod_{i=1}^n \mathbb{E}\left[\exp\left(s(Z_i - \mathbb{E}[Z_i])\right)\right] \quad (\text{Since all } Z_i - \mathbb{E}[Z_i] \text{ are independent}) \\
&\leq \inf_{s>0} e^{-st} \prod_{i=1}^n \exp\left(\frac{s^2(b_i - a_i)^2}{8}\right) \quad (\text{By Hoeffding's lemma}) \\
&= \inf_{s>0} \exp\left(-st + \sum_{i=1}^n \frac{s^2(b_i - a_i)^2}{8}\right)
\end{aligned}$$

In order for the above to be minimized, we differentiate the term inside the exponential and set the derivative to 0 to find the optimal  $s > 0$ . We have:

$$-t + s \sum_{i=1}^n \frac{(b_i - a_i)^2}{4} = 0 \implies s = \frac{4t}{\sum_{i=1}^n (b_i - a_i)^2}$$

Letting  $c = \sum_{i=1}^n (b_i - a_i)^2$ , we now can derive the tightest Chernoff's bound as followed:

$$\begin{aligned}
P\left(\left|S_n - \mathbb{E}[S_n]\right| \geq t\right) &\leq \exp\left(-\frac{4t^2}{c} + \frac{16t^2}{c^2} \cdot \frac{c}{8}\right) = \exp\left(-\frac{2t^2}{c}\right) \\
&= \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)
\end{aligned}$$

□.

### 3.3 Convergence of Empirical Risk

**Definition 3.1** (Empirical Risk  $(\widehat{R}_n)$ ).

---

Suppose we are given training data  $\{(X_i, Y_i)_{i=1}^n\}$  such that each pair  $(X_i, Y_i) \sim P_{XY}$  are independently identically distributed. Let  $h : \mathcal{X} \rightarrow \mathcal{Y}$  be a classifier. We define the **empirical risk** to be:

$$\widehat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}}$$

Note that  $\mathbb{E}[\widehat{R}_n(h)] = R(h)$  and  $n\widehat{R}_n(h) \sim \text{Binomial}(n, R(h))$ . In the following corollary of the Hoeffding's inequality, we will answer the question **how close the empirical risk is as an estimate of true risk** or **how fast the empirical risk converges to the true risk**.



### Corollary 3.3: Convergence of Empirical Risk

Given training data  $\{(X_i, Y_i)_{i=1}^n\}$  such that each pair  $(X_i, Y_i) \sim P_{XY}$  are independently identically distributed. Let  $h : \mathcal{X} \rightarrow \mathcal{Y}$  be a classifier, we have:

$$P\left(\left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \leq 2e^{-2n\epsilon^2}, \quad \epsilon > 0$$

**Proof** (Corollary 3.3).

For all  $1 \leq i \leq n$ , we have  $\mathbf{1}_{\{h(X_i) \neq Y_i\}} \in \{0, 1\}$ . Hence, with probability one,  $0 \leq \mathbf{1}_{\{h(X_i) \neq Y_i\}} \leq 1$  and  $b_i = 1, a_i = 0$  for all  $1 \leq i \leq n$ .

Using the Hoeffding's inequality, we have:

$$\begin{aligned} P\left(\left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) &= P\left(\left|\widehat{R}_n(h) - \mathbb{E}[\widehat{R}_n(h)]\right| \geq \epsilon\right) \\ &= P\left(\left|n\widehat{R}_n(h) - \mathbb{E}[n\widehat{R}_n(h)]\right| \geq n\epsilon\right) \\ &\leq \exp\left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad (\text{Hoeffding's inequality}) \\ &= e^{-2n\epsilon^2} \end{aligned}$$

□.

## 3.4 KL-divergence & Hypothesis Testing

**Set-up (Hypothesis Testing)** : Suppose that we have  $\mathcal{Y} = \{0, 1\}$  and  $P_{XY}$  is a distribution on  $\mathcal{X} \times \mathcal{Y}$ . Let's assume that:

- The prior probabilities  $\pi_y$  are equal.
- The supports of likelihoods  $p_0, p_1$  are the same.
- $0 < \alpha \leq p_y(x) \leq \beta < \infty$  for all  $x \in \mathcal{X}$  such that  $p_y(x) > 0$  and for all  $y \in \{0, 1\}$ .

Now suppose  $X_1, \dots, X_n \sim p_y$  are independently identically distributed where  $y \in \{0, 1\}$  is unknown. Can we guess  $y$  and how good our guess would be?

### Proposition 3.2: KL-divergence hypothesis testing

From the above settings, the optimal classifier is given by the likelihood ratio test:

$$\widehat{h}_n(x) = \begin{cases} 1 & \text{if } \frac{\prod_{i=1}^n p_1(x_i)}{\prod_{i=1}^n p_0(x_i)} \geq \frac{\pi_0}{\pi_1} \quad (= 1) \\ 0 & \text{otherwise} \end{cases}$$

Where  $x = (x_1, \dots, x_n)$  is an observation of the random vector  $X = (X_1, \dots, X_n)$ . Define the class-specific risk  $R_y(h)$  be the risk of misclassification when the true label is  $Y = y$ :

$$R_y(h) = P(h(X) \neq Y | Y = y)$$

Then, we have:

$$R_0(\widehat{h}_n) \leq e^{-2nD(p_0||p_1)^2/c}, \text{ where } c = 4(\log \beta - \log \alpha)^2$$

Where  $D(p_0||p_1)$  is the *KL*-divergence of  $p_1$  from  $p_0$ . We can prove a similar exponentially decaying bound for  $R_1(\widehat{h}_n)$ .

#### Proof.

*Proposition 3.2* We can rewrite the optimal classifier as:

$$\widehat{h}_n(X) = \begin{cases} 1 & \text{if } \widehat{S}_n(X_1, \dots, X_n) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Where we have:

$$\begin{aligned} \widehat{S}_n(X_1, \dots, X_n) &= \log \frac{\prod_{i=1}^n p_1(X_i)}{\prod_{i=1}^n p_0(X_i)} \\ &= \sum_{i=1}^n \log \frac{p_1(X_i)}{p_0(X_i)} \\ &= \sum_{i=1}^n Z_i \quad \left( \text{Letting } Z_i = \log \frac{p_1(X_i)}{p_0(X_i)} \right) \end{aligned}$$

Since the likelihoods are bounded, we have:

$$a_i = \log \frac{\alpha}{\beta} \leq Z_i \leq \log \frac{\beta}{\alpha} = b_i, \quad 1 \leq i \leq n$$

Now, we have:

$$\begin{aligned} R_0(\widehat{h}_n) &= P(h(X) \neq Y | Y = 0) \\ &= P(\widehat{S}_n \geq 0 | Y = 0) \\ &= P(\widehat{S}_n - \mathbb{E}[S_n | Y = 0] \geq -\mathbb{E}[S_n | Y = 0] | Y = 0) \end{aligned}$$

To calculate the conditional expectation  $\mathbb{E}[S_n | Y = 0]$ , we have:

$$\begin{aligned} \mathbb{E}[S_n | Y = 0] &= n\mathbb{E}[Z_1 | Y = 0] \\ &= n \int \log \frac{p_1(x)}{p_0(x)} p_0(x) dx \\ &= -n \int \log \frac{p_0(x)}{p_1(x)} p_0(x) dx = -nD(p_0||p_1) \end{aligned}$$

Therefore, we have:

$$\begin{aligned} R_0(\widehat{h}_n) &= P(\widehat{S}_n - \mathbb{E}[S_n|Y=0] \geq nD(p_0||p_1)|Y=0) \\ &\leq \exp\left(-\frac{2n^2D(p_0||p_1)^2}{\sum_{i=1}^n(b_i - a_i)^2}\right) \quad (\text{Hoeffding's inequality}) \end{aligned}$$

For every  $1 \leq i \leq n$ , we have:

$$\begin{aligned} b_i - a_i &= \log \frac{\beta}{\alpha} - \log \frac{\alpha}{\beta} \\ &= \log \frac{\beta^2}{\alpha^2} = 2 \log \frac{\beta}{\alpha} = 2(\log \beta - \log \alpha) \\ \implies \sum_{i=1}^n (b_i - a_i)^2 &= 4n(\log \beta - \log \alpha)^2 \end{aligned}$$

Finally, we have:

$$R_0(\widehat{h}_n) \leq \exp\left(-\frac{2nD(p_0||p_1)^2}{4(\log \beta - \log \alpha)^2}\right)$$

Similarly, for  $R_1(\widehat{h}_n)$ , we have:

$$R_1(\widehat{h}_n) \leq \exp\left(-\frac{2nD(p_1||p_0)^2}{4(\log \beta - \log \alpha)^2}\right)$$

□.

### 3.5 End of chapter exercises

#### Exercise 3.1

- (i) Apply Chernoff's bounding method to obtain an exponential bound on the tail probability  $P(Z \geq t)$  for a Gaussian random variable  $Z \sim \mathcal{N}(\mu, \sigma^2)$ .
- (ii) Appealing to the central limit theorem, use part (i) to give an approximate bound on the binomial tail. This should not only match the exponential decay given by Hoeffding's inequality, but also reveal the dependence on the variance of the binomial.

