Statistical Learning Theory Notes

Nong Minh $\mathrm{Hieu^1}$

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

Contents

| 1 | Probability settings 3 | | | | | | |
|---|---|--|----|--|--|--|--|
| | 1.1 Classification problem | | 3 | | | | |
| | 1.2 Goal of classification | | 5 | | | | |
| 2 | Bayes classifier | | 6 | | | | |
| | 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 | | | | |
| 3 | Hoeffding's inequality 14 | | | | | | |
| | 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 | | | | |
| 4 | Empirical Risk Minimization 22 | | | | | | |
| | 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 | | | | |
| 5 | Vapnik-Chevronenkis Theory 31 | | | | | | |
| | 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 | | | | |
| 6 | Rademacher Complexity | | 44 | | | | |
| - | 6.1 Bounded Difference Inequality | | 44 | | | | |
| | 6.2 Rademacher Complexity | | 46 | | | | |
| | 6.3 Bounds for binary classification | | 47 | | | | |
| | 6.4 Proof of VC Inequality | | 47 | | | | |
| | 1 T T T T T T T T T T T T T T T T T T T | | | | | | |

| \mathbf{A} | Related topics | | | | | |
|--------------|------------------------|-------|--|----|--|--|
| | A.1 | Neyma | an-Pearson Lemma | 48 | | |
| | | A.1.1 | Type I & Type II errors | 48 | | |
| | | A.1.2 | Neyman-Pearson Lemma | 49 | | |
| | A.2 | Raden | nacher Complexity bound for linear function classes | 50 | | |
| | | A.2.1 | Problem Statement | 50 | | |
| | | A.2.2 | Covering Number | 50 | | |
| | | A.2.3 | Massart's Lemma | 52 | | |
| | | A.2.4 | Dudley's Theorem | 52 | | |
| | | A.2.5 | Bound on covering number of linear function class | 54 | | |
| | | A.2.6 | Rademacher Complexity bound for linear functions class | 54 | | |
| В | List of Definitions 5 | | | | | |
| \mathbf{C} | Important Theorems | | | | | |
| D | Important Corollaries | | | | | |
| \mathbf{E} | Important Propositions | | | | | |
| \mathbf{F} | References | | | | | |

1 Probability settings

1.1 Classification problem

- classification problems, we consider pairs (x,y) where $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Wh
- *Y* is the space of labels.

• \mathcal{X} is the space of **feature vectors**.

A classifier is a function $h: \mathcal{X} \to \mathcal{Y}$ 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:

$$P_{XY} = P_{X|Y}P_Y$$
$$P_{XY} = P_{Y|X}P_X$$

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

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

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:

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

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

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

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

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

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

$$\mathcal{H} = \Big\{ h: \mathcal{X} \rightarrow \mathcal{Y} \Big| some \ common \ properties \Big\}$$

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

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

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

To learn a classifier $h: \mathcal{X} \to \mathcal{Y}$, 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 $\hat{h}_n: \mathcal{X} \to \mathcal{Y}$ from the training data.

$$\mathcal{L}_n: (\mathcal{X} \times \mathcal{Y})^n \to \mathcal{H}$$

 $\hat{h}_n = \mathcal{L}_n((X_1, Y_1), \dots, (X_n, Y_n))$

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 $h: \mathcal{X} \to \mathcal{Y}$, not just for $h \in \mathcal{H}$:

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

Definition 1.7 (Consistency of learning algorithms).

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

• Weakly consistent if $R(\hat{h}_n) \stackrel{p}{\rightarrow} R^*$:

$$\lim_{n \to \infty} P(R(\hat{h}_n) \le r) = P(R^* \le r), \ \forall r \ge 0$$

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

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

• 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:

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

Theorem 2.1: Properties of Bayes classifier

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

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

Proof (Theorem 2.1).

Proving each point:

(i)
$$R(h^*) = \inf_{h:\mathcal{X}\to\mathcal{Y}} \{R(h)\} = R^*$$
.
For all $h:\mathcal{X}\to\mathcal{Y}$, we have:

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

$$= \mathbb{E}_{x \sim X} \left[\mathbb{E}_{Y|X=x} \left[\mathbf{1}_{\{Y \neq h(x)\}} \right] \right]$$

$$= \mathbb{E}_{x \sim X} \left[\sum_{y \in \{0,1\}} \mathbf{1}_{\{y \neq h(x)\}} \right]$$

$$= \mathbb{E}_{x \sim X} \left[\eta(x) \mathbf{1}_{\{h(x)=0\}} + (1 - \eta(x)) \mathbf{1}_{\{h(x)=1\}} \right]$$

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 $\eta(x) \ge 1 - \eta(x) \implies \eta(x) \ge \frac{1}{2}$. Therefore, we have:

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

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

We have:

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

(iii) $R^* = \mathbb{E}\Big[\min(\eta(X), 1 - \eta(x))\Big]$. From (i) we have:

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

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:

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

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

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

 \Box .

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:

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

Hence, we have:

$$\eta(x) \ge \frac{1}{2} \iff \frac{\pi_0 p_0(x)}{\pi_1 p_1(x)}$$

$$\iff \frac{p_1(x)}{p_0(x)} \ge \frac{\pi_0}{\pi_1}$$

Proposition 2.1: Likelihood ratio test

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

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

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 η based on n training samples $\{(X_i,Y_i)\}_{i=1}^n$. We define $\widehat{h_n}$ as:

$$\widehat{h_n} = \begin{cases} 1 & \text{if } \widehat{\eta_n}(x) \ge \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Corollary 2.1: Excess risk of plug-in classifier

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

$$R(\widehat{h_n}) - R^* \le 2\mathbb{E}_X \left[\left| \eta(X) - \widehat{\eta_n}(X) \right| \right]$$

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:

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

Case 1: $\widehat{\eta_n}(X) \ge \frac{1}{2}, \eta(X) < \frac{1}{2}$ We have:

$$\begin{split} \eta(X) - \widehat{\eta_n}(X) &\leq \eta(X) - \frac{1}{2} \quad (Both \ sides \ negative) \\ \Longrightarrow \left| \eta(X) - \widehat{\eta_n}(X) \right| &\geq \left| \eta(X) - \frac{1}{2} \right| \end{split}$$

Case 2: $\widehat{\eta_n}(X) < \frac{1}{2}, \eta(X) \ge \frac{1}{2}$

 $We\ have:$

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

Therefore, we have:

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

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

$$R(\widehat{h_n}) - R^* \le 2\mathbb{E}_X \left[\left| \eta(X) - \widehat{\eta_n}(X) \right| \right]$$

2.4 End of chapter exercises

Exercise 2.1

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

Solution (Exercise 2.1).