**Solution** (Exercise 3.1). \_\_\_\_\_

(i) **Chernoff's bounds for  $Z \sim \mathcal{N}(\mu, \sigma^2)$**

Using the Chernoff's bounding method, we have:

$$\begin{aligned} P(Z \geq t) &\leq \inf_{s>0} e^{-st} M_Z(s) \\ &= \inf_{s>0} \exp \left( -st + \mu s + \frac{1}{2} \sigma^2 s^2 \right) \end{aligned}$$

The above bound is the tightest when the derivative of the term inside the exponential equals zero. Hence, we have:

$$-t + \mu + s\sigma^2 = 0 \implies s = \frac{t - \mu}{\sigma^2}$$

From the above, we have the tightest Chernoff's bound as followed:

$$P(Z \geq t) \leq \exp \left( -\frac{(t - \mu)^2}{\sigma^2} + \frac{(t - \mu)^2}{2\sigma^2} \right) = \exp \left( -\frac{(t - \mu)^2}{2\sigma^2} \right)$$

(ii) **Binomial tail upper bound**

Let  $S_n$  be the binomial random variable such that:

$$S_n = \sum_{i=1}^n X_i, \quad X_i \sim \text{Bernoulli}(p)$$

For a positive  $\epsilon > 0$ , we want to know the upper tail bound  $P(S_n - \mathbb{E}[S_n] \geq \epsilon)$ . Letting  $\bar{X} = \frac{1}{n} S_n$ , we have:

$$\begin{aligned} P(S_n - \mathbb{E}[S_n] \geq \epsilon) &= P \left( \bar{X} - \frac{\mathbb{E}[S_n]}{n} \geq \frac{\epsilon}{n} \right) \\ &= P \left( \bar{X} - p \geq \frac{\epsilon}{n} \right) \\ &= P \left( \frac{\bar{X} - p}{\sqrt{pq}/\sqrt{n}} \geq \frac{\epsilon}{\sqrt{npq}} \right), \quad (q = 1 - p) \end{aligned}$$

By the Central Limit Theorem, we have:

$$\frac{\bar{X} - p}{\sqrt{pq}/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1)$$

Hence, as  $n \rightarrow \infty$ , the upper tail bound would be:

$$\begin{aligned} P(S_n - \mathbb{E}[S_n] \geq \epsilon) &= P\left(\frac{\bar{X} - p}{\sqrt{pq}/\sqrt{n}} \geq \frac{\epsilon}{\sqrt{npq}}\right) \\ &\leq \exp\left(-\frac{\epsilon^2}{2npq}\right) = \exp\left(-\frac{\epsilon^2}{2\text{Var}(S_n)}\right) \end{aligned}$$

Double-check the bound with Hoeffding's inequality, we have:

$$P(S_n - \mathbb{E}[S_n] \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{n}\right)$$

□.

### Exercise 3.2

Can you remove the assumption in  $0 < \alpha \leq p_y(x)$ ? Consider other restrictions on  $p_y$ , other concentration inequalities, or other  $f$ -divergences.

**Solution** (Exercise 3.2).

When we remove the assumption that  $0 < \alpha \leq p_y(x)$ , the class-conditional densities are not bounded below. Hence, we have:

$$\exp\left(-\frac{2nD(p_1||p_0)^2}{4(\log \beta - \log \alpha)^2}\right) \rightarrow 1 \text{ when } \alpha \rightarrow 0$$

In other words, the bound is no longer meaningful. We can instead use the Chernoff bounding method:

$$\begin{aligned} R_0(\widehat{h}_n) &= P(S_n \geq 0 | Y = 0) \\ &\leq \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[ e^{sZ_i} \right] \\ &= \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[ \exp\left(s \log \frac{p_1(X_i)}{p_0(X_i)}\right) \right] \\ &= \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[ \frac{p_1(X_i)^s}{p_0(X_i)^s} \right] \end{aligned}$$

Taking logarithm from both sides, we have:

$$\begin{aligned} \log R_0(\widehat{h}_n) &\leq \inf_{s>0} \sum_{i=1}^n \log \mathbb{E}_{q_0} \left[ \frac{p_1(X_i)^s}{p_0(X_i)^s} \right] \\ &= \inf_{s>0} \sum_{i=1}^n (s-1) R_s(p_1||p_0) \\ &= \inf_{s>0} n(s-1) R_s(p_1||p_0) \end{aligned}$$

Where  $R_s(p_1||p_0)$  is the Renyi divergence [7].

□.

## 4 Empirical Risk Minimization

### 4.1 Uniform Deviation Bounds

**Definition 4.1** (Empirical Risk Minimization ( $\widehat{h}_n$ )).

Let  $\{(X_i, Y_i)\}_{i=1}^n$  be independently identically distributed random variables sampled from  $P_{XY}$ . Let  $\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}$  be a set of classifiers. **Empirical Risk Minimization** is a learning algorithm such that:

$$\widehat{h}_n = \arg \min_{h \in \mathcal{H}} \widehat{R}_n(h)$$

Where  $\widehat{R}_n$  is the empirical risk and  $\widehat{h}_n$  is called the **Empirical Risk Minimizer**. An important question is how close  $\widehat{R}_n$  is to  $R_{\mathcal{H}}^* = \inf_{h \in \mathcal{H}} R(h)$ .

**Overview (Uniform Deviation Bounds)** : Previously, we proved the following bound using the Hoeffding's inequality:

$$P\left(\left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \leq \delta$$

Where  $\delta = 2e^{-2n\epsilon^2}$ . However, since we do not know  $\widehat{h}_n$  (the specific function in  $\mathcal{H}$  that minimizes the empirical risk), we look for a bound that is guaranteed to apply for all  $h \in \mathcal{H}$ . This is called the Uniform Deviation Bound.

**Definition 4.2** (Uniform Deviation Bounds (UDB)).

Given a set of classifiers  $\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}$ ,  $\epsilon > 0$ , the **Uniform Deviation Bounds** is the probability that for at least one  $h \in \mathcal{H}$ , the empirical risk deviates away from the true risk by  $\epsilon$  and has the following form:

$$P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \leq \epsilon\right) \geq 1 - \delta$$
$$\text{Or : } P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \leq \delta$$

The above bounds have the following interpretations:

- The probability that the deviation from the true risk is at most  $\epsilon$  for all functions in  $\mathcal{H}$  is at least  $1 - \delta$ .
- The probability that there exists at least a function in  $\mathcal{H}$  whose deviation from the true risk is at least  $\epsilon$  is at most  $\delta$ .

Basically, we want to **bound the probability that some function deviates too far from the true risk**.

#### Theorem 4.1: Uniform Deviation Bounds for finite $\mathcal{H}$

Assume that  $|\mathcal{H}| < \infty$ . We have:

$$P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \leq 2|\mathcal{H}|e^{-2n\epsilon^2}$$

**Proof** (Theorem 4.1). 

---

For  $h \in \mathcal{H}$ , define the following event:

$$\Omega_\epsilon(h) = \left\{ \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon \right\}$$

Which is the event that the function  $h$  deviates away from the true risk by  $\epsilon > 0$ . Now, define the following event:

$$\Omega_\epsilon(\mathcal{H}) = \bigcup_{h \in \mathcal{H}} \Omega_\epsilon(h)$$

Which is the event that at least one  $h \in \mathcal{H}$  deviates away from the true risk by  $\epsilon > 0$ . We have:

$$\begin{aligned} P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon\right) &= P(\Omega_\epsilon(\mathcal{H})) \\ &= P\left(\bigcup_{h \in \mathcal{H}} \Omega_\epsilon(h)\right) \\ &\leq \sum_{h \in \mathcal{H}} P(\Omega_\epsilon(h)) \\ &\leq \sum_{h \in \mathcal{H}} 2e^{-2n\epsilon^2} = 2|\mathcal{H}|e^{-2n\epsilon^2} \end{aligned}$$

□.

**Proposition 4.1: (Probabilistic) Bound on Excess Risk of  $\widehat{h}_n$**

Suppose that  $\mathcal{H}$  satisfies:

$$P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon\right) \leq \delta$$

Then, with probability of at least  $1 - \delta$ , we have the following **upper bound on the Excess Risk of the Empirical Risk Minimizer**:

$$R(\widehat{h}_n) - R_{\mathcal{H}}^* \leq 2\epsilon$$

In other words, **with probability  $1 - \delta$ , the empirical risk minimizer deviates from the true risk minimizer by at most  $2\epsilon$ .**

**Proof** (Proposition 4.1). 

---

We have:

$$P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon\right) \leq \delta \implies P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \leq \epsilon\right) \geq 1 - \delta$$

Hence, with probability  $1 - \delta$ , for all  $h \in \mathcal{H}$ , we have:

$$\begin{aligned} \left| \widehat{R}_n(h) - R(h) \right| \leq \epsilon &\implies -\epsilon \leq \widehat{R}_n(h) - R(h) \leq \epsilon \\ &\implies \begin{cases} \widehat{R}_n(h) &\leq R(h) + \epsilon \\ R(h) &\leq \widehat{R}_n(h) + \epsilon \end{cases} \end{aligned}$$

Therefore:

$$\begin{aligned}
R(\widehat{h}_n) &\leq \widehat{R}_n(\widehat{h}_n) + \epsilon \\
&\leq \widehat{R}_n(h) + \epsilon \quad (\text{Since } \widehat{h}_n \text{ minimizes the Empirical Risk}) \\
&\leq (R(h) + \epsilon) + \epsilon = R(h) + 2\epsilon
\end{aligned}$$

Since  $h \in \mathcal{H}$  is an arbitrary choice, we take the infimum over  $\mathcal{H}$  to get the tightest bound. We have:

$$\begin{aligned}
R(\widehat{h}_n) &\leq \inf_{h \in \mathcal{H}} R(h) + 2\epsilon \\
&= R_{\mathcal{H}}^* + 2\epsilon
\end{aligned}$$

□.