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

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

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

We have:

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

(i) Calculate Bayes risk R*

For any classifier $h: \mathcal{X} \to \mathcal{Y}$, 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 maxmized for all $x \in \mathcal{X}$. Hence, we have:

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

Therefore, we have $h^*(x) = \arg \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\}$ is the Bayes classifier and the Bayes risk $R^* = \mathbb{E}_{x \sim X} \left[1 - \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\} \right]$.

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

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

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

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

$$\sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\} = \sum_{y \in \mathcal{Y}; y \neq y_x} \eta_y(x) + \eta_{y_x^*}(x)$$

$$= \sum_{y \in \mathcal{Y}; y \neq y_x^*} \eta_y(x) + \eta_{y_x^*}(x)$$

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

$$\implies \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\} - 1 = 0$$

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

$$\sum_{y \in \mathcal{Y}} \mathbf{1}_{\{h(x) \neq y\}} \eta_y(x) + \max_{y \in \mathcal{Y}} \left\{ \eta_y(x) \right\} - 1 = \sum_{y \in \mathcal{Y}; y \neq y_x} \eta_y(x) + \eta_{y_x^*}(x) - 1$$

$$= 2\eta_{y_x^*}(x) - 1 + \sum_{y \in \mathcal{Y} \setminus \{y_x, y_x^*\}} \eta_y(x)$$

$$= 2\eta_{y_x^*}(x) - \left(\eta_{y_x}(x) + \eta_{y_x^*}(x) \right)$$

$$= \eta_{y_x^*}(x) - \eta_{y_x}(x).$$

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:

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

 \Box .

Where $\overline{\eta_y}(x) = P(Y \neq y | X = x)$.

Exercise 2.2

Define the α -cost-sensitive risk of a classifier $h: \mathcal{X} \to \mathcal{Y}$ 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 and analogue of theorem 2.1.

Solution (Exercise 2.2).

Using the law of total expectation, we have:

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

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:

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

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

$$\eta(x) \ge \alpha \implies \frac{1}{1 + \frac{\pi_0 p_0(x)}{\pi_1 p_1(x)}} \ge \alpha$$

$$\implies 1 + \frac{\pi_0 \cdot p_0(x)}{\pi_1 \cdot p_1(x)} \le \frac{1}{\alpha}$$

$$\implies \frac{p_1(x)}{p_0(x)} \ge \frac{\alpha}{1 - \alpha} \cdot \frac{\pi_0}{\pi_1}$$

Hence, we can rewrite the Bayes classifier as followed:

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

(i) Bayes Risk R_{α}^* We have:

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

(ii) Excess Risk $R_{\alpha}(h) - R_{\alpha}^{*}$ For an arbitrary $h : \mathcal{X} \to \mathcal{Y}$, we have:

$$R_{\alpha}(h) - R_{\alpha}^{*} = \mathbb{E}_{x \sim X} \Big[(1 - \alpha) \eta(x) \Big(\mathbf{1}_{\{h(x) = 0\}} - \mathbf{1}_{\{h^{*}(x) = 0\}} \Big) + \alpha (1 - \eta(x)) \Big(\mathbf{1}_{\{h(x) = 1\}} - \mathbf{1}_{\{h^{*}(x) = 1\}} \Big) \Big]$$

$$= \mathbb{E}_{x \sim X} \Big[(1 - \alpha) \eta(x) \Big(\mathbf{1}_{\{h(x) = 0, h^{*}(x) = 1\}} - \mathbf{1}_{\{h(x) = 1, h^{*}(x) = 0\}} \Big)$$

$$+ \alpha (1 - \eta(x)) \Big(\mathbf{1}_{\{h(x) = 1, h^{*}(x) = 0\}} - \mathbf{1}_{\{h(x) = 0, h^{*}(x) = 1\}} \Big) \Big]$$

$$= \mathbb{E}_{x \sim X} \Big[\mathbf{1}_{\{h(x) = 0, h^{*}(x) = 1\}} (\eta(x) - \alpha) + \mathbf{1}_{\{h(x) = 1, h^{*}(x) = 0\}} (\alpha - \eta(x)) \Big]$$

$$= \mathbb{E}_{X} \Big[|\eta(X) - \alpha| \mathbf{1}_{\{h(X) \neq h^{*}(X)\}} \Big]$$

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:

$$P(U \ge t) \le \frac{1}{t} \mathbb{E}[U]$$

Proof (Proposition 3.1). _

We have:

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

Corollary 3.1: Chebyshev's Inequality

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

$$P(\left|Z - \mu\right| \ge t) \le \frac{\sigma^2}{t^2}$$

 \Box .

 \Box .

Proof (Corollary 3.1).

Using Markov's inequality, we have:

$$P(\left|Z - \mu\right| \ge t) = P(\left|Z - \mu\right|^2 \ge t^2)$$

$$\le \frac{\mathbb{E}\left[\left|Z - \mu\right|^2\right]}{t^2} = \frac{\sigma^2}{t^2}$$

Corollary 3.2: Chernoff's bounding method

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

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

Proof (Corollary 3.2).