**Remark :** We can express the above proposition verbally as "If the UDB is at most  $\delta$ , then with probability  $1 - \delta$ , the Excess Risk of the Empirical Risk Minimizer is at most  $2\epsilon$ ".

**Remark :** Note that the above proof assumes that *there exists an empirical risk minimizer*. This is not guaranteed when  $|\mathcal{H}|$  is infinite.

#### Proposition 4.2: (Non-probabilistic) Bound on Excess Risk of $\widehat{h}_n$

We have the following inequality:

$$R(\widehat{h}_n) - R_{\mathcal{H}}^* \leq 2 \sup_{h \in \mathcal{H}} |\widehat{R}_n(h) - R(h)|$$

**Proof** (Proposition 4.2).

Let  $h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} R(h)$ . We have:

$$R(\widehat{h}_n) - R_{\mathcal{H}}^* \leq |R(\widehat{h}_n) - \widehat{R}_n(\widehat{h}_n)| + \widehat{R}_n(\widehat{h}_n) - \widehat{R}_n(h_{\mathcal{H}}^*) + |\widehat{R}_n(h_{\mathcal{H}}^*) - R_{\mathcal{H}}^*|$$

Since  $\widehat{h}_n$  is the Empirical Risk Minimizer, we have  $\widehat{R}_n(\widehat{h}_n) - \widehat{R}_n(h_{\mathcal{H}}^*) \leq 0$ . Hence:

$$\begin{aligned}
R(\widehat{h}_n) - R_{\mathcal{H}}^* &\leq |R(\widehat{h}_n) - \widehat{R}_n(\widehat{h}_n)| + |\widehat{R}_n(h_{\mathcal{H}}^*) - R_{\mathcal{H}}^*| \\
&\leq 2 \sup_{h \in \mathcal{H}} |\widehat{R}_n(h) - R(h)|
\end{aligned}$$

□.

#### Corollary 4.1: Excess Risk of $\widehat{h}_n$ - $\delta \rightarrow \epsilon$ relation

This is a Corollary for both proposition 4.1 and proposition 4.2. If  $\mathcal{H}$  is finite, then:

$$P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) \leq \underbrace{2|\mathcal{H}|e^{-n\epsilon^2/2}}_{\delta}$$

Equivalently, with probability of at least  $1 - \delta$ , we have:

$$R(\widehat{h}_n) \leq R_{\mathcal{H}}^* + \sqrt{\frac{2}{n} \left( \log |\mathcal{H}| - \log \frac{\delta}{2} \right)}$$



**Proof** (Corollary 4.1).

By proposition 4.2, we have:

$$\begin{aligned} P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) &\leq P\left(2 \sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \\ &= P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \frac{\epsilon}{2}\right) \\ &\leq 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2}\right) \end{aligned}$$

Now, let:

$$\delta = 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2}\right) \implies \epsilon = \sqrt{\frac{2}{n} \left(\log |\mathcal{H}| - \log \frac{\delta}{2}\right)}$$

By proposition 4.1, with at least probability  $1 - \delta$ , we have:

$$R(\widehat{h}_n) \leq R_{\mathcal{H}}^* + \epsilon = R_{\mathcal{H}}^* + \sqrt{\frac{2}{n} \left(\log |\mathcal{H}| - \log \frac{\delta}{2}\right)}$$

□.

## 4.2 PAC Learning & Sample Complexity

**Definition 4.3** (PAC & Sample Complexity  $(N(\epsilon, \delta))$ ).

We say that  $\widehat{h}_n$  is a  $(\epsilon, \delta)$ -**learning algorithm** for  $\mathcal{H}$  if there exists a function  $N(\epsilon, \delta)$  such that:

$$\forall \epsilon, \delta > 0 : n \geq N(\epsilon, \delta) \implies P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) \leq \delta$$

Where we have:

- $N(\epsilon, \delta)$  is called the **Sample Complexity**.
- $\mathcal{H}$  is called **Uniformly Learnable**.
- $\widehat{h}_n$  is called **Probably Approximately Correct (PAC)**.

**Remark :** By corollary 4.1, we have  $\delta = 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2}\right)$ . Solving for  $n$ , we have:

$$N(\epsilon, \delta) = \frac{2}{\epsilon^2} \left(\log |\mathcal{H}| - \log \frac{\delta}{2}\right)$$

## 4.3 Zero-error case

In the following proposition, we can obtain a tighter bound for the zero empirical risk case. However, it is not particularly useful in many cases.

**Proposition 4.3: Zero-error case bound**

If  $\widehat{R}_n(\widehat{h}_n) = 0$  and  $|\mathcal{H}| < \infty$ , we have:

$$P\left(\exists h \in \mathcal{H} : \widehat{R}_n(h) = 0, R(h) \geq \epsilon\right) \leq \underbrace{|\mathcal{H}|e^{-n\epsilon}}_{\delta}$$

Meaning, with probability of at least  $1 - \delta$ , if  $\widehat{R}_n(h) = 0$  then  $R(h) \leq \frac{1}{n}(\log |\mathcal{H}| - \log \delta)$ .

**Proof** (Proposition 4.3).

Let  $\Omega_0(h) = \{\widehat{R}_n(h) = 0\}$  and define the event  $\Omega_\epsilon$  as:

$$\Omega_\epsilon = \bigcup_{h \in \mathcal{H}; R(h) \geq \epsilon} \Omega_0(h) = \left\{ \exists h \in \mathcal{H} : \widehat{R}_n(h) = 0, R(h) \geq \epsilon \right\}$$

For any  $h \in \mathcal{H}$  such that  $R(h) \geq \epsilon$ , we have:

$$\begin{aligned} P(\Omega_0(h)) &= P\left(\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}} = 0\right) \\ &= P\left(\sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}} = 0\right) \\ &= P\left(\bigcup_{i=1}^n \{h(X_i) = Y_i\}\right) \\ &= \prod_{i=1}^n P(h(X_i) = Y_i) \quad (\text{Since all } (X_i, Y_i) \text{ pairs are independent}) \end{aligned}$$

Each  $\mathbf{1}_{\{h(X_i) \neq Y_i\}}$  is a Bernoulli variable with hit probability  $p_i = 1 - \mathbb{E}[h(X_i) \neq Y_i] = 1 - R(h)$ . Hence, we have:

$$\begin{aligned} P(\Omega_0(h)) &= \prod_{i=1}^n P(h(X_i) = Y_i) \\ &= (1 - R(h))^n \\ &\leq (1 - \epsilon)^n \end{aligned}$$

Using the inequality  $\log(1 - \epsilon) \leq -\epsilon$ , we have:

$$\begin{aligned} P(\Omega_0(h)) &\leq (1 - \epsilon)^n = e^{n \log(1 - \epsilon)} \\ &\leq e^{-n\epsilon} \end{aligned}$$

Finally, we have:

$$\begin{aligned} P(\Omega_\epsilon) &= P\left(\bigcup_{h \in \mathcal{H}; R(h) \geq \epsilon} \Omega_0(h)\right) \\ &\leq \sum_{h \in \mathcal{H}; R(h) \geq \epsilon} P(\Omega_0(h)) \\ &\leq \sum_{h \in \mathcal{H}; R(h) \geq \epsilon} e^{-n\epsilon} \\ &\leq |\mathcal{H}|e^{-n\epsilon} \end{aligned}$$

□.

**Remark :** Note that the bound obtained in proposition 4.3 is NOT the Uniform Deviation Bound (UDB) because we have:

$$\left\{ \sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon \right\} = \left\{ \exists h \in \mathcal{H} : \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon \right\}$$

Therefore, we have:

$$\left\{ \exists h \in \mathcal{H} : \widehat{R}_n(h) = 0, R(h) \geq \epsilon \right\} \subseteq \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon \right\}$$

**Remark :** This is trivial improvement. However, define the following subset of  $\mathcal{H}$ :

$$H_\epsilon^+ = \left\{ h \in \mathcal{H} : R(h) \geq \epsilon \right\}$$

We can improve the bound in proposition 4.3 as followed:

$$P(\Omega_\epsilon) \leq |H_\epsilon^+| e^{-n\epsilon}$$

#### 4.4 End of chapter exercises

##### Exercise 4.1: Neyman-Pearson Criterion

The probability of error is not the only performance measure for binary classification. Indeed, the probability of error depends on the prior probability of the class label  $Y$ , and it may be that the frequency of the classes changes from training to testing data. In such cases, it is desirable to have a performance measure that does not require knowledge of the prior class probability. Let  $P_y$  be the class conditional distribution of class  $y \in \{0, 1\}$ . Define  $R_y(h) = P_y(h(X) \neq y)$ . Also let  $\alpha \in (0, 1)$ . For  $\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}$ , define:

$$R_{\mathcal{H},1}^* = \inf_{h \in \mathcal{H}} R_1(h) \\ \text{s.t. } R_0(h) \leq \alpha$$

In this problem you will investigate a discrimination rule that is probably approximately correct with respect to the above criterion, which is sometimes called the Neyman-Pearson criterion based on connections to the Neyman-Pearson lemma in hypothesis testing. Suppose we observe  $X_1^y, X_2^y, \dots, X_{n_y}^y \sim P_y$  for  $y \in \{0, 1\}$ . Define the empirical errors:

$$\widehat{R}_y(h) = \frac{1}{n_y} \sum_{i=1}^{n_y} \mathbf{1}_{\{h(X_i^y) \neq y\}}$$

Fix  $\epsilon > 0$  and consider the discrimination rule:

$$\widehat{h}_n = \arg \min_{h \in \mathcal{H}} \widehat{R}_1(h) \\ \text{s.t. } \widehat{R}_0(h) \leq \alpha + \frac{\epsilon}{2}$$

Suppose  $\mathcal{H}$  is finite. Show that with high probability:

$$R_0(\widehat{h}_n) \leq \alpha + \epsilon \text{ and } R_1(\widehat{h}_n) \leq R_{\mathcal{H},1}^* + \epsilon$$

**Solution** (Exercise 4.1).

We will prove each point one by one:

- (i)  $R_0(\widehat{h}_n) \leq \alpha + \epsilon$  **with high probability**.