We have:

$$P(Z \ge t) = P(sZ \ge st), \quad (t > 0)$$

$$= P(e^{sZ} \ge e^{st})$$

$$\le \frac{\mathbb{E}\left[e^{sZ}\right]}{e^{st}} = e^{-st}M_Z(s) \quad (Markov's inequality)$$

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

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

 \Box .

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:

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

Proof (Lemma 3.1).

(The proof for this lemma can be found here [6]).

 \Box .

Theorem 3.1: Hoeffding's Inequality

Let Z_1, Z_2, \ldots, 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:

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

Proof (Theorem 3.1).

Using the Chernoff's bounds, we have:

$$P(\left|S_{n} - \mathbb{E}[S_{n}]\right| \geq t) \leq \inf_{s>0} e^{-st} M_{S_{n} - \mathbb{E}[S_{n}]}(s)$$

$$= \inf_{s>0} e^{-st} \mathbb{E}\left[e^{s(S_{n} - \mathbb{E}[S_{n}])}\right]$$

$$= \inf_{s>0} e^{-st} \mathbb{E}\left[\exp\left(s\sum_{i=1}^{n} (Z_{i} - \mathbb{E}[Z_{i}])\right)\right]$$

$$= \inf_{s>0} e^{-st} \mathbb{E}\left[\prod_{i=1}^{n} \exp\left(s(Z_{i} - \mathbb{E}[Z_{i}])\right)\right]$$

$$= \inf_{s>0} e^{-st} \prod_{i=1}^{n} \mathbb{E}\left[\exp\left(s(Z_{i} - \mathbb{E}[Z_{i}])\right)\right] \quad (Since \ all \ Z_{i} - \mathbb{E}[Z_{i}] \ are \ independent)$$

$$\leq \inf_{s>0} e^{-st} \prod_{i=1}^{n} \exp\left(\frac{s^{2}(b_{i} - a_{i})^{2}}{8}\right) \quad (By \ Hoeffding's \ lemma)$$

$$= \inf_{s>0} \exp\left(-st + \sum_{i=1}^{n} \frac{s^{2}(b_{i} - a_{i})^{2}}{8}\right)$$

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:

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

 \Box .

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 $h: \mathcal{X} \to \mathcal{Y}$ 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 $n\widehat{R}_n(h) \sim Binomial(n, R(h))$. 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 $h: \mathcal{X} \to \mathcal{Y}$ be a classifier, we have:

$$P(\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon) \le 2e^{-2n\epsilon^2}, \quad \epsilon > 0$$

Proof (Corollary 3.3).

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

Using the Hoeffding's inequality, we have:

$$P(\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon) = P(\left|\widehat{R_n}(h) - \mathbb{E}[\widehat{R_n}(h)]\right| \ge \epsilon)$$

$$= P\left(\left|n\widehat{R_n}(h) - \mathbb{E}[n\widehat{R_n}(h)]\right| \ge n\epsilon\right)$$

$$\le \exp\left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad (Hoeffding's inequality)$$

$$= e^{-2n\epsilon^2}$$

 \Box .

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

Now suppose $X_1, \ldots, X_n \sim p_y$ 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:

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

Where $x = (x_1, ..., x_n)$ is an observation of the random vector $X = (X_1, ..., X_n)$. 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:

$$R_0(\widehat{h_n}) \le e^{-2nD(p_0||p_1)^2/c}$$
, 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:

$$\widehat{h_n}(X) = \begin{cases} 1 & \text{if } \widehat{S_n}(X_1, \dots, X_n) \ge 0 \\ 0 & \text{otherwise} \end{cases}$$

Where we have:

$$\widehat{S_n}(X_1, \dots, X_n) = \log \frac{\prod_{i=1}^n p_1(X_i)}{\prod_{i=1}^n p_0(X_i)}$$

$$= \sum_{i=1}^n \log \frac{p_1(X_i)}{p_0(X_i)}$$

$$= \sum_{i=1}^n Z_i \quad \left(Letting \ Z_i = \log \frac{p_1(X_i)}{p_0(X_i)} \right)$$

Since the likelihoods are bounded, we have:

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

Now, we have:

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

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

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

Therefore, we have:

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

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

$$b_i - a_i = \log \frac{\beta}{\alpha} - \log \frac{\alpha}{\beta}$$

$$= \log \frac{\beta^2}{\alpha^2} = 2 \log \frac{\beta}{\alpha} = 2(\log \beta - \log \alpha)$$

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

Finally, we have:

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

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

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

3.5 End of chapter exercises

Exercise 3.1

- (i) Apply Chernoff's bounding method to obtain an exponential bound on the tail probability $P(Z \ge t)$ 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:

$$P(Z \ge t) \le \inf_{s>0} e^{-st} M_Z(s)$$
$$= \inf_{s>0} \exp\left(-st + \mu s + \frac{1}{2}\sigma^2 s^2\right)$$

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:

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

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

$$P(S_n - \mathbb{E}[S_n] \ge \epsilon) = P\left(\overline{X} - \frac{\mathbb{E}[S_n]}{n} \ge \frac{\epsilon}{n}\right)$$
$$= P\left(\overline{X} - p \ge \frac{\epsilon}{n}\right)$$
$$= P\left(\frac{\overline{X} - p}{\sqrt{pq}/\sqrt{n}} \ge \frac{\epsilon}{\sqrt{npq}}\right), \quad (q = 1 - p)$$

By the Central Limit Theorem, we have:

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

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

$$P(S_n - \mathbb{E}[S_n] \ge \epsilon) = P\left(\frac{\overline{X} - p}{\sqrt{pq}/\sqrt{n}} \ge \frac{\epsilon}{\sqrt{npq}}\right)$$

$$\le \exp\left(-\frac{\epsilon^2}{2npq}\right) = \exp\left(-\frac{\epsilon^2}{2Var(S_n)}\right)$$

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

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

Exercise 3.2

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

Solution (Exercise 3.2).

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

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

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

$$R_0(\widehat{h_n}) = P(S_n \ge 0 | Y = 0)$$

$$\le \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[e^{sZ_i} \right]$$

$$= \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[\exp \left(s \log \frac{p_1(X_i)}{p_0(X_i)} \right) \right]$$

$$= \inf_{s>0} \prod_{i=1}^n \mathbb{E}_{q_0} \left[\frac{p_1(X_i)^s}{p_0(X_i)^s} \right]$$

Taking logarithm from both sides, we have:

$$\log R_0(\widehat{h_n}) \le \inf_{s>0} \sum_{i=1}^n \log \mathbb{E}_{q_0} \left[\frac{p_1(X_i)^s}{p_0(X_i)^s} \right]$$
$$= \inf_{s>0} \sum_{i=1}^n (s-1) R_s(p_1||p_0)$$
$$= \inf_{s>0} n(s-1) R_s(p_1||p_0)$$

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:

$$P(\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon) \le \delta$$

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 ϵ and has the following form:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \le \epsilon\right) \ge 1 - \delta$$

$$Or: P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon\right) \le \delta$$

The above bounds have the following interpretations:

- The probability that the deviation from the true risk is at most ϵ 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 ϵ is at most δ .

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:

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

Proof (Theorem 4.1). _

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

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

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:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon\right) = P(\Omega_{\epsilon}(\mathcal{H}))$$

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

$$\le \sum_{h\in\mathcal{H}}P(\Omega_{\epsilon}(h))$$

$$\le \sum_{h\in\mathcal{H}}2e^{-2n\epsilon^2} = 2|\mathcal{H}|e^{-2n\epsilon^2}$$

 \Box .

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

Suppose that \mathcal{H} satisfies:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon\right) \le \delta$$

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

$$R(\widehat{h_n}) - R_{\mathcal{H}}^* \le 2\epsilon$$

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

Proof (Proposition 4.1).

We have:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon\right) \le \delta \implies P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \le \epsilon\right) \ge 1 - \delta$$

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

$$\left| \widehat{R_n}(h) - R(h) \right| \le \epsilon \implies -\epsilon \le \widehat{R_n}(h) - R(h) \le \epsilon$$

$$\implies \begin{cases} \widehat{R_n}(h) & \le R(h) + \epsilon \\ \\ R(h) & \le \widehat{R_n}(h) + \epsilon \end{cases}$$

Therefore:

$$R(\widehat{h_n}) \leq \widehat{R_n}(\widehat{h_n}) + \epsilon$$

$$\leq \widehat{R_n}(h) + \epsilon \quad (Since \ \widehat{h_n} \ minimizes \ the \ Empirical \ Risk)$$

$$\leq \left(R(h) + \epsilon\right) + \epsilon = R(h) + 2\epsilon$$

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

$$R(\widehat{h_n}) \le \inf_{h \in \mathcal{H}} R(h) + 2\epsilon$$

= $R_{\mathcal{H}}^* + 2\epsilon$

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

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:

$$R(\widehat{h_n}) - R_{\mathcal{H}}^* \le 2 \sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right|$$

Proof (Proposition 4.2).

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

$$R(\widehat{h_n}) - R_{\mathcal{H}}^* \le \left| R(\widehat{h_n}) - \widehat{R_n}(\widehat{h_n}) \right| + \widehat{R_n}(\widehat{h_n}) - \widehat{R_n}(h_{\mathcal{H}}^*) + \left| \widehat{R_n}(h_{\mathcal{H}}^*) - R_{\mathcal{H}}^* \right|$$

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:

$$R(\widehat{h_n}) - R_{\mathcal{H}}^* \le \left| R(\widehat{h_n}) - \widehat{R_n}(\widehat{h_n}) \right| + \left| \widehat{R_n}(h_{\mathcal{H}}^*) - R_{\mathcal{H}}^* \right|$$
$$\le 2 \sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right|$$

Corollary 4.1: Excess Risk of $\widehat{h_n}$ - $\delta \to \epsilon$ relation

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

$$P(R(\widehat{h_n}) - R_{\mathcal{H}}^* \ge \epsilon) \le 2|\mathcal{H}|e^{-n\epsilon^2/2}$$

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

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

 \Box .

Proof (Corollary 4.1).

By proposition 4.2, we have:

$$P(\widehat{R(h_n)} - R_{\mathcal{H}}^* \ge \epsilon) \le P\left(2 \sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \epsilon)\right)$$

$$= P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \frac{\epsilon}{2})\right)$$

$$\le 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2}\right)$$

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:

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

 \Box .

4.2 PAC Learning & Sample Complexity

$$\forall \epsilon, \delta > 0 : n \ge N(\epsilon, \delta) \implies P(R(\widehat{h_n}) - R_{\mathcal{H}}^* \ge \epsilon) \le \delta$$

Where we have:

- $N(\epsilon, \delta)$ is called the **Sample Complexity**.
- 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:

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

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 $\Omega_0(h) = \left\{\widehat{R_n}(h) = 0\right\}$ and define the event Ω_{ϵ} as:

$$\Omega_{\epsilon} = \bigcup_{h \in \mathcal{H}: R(h) > \epsilon} \Omega_0(h) = \left\{ \exists h \in \mathcal{H} : \widehat{R_n}(h) = 0, R(h) \ge \epsilon \right\}$$

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

$$P(\Omega_0(h)) = P\left(\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}} = 0\right)$$

$$= P\left(\sum_{i=1}^n \mathbf{1}_{\{h(X_i) \neq Y_i\}} = 0\right)$$

$$= P\left(\bigcup_{i=1}^n \left\{h(X_i) = Y_i\right\}\right)$$

$$= \prod_{i=1}^n P(h(X_i) = Y_i) \quad (Since all (X_i, Y_i) pairs are independent)$$