**Claim 1 :**  $\forall y \in \{0, 1\}, \epsilon > 0 : P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_y(h) - R_y(h)\right| \geq \epsilon\right) \leq 2|\mathcal{H}|e^{-2n\epsilon^2}$

We have that  $n\widehat{R}_n(h) \sim \text{Binomial}(n, R_y(h))$  for all  $h \in \mathcal{H}$ . Hence, we have:

$$\begin{aligned} P\left(\left|\widehat{R}_y(h) - R_y(h)\right| \geq \epsilon\right) &= P\left(\left|n\widehat{R}_y(h) - nR_y(h)\right| \geq n\epsilon\right) \\ &= P\left(\left|n\widehat{R}_y(h) - \mathbb{E}[n\widehat{R}_y(h)]\right| \geq n\epsilon\right) \\ &\leq 2 \exp\left(-\frac{2n^2\epsilon^2}{n}\right) = 2e^{-2n\epsilon^2} \quad (\text{Hoeffding's inequality}) \end{aligned}$$

From the above, we have:

$$\begin{aligned}
P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_y(h) - R_y(h) \right| \geq \epsilon\right) &= P\left(\bigcup_{h \in \mathcal{H}} \left\{ \left| \widehat{R}_y(h) - R_y(h) \right| \geq \epsilon \right\}\right) \\
&\leq \sum_{h \in \mathcal{H}} P\left(\left| \widehat{R}_y(h) - R_y(h) \right| \geq \epsilon\right) \\
&\leq \sum_{h \in \mathcal{H}} 2e^{-2n\epsilon^2} = 2|\mathcal{H}|e^{-2n\epsilon^2}
\end{aligned}$$

From the assumption, we have:

$$\widehat{R}_0(\widehat{h}_n) \leq \alpha + \frac{\epsilon}{2}$$

Hence, we have:

$$\begin{aligned}
R_0(\widehat{h}_n) &= \widehat{R}_0(\widehat{h}_n) + R_0(\widehat{h}_n) - \widehat{R}_0(\widehat{h}_n) \\
&\leq \alpha + \frac{\epsilon}{2} + \left| R_0(\widehat{h}_n) - \widehat{R}_0(\widehat{h}_n) \right| \\
&\leq \alpha + \frac{\epsilon}{2} + \sup_{h \in \mathcal{H}} \left| R_0(h) - \widehat{R}_0(h) \right|
\end{aligned}$$

From **Claim 1**, we know that:

$$\begin{aligned}
P\left(\sup_{h \in \mathcal{H}} \left| R_0(h) - \widehat{R}_0(h) \right| \geq \frac{\epsilon}{2}\right) &\leq 2|\mathcal{H}|e^{-n\epsilon^2/2} \\
\implies P\left(\sup_{h \in \mathcal{H}} \left| R_0(h) - \widehat{R}_0(h) \right| \leq \frac{\epsilon}{2}\right) &\geq 1 - 2|\mathcal{H}|e^{-n\epsilon^2/2}
\end{aligned}$$

Hence, with probability of at least  $1 - 2|\mathcal{H}|e^{-n\epsilon^2/2}$ , we have:

$$R_0(\widehat{h}_n) \leq \alpha + \frac{\epsilon}{2} + \frac{\epsilon}{2} = \alpha + \epsilon$$

- (ii)  $R_1(\widehat{h}_n) \leq R_{\mathcal{H},1}^* + \epsilon$  **with high probability.**

**Claim 2 :**  $R_1(\widehat{h}_n) - R_{\mathcal{H},1}^* \leq 2 \sup_{h \in \mathcal{H}} \left| \widehat{R}_1(h) - R_1(h) \right|$

Let  $h' \in \mathcal{H}$  be any function such that  $\widehat{R}_0(h') \leq \alpha + \frac{\epsilon}{2}$ . We have:

$$\begin{aligned}
R_1(\widehat{h}_n) - R_{\mathcal{H},1}^* &= R_1(\widehat{h}_n) - \widehat{R}_1(\widehat{h}_n) + \widehat{R}_1(\widehat{h}_n) - \widehat{R}_1(h') + \widehat{R}_1(h') - R_{\mathcal{H},1}^* \\
&\leq \left| R_1(\widehat{h}_n) - \widehat{R}_1(\widehat{h}_n) \right| + \underbrace{\left| \widehat{R}_1(\widehat{h}_n) - \widehat{R}_1(h') \right|}_{\leq 0} + \left| \widehat{R}_1(h') - R_{\mathcal{H},1}^* \right| \\
&\leq \left| R_1(\widehat{h}_n) - \widehat{R}_1(\widehat{h}_n) \right| + \left| \widehat{R}_1(h') - R_{\mathcal{H},1}^* \right| \\
&\leq 2 \sup_{h \in \mathcal{H}} \left| \widehat{R}_1(h) - R_1(h) \right|
\end{aligned}$$

From **Claim 2**, we have:

$$\begin{aligned}
P\left(R_1(\widehat{h}_n) - R_{\mathcal{H},1}^* \geq \epsilon\right) &\leq P\left(2 \sup_{h \in \mathcal{H}} \left| \widehat{R}_1(h) - R_1(h) \right| \geq \epsilon\right) \\
&= P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_1(h) - R_1(h) \right| \geq \frac{\epsilon}{2}\right) \\
&\leq 2|\mathcal{H}|e^{-n\epsilon^2/2} \quad (\text{From **Claim 1**}) \\
\implies P\left(R_1(\widehat{h}_n) - R_{\mathcal{H},1}^* \leq \epsilon\right) &\geq 1 - 2|\mathcal{H}|e^{-n\epsilon^2/2}
\end{aligned}$$

Hence, with probability of at least  $1 - 2|\mathcal{H}|e^{-n\epsilon^2/2}$ , we have that  $R_1(\widehat{h_n}) \leq R_{\mathcal{H},1}^* + \epsilon$ .

□.

## 5 Vapnik-Chevronenkis Theory

In the following section, we will review a notion for measuring complexity of function class called **VC dimension**. Later on in this note, we will show that VC dimension is an upper-bound for the **Rademacher Complexity**.

### 5.1 VC Dimension

**Definition 5.1** (Restriction ( $N_{\mathcal{H}}$ )).

Let  $\mathcal{H} \in \{0,1\}^{\mathcal{X}}$  be a set of classifiers. The **restriction** of  $\mathcal{H}$  to a finite subset  $C \subset \mathcal{X}$  where  $|C| = n$  is the set of binary vectors defined by:

$$N_{\mathcal{H}}(C) = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H}, x_i \in C \right\}$$

Clearly, we have  $|N_{\mathcal{H}}(C)| \leq 2^n$  (cardinality of powerset of  $C$ ).

**Definition 5.2** (Shattering Coefficient ( $S_{\mathcal{H}}$ )).

The  $n^{\text{th}}$  **Shattering coefficient** (sometimes called the **Growth function**) is defined as:

$$S_{\mathcal{H}}(n) = \sup_{C \subset \mathcal{X}; |C|=n} |N_{\mathcal{H}}(C)|$$

Hence, we have:

$$|N_{\mathcal{H}}(C)| \leq S_{\mathcal{H}}(n) \leq 2^n, \forall C \subset \mathcal{X}$$

Intuitively, the  $n^{\text{th}}$  shattering coefficient is the size of the largest  $n$ -element restriction of  $\mathcal{H}$ . It measures the **richness** of  $\mathcal{H}$ .

If  $S_{\mathcal{H}}(n) = 2^n$ . Then  $\exists C \subset \mathcal{X}, |C| = n$  such that  $|N_{\mathcal{H}}(C)| = 2^n$ . We then say that  $\mathcal{H}$  **shatters the points in  $C$** .

**Definition 5.3** (VC-dimension ( $V_{\mathcal{H}}$ )).

The **VC dimension** of  $\mathcal{H}$  is defined as:

$$V_{\mathcal{H}} = \sup \left\{ n : S_{\mathcal{H}}(n) = 2^n \right\} = \sup_{C \subset \mathcal{X}} \left\{ |C| : \mathcal{H} \text{ shatters } C \right\}$$

If  $S_{\mathcal{H}}(n) = 2^n, \forall n \geq 1$  then  $V_{\mathcal{H}} = \infty$ .

**Remark :** Note that when  $|\mathcal{H}| < \infty$ , we have:

$$\begin{aligned} |N_{\mathcal{H}}(C)| &= \left| \{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H}, x_i \in C \} \right| \leq |\mathcal{H}| \\ \implies S_{\mathcal{H}}(n) &\leq |\mathcal{H}| \\ \implies V_{\mathcal{H}} &\leq \log_2 |\mathcal{H}| \end{aligned}$$

**Remark :** To show that  $V_{\mathcal{H}} = n$ , we must show that there exists at least  $n$  points  $x_1, \dots, x_n$  shattered by  $\mathcal{H}$  and no set of  $n+1$  points can be shattered by  $\mathcal{H}$ .