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

$$P(\Omega_0(h)) = \prod_{i=1}^n P(h(X_i) = Y_i)$$
$$= (1 - R(h))^n$$
$$< (1 - \epsilon)^n$$

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

$$P(\Omega_0(h)) \le (1 - \epsilon)^n = e^{n \log(1 - \epsilon)}$$

$$< e^{-n\epsilon}$$

Finally, we have:

$$P(\Omega_{\epsilon}) = P\left(\bigcup_{h \in \mathcal{H}; R(h) \ge \epsilon} \Omega_{0}(h)\right)$$

$$\leq \sum_{h \in \mathcal{H}; R(h) \ge \epsilon} P(\Omega_{0}(h))$$

$$\leq \sum_{h \in \mathcal{H}; R(h) \ge \epsilon} e^{-n\epsilon}$$

$$< |\mathcal{H}| e^{-n\epsilon}$$

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

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

Therefore, we have:

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

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

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

We can improve the bound in proposition 4.3 as followed:

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

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:

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

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 $X_1^y, X_2^y, \ldots, X_{n_y}^y \sim P_y$ 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:

$$\widehat{h_n} = \arg\min_{h \in \mathcal{H}} \widehat{R_1}(h)$$
s.t.
$$\widehat{R_0}(h) \le \alpha + \frac{\epsilon}{2}$$

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

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

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:
$$\forall y \in \{0,1\}, \epsilon > 0: P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R_y}(h) - R_y(h) \right| \ge \epsilon \right) \le 2|\mathcal{H}|e^{-2n\epsilon^2}$$

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

$$P\left(\left|\widehat{R_y}(h) - R_y(h)\right| \ge \epsilon\right) = P\left(\left|n\widehat{R_y}(h) - nR_y(h)\right| \ge n\epsilon\right)$$

$$= P\left(\left|n\widehat{R_y}(h) - \mathbb{E}\left[n\widehat{R_y}(h)\right]\right| \ge n\epsilon\right)$$

$$\le 2\exp\left(-\frac{2n^2\epsilon^2}{n}\right) = 2e^{-2n\epsilon^2} \quad (Hoeffding's inequality)$$

From the above, we have:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R}_{y}(h)-R_{y}(h)\right|\geq\epsilon\right)=P\left(\bigcup_{h\in\mathcal{H}}\left\{\left|\widehat{R}_{y}(h)-R_{y}(h)\right|\geq\epsilon\right\}\right)$$

$$\leq\sum_{h\in\mathcal{H}}P\left(\left|\widehat{R}_{y}(h)-R_{y}(h)\right|\geq\epsilon\right)$$

$$\leq\sum_{h\in\mathcal{H}}2e^{-2n\epsilon^{2}}=2|\mathcal{H}|e^{-2n\epsilon^{2}}$$

From the assumption, we have:

$$\widehat{R_0}(\widehat{h_n}) \le \alpha + \frac{\epsilon}{2}$$

Hence, we have:

$$R_{0}(\widehat{h_{n}}) = \widehat{R_{0}}(\widehat{h_{n}}) + R_{0}(\widehat{h_{n}}) - \widehat{R_{0}}(\widehat{h_{n}})$$

$$\leq \alpha + \frac{\epsilon}{2} + \left| R_{0}(\widehat{h_{n}}) - \widehat{R_{0}}(\widehat{h_{n}}) \right|$$

$$\leq \alpha + \frac{\epsilon}{2} + \sup_{h \in \mathcal{U}} \left| R_{0}(\widehat{h_{n}}) - \widehat{R_{0}}(\widehat{h_{n}}) \right|$$

From Claim 1, we know that:

$$P\left(\sup_{h \in \mathcal{H}} \left| R_0(h) - \widehat{R_0}(h) \right| \ge \frac{\epsilon}{2} \right) \le 2|\mathcal{H}|e^{-n\epsilon^2/2}$$

$$\implies P\left(\sup_{h \in \mathcal{H}} \left| R_0(h) - \widehat{R_0}(h) \right| \le \frac{\epsilon}{2} \right) \ge 1 - 2|\mathcal{H}|e^{-n\epsilon^2/2}$$

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

$$R_0(\widehat{h_n}) \le \alpha + \frac{\epsilon}{2} + \frac{\epsilon}{2} = \alpha + \epsilon$$

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

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

Let $h' \in \mathcal{H}$ be any function such that $\widehat{R_0}(h') \le \alpha + \frac{\epsilon}{2}$. We have:

$$R_{1}(\widehat{h_{n}}) - R_{\mathcal{H},1}^{*} = R_{1}(\widehat{h_{n}}) - \widehat{R_{1}}(\widehat{h_{n}}) + \widehat{R_{1}}(\widehat{h_{n}}) - \widehat{R_{1}}(h') + \widehat{R_{1}}(h') - R_{\mathcal{H},1}^{*}$$

$$\leq \left| R_{1}(\widehat{h_{n}}) - \widehat{R_{1}}(\widehat{h_{n}}) \right| + \underbrace{\widehat{R_{1}}(\widehat{h_{n}}) - \widehat{R_{1}}(h')}_{\leq 0} + \left| \widehat{R_{1}}(h') - R_{\mathcal{H},1}^{*} \right|$$

$$\leq \left| R_{1}(\widehat{h_{n}}) - \widehat{R_{1}}(\widehat{h_{n}}) \right| + \left| \widehat{R_{1}}(h') - R_{\mathcal{H},1}^{*} \right|$$

$$\leq 2 \sup_{h \in \mathcal{H}} \left| \widehat{R_{1}}(h) - R_{1}(h) \right|$$

From Claim 2, we have:

$$P(R_{1}(\widehat{h_{n}}) - R_{\mathcal{H},1}^{*} \ge \epsilon) \le P(2 \sup_{h \in \mathcal{H}} \left| \widehat{R_{1}}(h) - R_{1}(h) \right| \ge \epsilon)$$

$$= P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R_{1}}(h) - R_{1}(h) \right| \ge \frac{\epsilon}{2} \right)$$

$$\le 2|\mathcal{H}|e^{-n\epsilon^{2}/2} \quad (From \ Claim \ 1)$$

$$\implies P(R_{1}(\widehat{h_{n}}) - R_{\mathcal{H},1}^{*} \le \epsilon) \ge 1 - 2|\mathcal{H}|e^{-n\epsilon^{2}/2}$$

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.