**Remark :** From the above definitions, we can understand  $N_{\mathcal{H}}$ ,  $S_{\mathcal{H}}$  and  $V_{\mathcal{H}}$  as followed:

- $N_{\mathcal{H}}(C)$  : Number of ways to assign labels to  $C \subset \mathcal{X}$  of size  $n \geq 1$ .
- $S_{\mathcal{H}}(n)$  : Maximum number of ways to assign labels to subsets of size  $n \geq 1$ .
- $V_{\mathcal{H}}$  : Maximum subset size  $n \geq 1$  such that we have  $2^n$  ways to assign labels (fully labelled).

## 5.2 Sauer's Lemma

### Theorem 5.1: Sauer's Lemma

This is a bound on the shatter coefficient. Let  $d = V_{\mathcal{H}} \leq \infty$ . For all  $n \geq 1$ , we have:

$$S_{\mathcal{H}}(n) \leq \sum_{k=0}^d \binom{n}{k}$$

**Proof** (Theorem 5.1, cited [3]).

Given a function class  $\mathcal{H}$  and a subset  $C \subset \mathcal{X}$ . For brevity, denote the restriction of  $\mathcal{H}$  to  $C$  as:

$$N_{\mathcal{H}}(C) = \mathcal{H}_C$$

To prove the above theorem, we prove a stronger result: For all subset  $C \subset \mathcal{X}$  where  $|C| = n$ , we have

$$|\mathcal{H}_C| \leq \left| \left\{ B \subset C : \mathcal{H} \text{ shatters } B \right\} \right| \leq \sum_{k=0}^d \binom{n}{k}$$

The second inequality holds because no set with size larger than  $d$  is shattered by  $\mathcal{H}$ . To prove that the first inequality holds for subsets of any size  $n \geq 1$ , we prove by induction:

- **Base case :** Let  $n = 1$ . Hence, we have  $C = \{x\}$  for  $x \in \mathcal{X}$ . Denote that:

$$\Phi_C = \left\{ B \subset C : \mathcal{H} \text{ shatters } B \right\}$$

We have:

$$\begin{cases} \mathcal{H} \text{ shatters } C & \implies \mathcal{H}_C = \{0, 1\}, \Phi_C = \{\emptyset, C\} \\ \mathcal{H} \text{ not shatter } C & \implies \mathcal{H}_C = \{0\} \text{ or } \{1\}, \Phi_C = \{\emptyset\} \end{cases}$$

For both cases, we have  $|\mathcal{H}_C| = |\Phi_C|$ .

- **Inductive case :** Assume that the first inequality holds for  $n = m - 1, m \geq 2$ . We have to prove that it holds for  $n = m$ . Let  $C = \{c_1, \dots, c_m\}$  and  $C' = \{c_2, \dots, c_m\}$ . Define the following label sets:

$$Y_0 = \{(y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in \mathcal{H}_C \vee (1, y_2, \dots, y_m) \in \mathcal{H}_C\}$$

$$Y_1 = \{(y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in \mathcal{H}_C \wedge (1, y_2, \dots, y_m) \in \mathcal{H}_C\}$$

First, we notice that  $Y_0 = \mathcal{H}_{C'}$ . Hence, we have:

$$\begin{aligned} |Y_0| &= |\mathcal{H}_{C'}| \\ &\leq \left| \left\{ B \subset C' : \mathcal{H} \text{ shatters } B \right\} \right| \\ &= \left| \left\{ B \subset C : c_1 \notin B, \mathcal{H} \text{ shatters } B \right\} \right| \end{aligned}$$



Next, we define the following sub-class of  $\mathcal{H}$ :

$$\mathcal{H}' = \left\{ h \in \mathcal{H} : \exists h' \in \mathcal{H} \text{ s.t. } h'(c) = \begin{cases} 1 - h(c) & \text{if } c = c_1 \\ h(c) & \text{otherwise} \end{cases} \right\}$$

Note that  $Y_1 = \mathcal{H}'_{C'}$ , and  $\mathcal{H}'$  shatters  $B \in C'$  implies  $\mathcal{H}'$  shatters  $B \cup \{c_1\}$  because for any  $h' \in \mathcal{H}'$ , there is always another function in  $\mathcal{H}'$  that gives the opposite label to  $c_1$ . Hence, we have:

$$\begin{aligned} |Y_1| &= |\mathcal{H}'_{C'}| \\ &\leq \left| \left\{ B \subset C' : \mathcal{H}' \text{ shatters } B \right\} \right| \\ &= \left| \left\{ B \subset C' : \mathcal{H}' \text{ shatters } B \cup \{c_1\} \right\} \right| \\ &= \left| \left\{ B \subset C : c_1 \in B, \mathcal{H}' \text{ shatters } B \right\} \right| \\ &\leq \left| \left\{ B \subset C : c_1 \in B, \mathcal{H} \text{ shatters } B \right\} \right| \end{aligned}$$

From the above, we have:

$$\begin{aligned} |\mathcal{H}_C| &= |Y_0| + |Y_1| \\ &\leq \left| \left\{ B \subset C : c_1 \notin B, \mathcal{H} \text{ shatters } B \right\} \right| + \left| \left\{ B \subset C : c_1 \in B, \mathcal{H} \text{ shatters } B \right\} \right| \\ &= \left| \left\{ B \subset C : \mathcal{H} \text{ shatters } B \right\} \right| \\ &\leq \sum_{k=0}^d \binom{n}{k} \end{aligned}$$

Taking the supremum over  $C \subset \mathcal{X}$  for both sides, we have:

$$S_{\mathcal{H}}(n) \leq \sum_{k=0}^d \binom{n}{k}$$

□.

#### Corollary 5.1: Sauer's lemma - bound on $S_{\mathcal{H}}(n)$ I

If  $d = V_{\mathcal{H}} < \infty$ , for all  $n \geq 1$ , we have:

$$S_{\mathcal{H}}(n) \leq (n+1)^d$$

**Proof** (Corollary 5.1). \_\_\_\_\_

By Binomial theorem, we have:

$$\begin{aligned} (n+1)^d &= \sum_{k=1}^d n^k \binom{d}{k} \\ &= \sum_{k=1}^d n^k \frac{d!}{k!(d-k)!} \\ &\geq \sum_{k=1}^d \frac{n^k}{k!} \geq \sum_{k=1}^d \frac{n!}{(n-k)!k!} = \sum_{k=1}^d \binom{n}{k} \geq S_{\mathcal{H}}(n) \end{aligned}$$

□.

**Corollary 5.2: Sauer's lemma - bound on  $S_{\mathcal{H}}(n)$  II**

For all  $n \geq d = V_{\mathcal{H}}$ , we have:

$$S_{\mathcal{H}}(n) \leq \left(\frac{ne}{d}\right)^d$$

**Proof** (Corollary 5.2). \_\_\_\_\_

For  $\frac{d}{n} < 1$ , we have:

$$\begin{aligned} \left(\frac{d}{n}\right)^d \sum_{k=0}^d \binom{n}{k} &\leq \sum_{k=0}^d \left(\frac{d}{n}\right)^k \binom{n}{k} \\ &\leq \sum_{k=0}^n \left(\frac{d}{n}\right)^k \binom{n}{k} \\ &= \left(1 + \frac{d}{n}\right)^n \leq e^d \end{aligned}$$

Hence, we have:

$$\left(\frac{en}{d}\right)^d \geq \sum_{k=0}^d \binom{n}{k} \geq S_{\mathcal{H}}(n)$$

□.

**Corollary 5.3: Sauer's lemma - bound on  $S_{\mathcal{H}}(n)$  III**

If  $V_{\mathcal{H}} = d > 2$ , for all  $n \geq d$ , we have:

$$S_{\mathcal{H}}(n) \leq n^d$$

**Proof** (Corollary 5.3). \_\_\_\_\_

If  $d > 2$  then by corollary 5.2, we have:

$$\frac{e}{d} < 1 \implies S_{\mathcal{H}}(n) \leq \left(\frac{en}{d}\right)^d \leq n^d$$

□.

**5.3 VC Theorem for classifiers****Theorem 5.2: VC Theorem (for classifiers)**

For any  $n \geq 1$  and  $\epsilon > 0$ , we have:

$$P\left(\sup_{h \in \mathcal{H}} |\widehat{R}_n(h) - R(h)| \geq \epsilon\right) \leq 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/32}$$

**Proof** (Theorem 5.2). \_\_\_\_\_

The proof for this theorem will be mentioned later in section 6.4. For now, we will assume that it is true to prove the following corollaries. □.

#### Corollary 5.4: Convergence of Empirical Risk (VC-Theorem)

If  $\widehat{h}_n$  is an empirical risk minimizer over  $\mathcal{H}$  then:

$$P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) \leq 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128}$$

**Proof** (Corollary 5.4). \_\_\_\_\_

We have:

$$\begin{aligned} P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) &\leq P\left(2 \sup_{h \in \mathcal{H}} |\widehat{R}_n(h) - R(h)| \geq \epsilon\right) \\ &= P\left(\sup_{h \in \mathcal{H}} |\widehat{R}_n(h) - R(h)| \geq \frac{\epsilon}{2}\right) \\ &\leq 8S_{\mathcal{H}}(n)e^{-n(\epsilon/2)^2/32} = 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128} \end{aligned}$$

□.

#### Corollary 5.5: Excess Risk of $\widehat{h}_n$ - $\delta \rightarrow \epsilon$ relation (VC-Theorem)

If  $V_{\mathcal{H}} < \infty$  then  $\mathcal{H}$  is uniformly learnable by ERM. Specifically, we can define the sample complexity as:

$$N(\epsilon, \delta) = \frac{128}{\epsilon^2} \left( V_{\mathcal{H}} \log(n+1) - \log \frac{\delta}{8} \right)$$

In other words, with probability of at least  $1 - \delta$ , we have:

$$R(\widehat{h}_n) \leq R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left( V_{\mathcal{H}} \log(n+1) - \log \frac{\delta}{8} \right)}$$

**Proof** (Corollary 5.5). \_\_\_\_\_

Let  $\delta = 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128}$ . By corollary 5.4, with probability of at least  $1 - \delta$ , we have:

$$\begin{aligned} R(\widehat{h}_n) &\leq R_{\mathcal{H}}^* + \epsilon \\ &= R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left( \log S_{\mathcal{H}}(n) - \log \frac{\delta}{8} \right)} \end{aligned}$$

By Sauer's lemma, we have that  $(n+1)^{V_{\mathcal{H}}} \geq S_{\mathcal{H}}(n)$  for all  $n \geq 1$ . Hence, we have:

$$\begin{aligned} R(\widehat{h}_n) &\leq R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left( \log S_{\mathcal{H}}(n) - \log \frac{\delta}{8} \right)} \\ &\leq R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left( V_{\mathcal{H}} \log(n+1) - \log \frac{\delta}{8} \right)} \end{aligned}$$

Hence, we conclude that  $\mathcal{H}$  is PAC-learnable by ERM when  $V_{\mathcal{H}} < \infty$  with the following sample complexity:

$$N(\epsilon, \delta) = \frac{128}{\epsilon^2} \left( V_{\mathcal{H}} \log(n+1) - \log \frac{\delta}{8} \right)$$

□.

## 5.4 VC Classes

**Definition 5.4** (VC Class).

A **VC Class** is a set of classifiers  $\mathcal{H}$  such that  $V_{\mathcal{H}} < \infty$ . In the following section, we will look at some class of classifiers where the VC dimension can be established or bounded.

**Example 1** (Hypercubes) : Consider the set of classifiers

$$\mathcal{H} = \left\{ \mathbf{1}_{\{x \in R\}} : R = \prod_{i=1}^d [a_i, b_i], a_i < b_i \right\}$$

In this case, for any  $d \geq 1$ , no more than  $2d + 1$  points can be shattered by  $\mathcal{H}$ . Hence,  $V_{\mathcal{H}} = 2d$ .



Figure 1: Example when  $d = 2$ . Four points can be shattered by  $\mathcal{H}$  but no five points can be shattered by  $\mathcal{H}$ .

**Example 2** (Convex sets in  $\mathbb{R}^2$ ) : Let  $\mathcal{X} = \mathbb{R}^2$ . Consider the set of classifiers

$$\mathcal{H} = \left\{ \mathbf{1}_{\{x \in C\}} : C \text{ is convex in } \mathbb{R}^2 \right\}$$

In this case,  $V_{\mathcal{H}} = \infty$  because for any  $n$  points on a circle and for any  $1 \leq k \leq n$ , we can draw a polygon that includes  $k$  points but not the remaining  $n - k$  points for any selection of  $k$  in  $n$  points (Figure 2).

**Example 3** (Finite  $|\mathcal{H}|$ ) : For any function class where  $|\mathcal{H}| < \infty$ , we have:

$$\begin{aligned} N_{\mathcal{H}}(x_1, \dots, x_n) &= \left| \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} \right| \leq |\mathcal{H}| \\ \implies S_{\mathcal{H}}(n) &\leq |\mathcal{H}| \\ \implies V_{\mathcal{H}} &\leq \log_2 |\mathcal{H}| \end{aligned}$$



Figure 2:  $\mathcal{H}$  can shatter any  $n$  points on a circle.

**Proposition 5.1: Steele & Dudley**

Let  $\mathcal{F}$  be the set of real-valued functions of the form:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} \right\}, \dim(\mathcal{F}) = m$$

Then, the following set of classifiers:

$$\mathcal{H} = \left\{ \mathbf{1}_{\{f(x) \geq 0\}} : f \in \mathcal{F} \right\}$$

Has finite VC dimension. Specifically,  $V_{\mathcal{H}} \leq m$ .

**Proof** (Proposition 5.1).

Suppose that  $\mathcal{H}$  shatters  $m + 1$  points in  $C = (x_1, \dots, x_{m+1})$ . Define the linear mapping  $L_C : \mathcal{F} \rightarrow \mathbb{R}^{m+1}$  such that:

$$L_C(f) = (f(x_1), \dots, f(x_{m+1}))^T$$

**Claim :**  $L_C(\mathcal{F})$  is a closed subspace in  $\mathbb{R}^{m+1}$

Let  $\{l_n\}_{n=1}^{\infty} \subset L_C(\mathcal{F})$  be a sequence and let  $l_n \rightarrow l$  as  $n \rightarrow \infty$ . We can always choose a function  $f \in \mathcal{F}$  such that:

$$f(x) = l, \forall x \in \mathbb{R}^m$$

Hence, for all  $\{l_n\}_{n=1}^{\infty}$  such that  $l_n \rightarrow l$ , the limit  $l \in L_C(\mathcal{F})$ . Therefore,  $L_C(\mathcal{F})$  is closed in  $\mathbb{R}^{m+1}$ .

By the Hilbert Projection Theorem [5], we have:

$$\mathbb{R}^{m+1} = L_C(\mathcal{F}) \oplus L_C(\mathcal{F})^{\perp}$$

Since  $\dim(\mathcal{F}) = m$ , we have  $\dim(L_C(\mathcal{F})) \leq m$ . Therefore,  $\dim(L_C(\mathcal{F})^{\perp}) \geq 1$  and we have:

$$\forall f \in \mathcal{F}, \exists \gamma \in \mathbb{R}^{m+1} \setminus \{0\} : \gamma^T L(f) = 0$$

Equivalently,

$$\sum_{i=1}^{m+1} \gamma_i f(x_i) = 0 \implies \sum_{i, \gamma_i \geq 0} \gamma_i f(x_i) = \sum_{j, \gamma_j < 0} -\gamma_j f(x_j)$$

Since  $\mathcal{H}$  shatters  $x_1, \dots, x_{m+1}$ . We define a classifier  $h \in \mathcal{H}$  such that:

$$h(x_i) = \begin{cases} 1 & \text{if } \gamma_i \geq 0 \\ 0 & \text{otherwise} \end{cases} \implies f(x_i) \geq 0 \iff \gamma_i \geq 0$$

This implies that  $\sum_{i: \gamma_i \geq 0} \gamma_i f(x_i) \geq 0$  and  $\sum_{j: \gamma_j < 0} -\gamma_j f(x_j) < 0$ , which is a contradiction. Therefore, we have  $V_{\mathcal{H}} < \infty$ .  $\square$ .

#### Corollary 5.6: Linear classifiers have finite $V_{\mathcal{H}}$

Let  $\mathcal{X} = \mathbb{R}^d$  and define a function space  $\mathcal{F}$  as:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = w^T x + b, w \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

Then, define the set of linear classifiers as:

$$\mathcal{H} = \left\{ \mathbf{1}_{\{w^T x + b \geq 0\}} \mid w \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

We have  $V_{\mathcal{H}} \leq d + 1$ .

**Proof.**

*Corollary 5.6* Since  $\dim(\mathcal{F}) = d + 1$ , the above corollary is a direct consequence of proposition 5.1.  $\square$ .

## 5.5 VC Theorem for sets

**Definition 5.5** (VC Theory for sets).

Let  $\mathcal{G}$  be a collection of subsets in  $\mathcal{X}$ . Let  $C \subset \mathcal{X}$  and  $C = \{x_1, \dots, x_n\}$ . We have the following definitions for restriction of  $\mathcal{G}$  to  $C$ , shattering coefficient and VC-dimension of  $\mathcal{G}$ :

$$\begin{aligned} N_{\mathcal{G}}(C) &= \left| \left\{ G \cap C : G \in \mathcal{G} \right\} \right| \\ S_{\mathcal{G}}(n) &= \sup_{C \subset \mathcal{X}; |C|=n} \left| N_{\mathcal{G}}(C) \right| \\ V_{\mathcal{G}} &= \sup \left\{ n : S_{\mathcal{G}}(n) = 2^n \right\} = \sup_{C \subset \mathcal{X}} \left\{ |C| : \mathcal{G} \text{ shatters } C \right\} \end{aligned}$$

**Remark :** In analogy to the definitions for classifier, sets and binary classifiers are equivalent via:

$$\begin{aligned} G &\rightarrow h_G(x) = \mathbf{1}_{\{x \in G\}} \\ h &\rightarrow G_h = \left\{ x : h(x) = 1 \right\} \end{aligned}$$

### Theorem 5.3: VC Theorem (for sets)

If  $X_1, \dots, X_n \sim Q$  are identically independently distributed samples. Then for any collection  $\mathcal{G}$  of subsets in  $\mathcal{X}$ ,  $\epsilon > 0$ , we have:

$$P\left(\sup_{G \in \mathcal{G}} \left| \widehat{Q}(G) - Q(G) \right| \geq \epsilon\right) \leq 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32}$$

Where  $\widehat{Q}(G)$  is defined (similar to the empirical CDF) as:

$$\widehat{Q}(G) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \in G\}}$$

**Proof** (Theorem 5.3). \_\_\_\_\_

Define the following class of classifiers:

$$\mathcal{H} = \left\{ h_G = \mathbf{1}_{\{x \in G\}} : G \in \mathcal{G} \right\}$$