5.1VC 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^{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:

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

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

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

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

$$\implies S_{\mathcal{H}}(n) \le |\mathcal{H}|$$

$$\implies V_{\mathcal{H}} \le \log_2 |\mathcal{H}|$$

Remark: To show that $V_{\mathcal{H}} = n$, we must show that there exists at least n points x_1, \ldots, 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:

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

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

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

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:

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

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 $C=\{c_1,\ldots,c_m\}$ and $C'=\{c_2,\ldots,c_m\}$. Define the following label sets:

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

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

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

$$|Y_0| = |\mathcal{H}_{C'}|$$

$$\leq \left| \left\{ B \subset C' : \mathcal{H} \text{ shatters } B \right\} \right|$$

$$= \left| \left\{ B \subset C : c_1 \notin B, \mathcal{H} \text{ shatters } B \right\} \right|$$

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:

$$|Y_{1}| = |\mathcal{H}'_{C'}|$$

$$\leq \left| \left\{ B \subset C' : \mathcal{H}' \text{ shatters } B \right\} \right|$$

$$= \left| \left\{ B \subset C' : \mathcal{H}' \text{ shatters } B \cup \{c_{1}\} \right\} \right|$$

$$= \left| \left\{ B \subset C : c_{1} \in B, \mathcal{H}' \text{ shatters } B \right\} \right|$$

$$\leq \left| \left\{ B \subset C : c_{1} \in B, \mathcal{H} \text{ shatters } B \right\} \right|$$

From the above, we have:

$$|\mathcal{H}_C| = |Y_0| + |Y_1|$$

$$\leq \left| \left\{ B \subset C : c_1 \notin B, \mathcal{H} \text{ shatters } B \right\} \right| + \left| \left\{ B \subset C : c_1 \in B, \mathcal{H} \text{ shatters } B \right\} \right|$$

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

$$\leq \sum_{i=1}^{d} \binom{n}{k}$$

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

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

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

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

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

Proof (Corollary 5.1).

By Binomial theorem, we have:

$$(n+1)^{d} = \sum_{k=1}^{d} n^{k} \binom{d}{k}$$

$$= \sum_{k=1}^{d} n^{k} \frac{d!}{k!(d-k)!}$$

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

 \Box .

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

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

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

Proof (Corollary 5.2).

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

$$\left(\frac{d}{n}\right)^{d} \sum_{k=0}^{d} \binom{n}{k} \leq \sum_{k=0}^{d} \left(\frac{d}{n}\right)^{k} \binom{n}{k}$$

$$\leq \sum_{k=0}^{n} \left(\frac{d}{n}\right)^{k} \binom{n}{k}$$

$$= \left(1 + \frac{d}{n}\right)^{n} \leq e^{d}$$

Hence, we have:

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

 \Box .

 \Box .

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:

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

Proof (Corollary 5.3).

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

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

5.3 VC Theorem for classifiers

Theorem 5.2: VC Theorem (for classifiers)

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

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

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. \Box .

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

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

$$P(R(\widehat{h_n}) - R_{\mathcal{H}}^* \ge \epsilon) \le 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128}$$

Proof (Corollary 5.4).

We have:

$$P(\widehat{R(h_n)} - R_{\mathcal{H}}^* \ge \epsilon) \le P\left(2 \sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \epsilon\right)$$

$$= P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \frac{\epsilon}{2}\right)$$

$$\le 8S_{\mathcal{H}}(n)e^{-n(\epsilon/2)^2/32} = 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128}$$

 \Box .

Corollary 5.5: Excess Risk of $\widehat{h_n}$ - $\delta \to \epsilon$ 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:

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

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:

$$R(\widehat{h_n}) \le R_{\mathcal{H}}^* + \epsilon$$

$$= R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left(\log S_{\mathcal{H}}(n) - \log \frac{\delta}{8}\right)}$$

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:

$$R(\widehat{h_n}) \le R_{\mathcal{H}}^* + \sqrt{\frac{128}{n} \left(\log S_{\mathcal{H}}(n) - \log \frac{\delta}{8} \right)}$$

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

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 $d \ge 1$, 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:

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

$$\Longrightarrow S_{\mathcal{H}}(n) \le |\mathcal{H}|$$

$$\Longrightarrow V_{\mathcal{H}} \le \log_2 |\mathcal{H}|$$



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:

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

Then, the following set of classifiers:

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

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

Proof (Proposition 5.1).

Suppose that \mathcal{H} shatters m+1 points in $C=(x_1,\ldots,x_{m+1})$. Define the linear mapping $L_C:\mathcal{F}\to\mathbb{R}^{m+1}$ 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 $l_n \to l$ as $n \to \infty$. 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 $l_n \to l$, 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,

$$\sum_{i=1}^{m+1} \gamma_i f(x_0) = 0 \implies \sum_{i, \gamma_i \ge 0} \gamma_i f(x_i) = \sum_{j, \gamma_j < 0} -\gamma_j f(x_j)$$

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

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

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 $V_{\mathcal{H}} < \infty$.