Define a density function over  $\mathcal{X} \times \{0, 1\}$  such that  $\pi_0 = 1$ ,  $P_{X|Y=0} = Q$  and  $P_{X|Y=1}$  is arbitrary. For any  $h_G \in \mathcal{H}$ , we have:

$$\begin{aligned} R(h_G) &= P(h_G(X) \neq Y) \\ &= \pi_0 P_{X|Y=0}(h_G(X) \neq 0) + \pi_1 P_{X|Y=1}(h_G(X) \neq 1) \\ &= \pi_0 P_{X|Y=0}(h_G(X) = 1) + \pi_1 P_{X|Y=1}(h_G(X) = 0) \\ &= P_{X|Y=0}(h_G(X) = 1) \quad (\text{Since } \pi_0 = 1, \pi_1 = 0) \\ &= Q(G) \end{aligned}$$

Similarly, we have:

$$\widehat{R}_n(h_G) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{h_G(X_i)=1\}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \in G\}} = \widehat{Q}(G)$$

Therefore:

$$\begin{aligned} P\left(\sup_{G \in \mathcal{G}} \left| \widehat{Q}(G) - Q(G) \right| \geq \epsilon\right) &= P\left(\sup_{h_G \in \mathcal{H}} \left| \widehat{R}_n(h_G) - R(h_G) \right| \geq \epsilon\right) \\ &\leq 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/32} \quad (\text{Theorem 5.2}) \\ &= 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32} \end{aligned}$$

□.

### Corollary 5.7: Dvoretzky-Kiefer-Wolfowitz Inequality

Let  $X \sim Q$  be a random variable on the real line  $\mathbb{R}$  and denote  $G_t = (-\infty, t]$ . Then,

$$\begin{aligned} Q(G_t) &= P(X \leq t) = F(t) \quad (\text{CDF}) \\ \widehat{Q}(G_t) &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{x_i \leq t\}} = \widehat{F}_n(t) \quad (\text{Empirical CDF}) \end{aligned}$$

For all  $n \geq 1, \epsilon > 0$ , we have:

$$P\left(\sup_{t \in \mathbb{R}} \left| \widehat{F}_n(t) - F(t) \right| \geq \epsilon\right) \leq 8(n+1)e^{-n\epsilon^2/32}$$

**Proof** (Corollary 5.7). 

---

Let  $\mathcal{G} = \{G_t : t \in \mathbb{R}\}$ , by theorem 5.3, we have:

$$\begin{aligned} P\left(\sup_{G \in \mathcal{G}} |\widehat{Q}(G) - Q(G)| \geq \epsilon\right) &= P\left(\sup_{t \in \mathbb{R}} |\widehat{Q}(G_t) - Q(G_t)| \geq \epsilon\right) \\ &= P\left(\sup_{t \in \mathbb{R}} |\widehat{F}_n(t) - F(t)| \geq \epsilon\right) \\ &\leq 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32} \end{aligned}$$

For any set of  $n$  points on the real line, there are  $n+1$  ways to label them using half-open intervals. Hence  $S_{\mathcal{G}}(n) = n+1$ . Therefore:

$$P\left(\sup_{t \in \mathbb{R}} |\widehat{F}_n(t) - F(t)| \geq \epsilon\right) \leq 8(n+1)e^{-n\epsilon^2/32}$$

□.



## 5.6 End of chapter exercises

### Exercise 5.1

Determine the sample complexity  $N(\epsilon, \delta)$  for ERM over a class  $\mathcal{H}$  with VC dimension  $V_{\mathcal{H}} < \infty$ .

**Solution** (Exercise 5.1). \_\_\_\_\_

We have:

$$\begin{aligned}
 P\left(R(\widehat{h}_n) - R_{\mathcal{H}}^* \geq \epsilon\right) &\leq P\left(2 \sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \epsilon\right) \\
 &= P\left(\sup_{h \in \mathcal{H}} \left|\widehat{R}_n(h) - R(h)\right| \geq \frac{\epsilon}{2}\right) \\
 &\leq 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128} \quad (\text{Theorem 5.2}) \\
 &\leq 8(n+1)^{V_{\mathcal{H}}}e^{-n\epsilon^2/128} \quad (\text{Corollary 5.1})
 \end{aligned}$$

Now let:

$$\begin{aligned}
 \delta &= 8(n+1)^{V_{\mathcal{H}}}e^{-n\epsilon^2/128} \\
 \implies \log \frac{\delta}{8} &= V_{\mathcal{H}} \log(n+1) - \frac{n\epsilon^2}{128} \\
 \implies N(\epsilon, \delta) &= \frac{128}{\epsilon^2} \left( V_{\mathcal{H}} \log(n+1) - \log \frac{\delta}{8} \right)
 \end{aligned}$$

□.

### Exercise 5.2

Show that the VC Theorem for sets implies the VC Theorem for classifiers.

*Hint : Consider the sets of the form  $G' = G \times \{0\} \cup G^c \times \{1\} \subset \mathcal{X} \times \mathcal{Y}$ .*

**Solution** (Exercise 5.2). \_\_\_\_\_

Given an arbitrary class of classifiers  $\mathcal{H}$ . Define the following class of sets:

$$\mathcal{G} = \left\{ G_h \times \{0\} \cup G_h^c \times \{1\} : h \in \mathcal{H} \right\}$$

Where for a given  $h \in \mathcal{H}$ , we have:

$$G_h = \left\{ x \in \mathcal{X} : h(x) = 1 \right\}$$

Let  $P_{XY}$  be the density function over  $\mathcal{X} \times \mathcal{Y}$ . For any  $G \in \mathcal{G}$ , we have:

$$\begin{aligned}
 P_{XY}(G) &= \pi_0 P_{X|Y=0}(G) + \pi_1 P_{X|Y=1}(G) \\
 &= \pi_0 P_{X|Y=0}(G_h \times \{0\} \cup G_h^c \times \{1\}) + \pi_1 P_{X|Y=1}(G_h \times \{0\} \cup G_h^c \times \{1\}) \\
 &= \pi_0 P_{X|Y=0}(G_h) + \pi_1 P_{X|Y=1}(G_h^c) \\
 &= \pi_0 P_{X|Y=0}(h(X) = 1) + \pi_1 P_{X|Y=1}(h(X) = 0) \\
 &= P(h(X) \neq Y) \\
 &= R(h)
 \end{aligned}$$

Let  $Q = P_{XY}$ . We also have:

$$\widehat{Q}(G) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{(X_i, Y_i) \in G_h\}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}} = \widehat{R}_n(h)$$

From the above, we have:

$$\begin{aligned} P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R}_n(h) - R(h) \right| \geq \epsilon\right) &= P\left(\sup_{G \in \mathcal{G}} \left| \widehat{Q}(G) - Q(G) \right| \geq \epsilon\right) \\ &\leq 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32} \\ &= 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/32} \end{aligned}$$

□.

### Exercise 5.3

Let  $\mathcal{G}_1$  and  $\mathcal{G}_2$  denote two classes of sets:

- (a)  $\mathcal{G}_1 \cap \mathcal{G}_2 = \left\{ G_1 \cap G_2 : G_1 \in \mathcal{G}_1, G_2 \in \mathcal{G}_2 \right\}$ .
- (b)  $\mathcal{G}_1 \cup \mathcal{G}_2 = \left\{ G_1 \cup G_2 : G_1 \in \mathcal{G}_1, G_2 \in \mathcal{G}_2 \right\}$ .

Show that:

## 6 Rademacher Complexity

### 6.1 Bounded Difference Inequality

### 6.2 Rademacher Complexity

### 6.3 Bounds for binary classification

### 6.4 Proof of VC Inequality

## A Related topics

### A.1 Neyman-Pearson Lemma

#### A.1.1 Type I & Type II errors

**Overview** : In a hypothesis test, we are interested in testing a given null hypothesis  $H_0$  against some alternative hypothesis  $H_1$ . Hence, we define some rejection region  $\mathcal{R} \subset \mathbb{R}$  such that:

$$x \in \mathcal{R} \implies \text{reject } H_0$$

Equivalently, denote that  $\bar{\mathcal{R}}$  is the acceptance region. We define the following conditional probability densities:

- $f_1(x)$  : Density given that  $H_1$  is true.
- $f_0(x)$  : Density given that  $H_0$  is true.

**Definition A.1** (Type I & Type II errors).

*In hypothesis testing, we define the type I error as the probability that we falsely reject the null hypothesis given that the null hypothesis is true. On the other hands, type II error is the probability that we falsely accept the null hypothesis given that the hypothesis is not true:*

$$\alpha = P_{H_0}(\mathcal{R}) = \int_{\mathcal{R}} f_0(x)dx$$
$$\beta = P_{H_1}(\bar{\mathcal{R}}) = 1 - \int_{\mathcal{R}} f_1(x)dx$$

There is a trade-off between type I and type II errors as illustrated in the figure below:

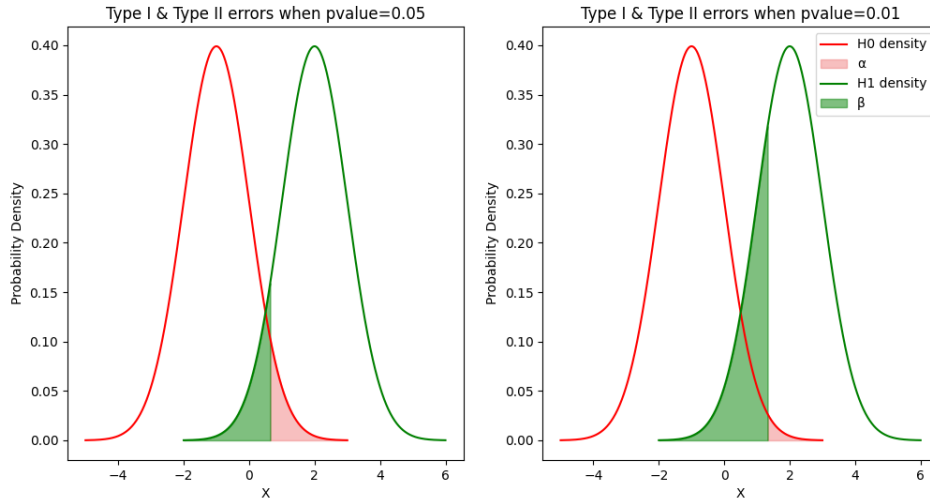


Figure 3: Trade-off between type I and type II errors

**Definition A.2** (Power of hypothesis test).

Given a hypothesis test used to test a null hypothesis  $H_0$  against an alternative hypothesis  $H_1$ . The probability:

$$P_{H_1}(\mathcal{R}) = \int_{\mathcal{R}} f_1(x)dx = 1 - \beta$$

Which denotes the probability that we correctly reject the null hypothesis given that  $H_1$  is true is called the **Power** of the hypothesis test. Later on we will see that using **Neyman-Pearson Lemma**, we can prove any hypothesis test has the power of at most the likelihood ratio test's power.

### A.1.2 Neyman-Pearson Lemma

**Overview** : The Neyman-Pearson Lemma is concerned with maximizing the power of hypothesis test subjected to a certain degree of type I error. Formally, we are trying to solve the following constrained optimization problem:

$$\begin{aligned} &\text{maximize : } P_{H_1}(\mathcal{R}) = \int_{\mathcal{R}} f_1(x)dx \\ &\text{subjected to : } P_{H_0}(\mathcal{R}) = \int_{\mathcal{R}} f_0(x)dx \leq \alpha \end{aligned}$$

#### Theorem A.1: Neyman-Pearson Lemma

Let  $H_0$  and  $H_1$  be simple hypotheses. For a constant  $c > 0$ , suppose the likelihood ratio test rejects  $H_0$  when  $L(X) > c$  has significance level  $\alpha \in (0, 1)$ . **Then for any other test of  $H_0$  with significance level of at most  $\alpha$ , its power against  $H_1$  is at most the power of the likelihood ratio test.**

$$\mathcal{R} = \left\{ x \in \mathbb{R} : L(x) = \frac{f_1(x)}{f_0(x)} > c \right\}$$

**Proof** (Theorem A.1).

Note that the rejection region  $\mathcal{R} = \left\{ x \in \mathbb{R} : L(x) = \frac{f_1(x)}{f_0(x)} > c \right\}$  maximizes the quantity:

$$\int_{\mathcal{R}} (f_1(x) - cf_0(x))dx$$

Because  $f_1(x) - cf_0(x) < 0$  for all  $x \notin \mathcal{R}$ . Therefore, for any other test with rejection region  $\mathcal{R}'$  with significance level of at most  $\alpha$ , we have:

$$\begin{aligned} \int_{\mathcal{R}} (f_1(x) - cf_0(x))dx &\geq \int_{\mathcal{R}'} (f_1(x) - cf_0(x))dx \\ \implies P_{H_1}(\mathcal{R}) - P_{H_1}(\mathcal{R}') &\geq c \left( \int_{\mathcal{R}} f_0(x)dx - \int_{\mathcal{R}'} f_0(x)dx \right) \\ &= c \left( \alpha - \int_{\mathcal{R}'} f_0(x)dx \right) \end{aligned}$$

Since  $\int_{\mathcal{R}'} f_0(x)dx \leq \alpha$ , we have:

$$P_{H_1}(\mathcal{R}) - P_{H_1}(\mathcal{R}') \geq 0 \implies P_{H_1}(\mathcal{R}') \leq P_{H_1}(\mathcal{R})$$

Hence, for any test with significance level of at most  $\alpha$ , the power is at most the power of the likelihood ratio test  $P_{H_1}(\mathcal{R})$ .  $\square$ .

## B List of Definitions

1.1	Definition (Classifier ( $h$ ))	3
1.2	Definition (Decomposition of $P_{XY}$ )	3
1.3	Definition (Hypothesis space ( $\mathcal{H}$ ))	4
1.4	Definition (Learning algorithm ( $\mathcal{L}_n$ ))	4
1.5	Definition (Risk ( $R(h)$ ))	5
1.6	Definition (Bayes Risk ( $R^*$ ))	5
1.7	Definition (Consistency of learning algorithms)	5
2.1	Definition (Plug-in classifier)	9
3.1	Definition (Empirical Risk ( $\widehat{R}_n$ ))	16
4.1	Definition (Empirical Risk Minimization ( $\widehat{h}_n$ ))	22
4.2	Definition (Uniform Deviation Bounds (UDB))	22
4.3	Definition (PAC & Sample Complexity ( $N(\epsilon, \delta)$ ))	25
5.1	Definition (Restriction ( $N_{\mathcal{H}}$ ))	31
5.2	Definition (Shattering Coefficient ( $S_{\mathcal{H}}$ ))	31
5.3	Definition (VC-dimension ( $V_{\mathcal{H}}$ ))	31
5.4	Definition (VC Class)	36
5.5	Definition (VC Theory for sets)	38
A.1	Definition (Type I & Type II errors)	44
A.2	Definition (Power of hypothesis test)	45

## C Important Theorems

2.1	Properties of Bayes classifier	6
2.2	Properties of Bayes classifier (Multi-class)	8
3.1	Hoeffding's Inequality	15
4.1	Uniform Deviation Bounds for finite $\mathcal{H}$	22
5.1	Sauer's Lemma	32
5.2	VC Theorem (for classifiers)	34
5.3	VC Theorem (for sets)	39
A.1	Neyman-Pearson Lemma	45

## D Important Corollaries

2.1	Excess risk of plug-in classifier	9
3.1	Chebyshev's Inequality	14
3.2	Chernoff's bounding method	14
3.3	Convergence of Empirical Risk	17
4.1	Excess Risk of $\widehat{h}_n$ - $\delta \rightarrow \epsilon$ relation	24
5.1	Sauer's lemma - bound on $S_{\mathcal{H}}(n)$ I	33
5.2	Sauer's lemma - bound on $S_{\mathcal{H}}(n)$ II	34
5.3	Sauer's lemma - bound on $S_{\mathcal{H}}(n)$ III	34
5.4	Convergence of Empirical Risk (VC-Theorem)	35
5.5	Excess Risk of $\widehat{h}_n$ - $\delta \rightarrow \epsilon$ relation (VC-Theorem)	35
5.6	Linear classifiers have finite $V_{\mathcal{H}}$	38
5.7	Dvoretzky-Kiefer-Wolfowitz Inequality	39

## E Important Propositions

1.1	Law of total expectation . . . . .	3
2.1	Likelihood ratio test . . . . .	9
3.1	Markov's Inequality . . . . .	14
3.2	KL-divergence hypothesis testing . . . . .	18
4.1	(Probabilistic) Bound on Excess Risk of $\widehat{h}_n$ . . . . .	23
4.2	(Non-probabilistic) Bound on Excess Risk of $\widehat{h}_n$ . . . . .	24
4.3	Zero-error case bound . . . . .	26
5.1	Steele & Dudley . . . . .	37



## F References

### References

- [1] Rick Durrett. *Probability: Theory and Examples*. 4th. USA: Cambridge University Press, 2010. ISBN: 0521765390.
- [2] Zhou Fan. *Statistics 200: Introduction to Statistical Inference - Lecture 6*. <https://web.stanford.edu/class/archive/stats/stats200/stats200.1172/Lecture06.pdf>. [Accessed 07-01-2024]. 2016.
- [3] Aurelien Garivier. *Lecture notes in Machine Learning Theory*. 2019. URL: [https://www.math.univ-toulouse.fr/~agarivie/sites/default/files/5\\_VC.pdf](https://www.math.univ-toulouse.fr/~agarivie/sites/default/files/5_VC.pdf).
- [4] Erhan undefinedinlar. *Probability and Stochastics*. Springer New York, 2011. ISBN: 9780387878591. DOI: [10.1007/978-0-387-87859-1](https://doi.org/10.1007/978-0-387-87859-1). URL: <http://dx.doi.org/10.1007/978-0-387-87859-1>.
- [5] Wikipedia. *Hilbert projection theorem* — *Wikipedia, The Free Encyclopedia*. <http://en.wikipedia.org/w/index.php?title=Hilbert%20projection%20theorem&oldid=1172787172>. [Online; accessed 11-January-2024]. 2024.
- [6] Wikipedia. *Hoeffding's lemma* — *Wikipedia, The Free Encyclopedia*. <http://en.wikipedia.org/w/index.php?title=Hoeffding's%20lemma&oldid=1114715065>. [Online; accessed 04-January-2024]. 2024.
- [7] Wikipedia. *Rényi entropy* — *Wikipedia, The Free Encyclopedia*. <http://en.wikipedia.org/w/index.php?title=R%C3%A9nyi%20entropy&oldid=1190869396>. [Online; accessed 05-January-2024]. 2024.
- [8] Wikipedia. *Vitali set* — *Wikipedia, The Free Encyclopedia*. <http://en.wikipedia.org/w/index.php?title=Vitali%20set&oldid=1187241923>. [Online; accessed 24-December-2023]. 2023.