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

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

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

Then, define the set of linear classifiers as:

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

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

 \mathbf{Proof}

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

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 $C = \{x_1, \ldots, x_n\}$. We have the following definitions for restriction of \mathcal{G} to C, shattering coefficient and VC-dimension of \mathcal{G} :

$$N_{\mathcal{G}}(C) = \left| \left\{ G \cap C : G \in \mathcal{G} \right\} \right|$$

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

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

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

$$G \to h_G(x) = \mathbf{1}_{\{x \in G\}}$$

 $h \to G_h = \left\{ x : h(x) = 1 \right\}$

Theorem 5.3: VC Theorem (for sets)

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

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

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

$$\widehat{Q}(G) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{\{X_1 \in G\}}$$

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:

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

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:

$$P\left(\sup_{G \in \mathcal{G}} \left| \widehat{Q}(G) - Q(G) \right| \ge \epsilon \right) = P\left(\sup_{h_G \in \mathcal{H}} \left| \widehat{R}_n(h_G) - R(h_G) \right| \ge \epsilon \right)$$

$$\le 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/32} \quad (Theorem 5.2)$$

$$= 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32}$$

 \Box .

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,

$$Q(G_t) = P(X \le t) = F(t) \quad (CDF)$$

$$\widehat{Q}(G_t) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{\{x_i \le t\}} = \widehat{F_n}(t) \quad (Empirical CDF)$$

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

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

Proof (Corollary 5.7). Let $\mathcal{G} = \{G_t : t \in \mathbb{R}\}$, by theorem 5.3, we have:

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

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:

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

 \Box .

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:

$$P(\widehat{R(h_n)} - R_{\mathcal{H}}^* \ge \epsilon) \le P\left(2 \sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \epsilon\right)$$

$$= P\left(\sup_{h \in \mathcal{H}} \left| \widehat{R_n}(h) - R(h) \right| \ge \frac{\epsilon}{2}\right)$$

$$\le 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/128} \quad (Theorem 5.2)$$

$$< 8(n+1)^{V_{\mathcal{H}}}e^{-n\epsilon^2/128} \quad (Corollary 5.1)$$

Now let:

$$\delta = 8(n+1)^{V_{\mathcal{H}}} e^{-n\epsilon^2/128}$$

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

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

 \Box .

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:

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

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

$$\widehat{Q}(G) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{\{(X_i, Y_i) \in G_h\}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{\{h(X_1) \neq Y_i\}} = \widehat{R}_n(h)$$

From the above, we have:

$$P\left(\sup_{h\in\mathcal{H}}\left|\widehat{R_n}(h) - R(h)\right| \ge \epsilon\right) = P\left(\sup_{G\in\mathcal{G}}\left|\widehat{Q}(G) - Q(G)\right| \ge \epsilon\right)$$

$$\le 8S_{\mathcal{G}}(n)e^{-n\epsilon^2/32}$$

$$= 8S_{\mathcal{H}}(n)e^{-n\epsilon^2/32}$$

 \Box .

Exercise 5.3

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

• (a)
$$\mathcal{G}_1 \cap \mathcal{G}_2 = \{G_1 \cap G_2 : G_1 \in \mathcal{G}_1, G_2 \in \mathcal{G}_2\}.$$

• (b)
$$\mathcal{G}_1 \cup \mathcal{G}_2 = \{G_1 \cup G_2 : G_1 \in \mathcal{G}_1, G_2 \in \mathcal{G}_2\}$$

Show that $S_{\mathcal{G}_1 \cap \mathcal{G}_2}(n) \leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$ and $S_{\mathcal{G}_1 \cup \mathcal{G}_2}(n) \leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$.

Solution (Exercise 5.3).

Proving each inequality one by one, we have:

Claim: $S_{\mathcal{G}_1 \cap \mathcal{G}_2}(n) \leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$ For any $\{x_1, \ldots, x_n\} \subset \mathcal{X}$, denote the following set:

$$\mathcal{F} = \left\{ G_1 \cap \{x_1, \dots, x_n\} : G_1 \in \mathcal{G}_1 \right\}$$

Then, \mathcal{F} is a collection of subsets of $\{x_1,\ldots,x_n\}$. Furthermore, the cardinality of \mathcal{F} is at most $S_{\mathcal{G}_1}(n)$. Now define the restriction of $\mathcal{G}_1 \cap \mathcal{G}_2$ to $\{x_1,\ldots,x_n\}$:

$$\mathcal{G}_1 \cap \mathcal{G}_{2\{x_1,\dots,x_n\}} = \left\{ G_1 \cap G_2 \cap \{x_1,\dots,x_n\} : G_1 \in \mathcal{G}_1, G_2 \in \mathcal{G}_2 \right\}$$
$$= \bigcup_{F \in \mathcal{F}} \left\{ G_2 \cap F : G_2 \in \mathcal{G}_2 \right\}$$

For each $F \in \mathcal{F}$, we have $|F| \leq n$. Hence, we have:

$$\left|\left\{G_2 \cap F : G_2 \in \mathcal{G}_2\right\}\right| \le S_{\mathcal{G}_2}(n), \ \forall F \in \mathcal{F}$$

Hence,

$$\left| \mathcal{G}_1 \cap \mathcal{G}_{2} \{ x_1, \dots, x_n \} \right| = \left| \bigcup_{F \in \mathcal{F}} \left\{ G_2 \cap F : G_2 \in \mathcal{G}_2 \right\} \right|$$

$$\leq \sum_{F \in \mathcal{F}} \left| \left\{ G_2 \cap F : G_2 \in \mathcal{G}_2 \right\} \right|$$

$$\leq \sum_{F \in \mathcal{F}} S_{\mathcal{G}_2}(n) = |\mathcal{F}| S_{\mathcal{G}_2}(n)$$

$$\leq S_{\mathcal{G}_1}(n) S_{\mathcal{G}_2}(n)$$

Since the above is a uniform bound, we can take the supremum over all $\{x_1, \ldots, x_n\} \subset \mathcal{X}$ and the inequality still holds. Hence,

$$\sup_{\{x_1,\ldots,x_n\}\subset\mathcal{X}}\left|\mathcal{G}_1\cap\mathcal{G}_{2}\{x_1,\ldots,x_n\}\right|=S_{\mathcal{G}_1\cap\mathcal{G}_2}(n)\leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$$

Claim: $S_{\mathcal{G}_1 \cup \mathcal{G}_2}(n) \leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$ For any $\{x_1, \ldots, x_n\} \subset \mathcal{X}$, we define the following collections of subsets:

$$\mathcal{F}_{1} = \left\{ G_{1} \cap \{x_{1}, \dots, x_{n}\} : G_{1} \in \mathcal{G}_{1} \right\}$$
$$\mathcal{F}_{2} = \left\{ G_{2} \cap \{x_{2}, \dots, x_{n}\} : G_{2} \in \mathcal{G}_{2} \right\}$$

Then we have:

$$\mathcal{G}_1 \cup \mathcal{G}_{2\{x_1,\dots,x_n\}} = \bigcup_{F_1 \in \mathcal{F}_1} \bigcup_{F_2 \in \mathcal{F}_2} \{F_1 \cup F_2\}$$

Since $|\mathcal{F}_1| \leq S_{\mathcal{G}_1}(n)$ and $|\mathcal{F}_2| \leq S_{\mathcal{G}_2}(n)$, we have:

$$\left| \mathcal{G}_1 \cup \mathcal{G}_{2\{x_1,\dots,x_n\}} \right| = \left| \bigcup_{F_1 \in \mathcal{F}_1} \bigcup_{F_2 \in \mathcal{F}_2} \{F_1 \cup F_2\} \right|$$

$$\leq |\mathcal{F}_1| |\mathcal{F}_2|$$

$$\leq S_{\mathcal{G}_1}(n) S_{\mathcal{G}_2}(n)$$

Taking the supremum over $\{x_1, \ldots, x_n\} \subset \mathcal{X}$ for both sides, we have:

$$\sup_{\{x_1,\dots,x_n\}\subset\mathcal{X}}\left|\mathcal{G}_1\cup\mathcal{G}_{2\{x_1,\dots,x_n\}}\right|=S_{\mathcal{G}_1\cup\mathcal{G}_2}(n)\leq S_{\mathcal{G}_1}(n)S_{\mathcal{G}_2}(n)$$

Exercise 5.4

Show that the following classes have finite VC dimension by exhibiting an explicit upperbound on the VC dimension.

- (a) $\mathcal{X} = \mathbb{R}^d$, $\mathcal{H} = \{\mathbf{1}_{\{f(x) \geq 0\}} : f \text{ inhomogeneous quadratic polynomial}\}.$
- (b) $\mathcal{X} = \mathbb{R}^d$, $\mathcal{H} = \{\mathbf{1}_{\{x \in C\}} : C \text{ is a closed sphere } \}$.
- (c) $\mathcal{X} = \mathbb{R}^2$, $\mathcal{H} = \{\mathbf{1}_{\{x \in P_k\}} : P_k \text{ is a convex polygon of at most } k \text{ sides}\}.$
- (d) $\mathcal{X} = \mathbb{R}^d$, $\mathcal{H} = \{\mathbf{1}_{\{x \in R_k\}} : R_k \text{ is a union of at most } k \text{ rectangles}\}.$

Solution (Exercise 5.4).

 \Box .

 \Box .

6 Rademacher Complexity

6.1 Bounded Difference Inequality

In the following section, we will discuss another concentration inequality that bounds the difference between functions of random sample and their mean given that the functions satisfy the **bounded** difference property.

Definition 6.1 (Bounded difference property).

Given a real-valued function $\phi: \mathcal{X}^n \to \mathbb{R}$. We say that ϕ satisfies the **bounded difference** property if $\exists c_1, \ldots, c_n \in \mathcal{X}$ such that $\forall 1 \leq i \leq n$:

$$\sup_{\{x_1,\ldots,x_n\}\subset\mathcal{X},x_i'\in\mathcal{X}}\left|\phi(x_1,\ldots,x_i,\ldots,x_n)-\phi(x_1,\ldots,x_i',\ldots,x_n)\right|\leq c_i$$

That is, substituting the value at the i^{th} coordinate x_i changes the value of ϕ by at most c_i .

Theorem 6.1: Bounded Difference (McDiarmid's) Inequality

Let X_1, \ldots, X_n be independent random variables (not necessarily identically distributed) and $\phi: \mathcal{X}^n \to \mathbb{R}$ be a function satisfying the bounded difference property:

$$\sup_{\{x_1,\ldots,x_n\}\subset\mathcal{X},x_i'\in\mathcal{X}}\left|\phi(x_1,\ldots,x_i,\ldots,x_n)-\phi(x_1,\ldots,x_i',\ldots,x_n)\right|\leq c_i,\ \forall 1\leq i\leq n$$

Then, we have:

$$P\left(\left|\phi(X_1,\ldots,X_n) - \mathbb{E}\left[\phi(X_1,\ldots,X_n)\right]\right| \ge t\right) \le 2\exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right), \ \forall t > 0$$

Remark: Assume that $X_i \in [a_i, b_i]$ and $\phi(X_1, \dots, X_n) = \sum_{i=1}^n X_i$. Then the bounded difference inequality recovers the Hoeffding's inequality 3.1.

Proof (Theorem 6.1). $_$

Define the following random variable:

$$V_i = \mathbb{E}\left[\phi \middle| X_1, \dots, X_i\right] - \mathbb{E}\left[\phi \middle| X_1, \dots, X_{i-1}\right]$$

Denote $\phi(X_1,\ldots,X_n)=\phi$ and $\mathbb{E}[\phi(X_1,\ldots,X_n)]=\mu_{\phi}$ for brevity, we have:

$$\phi - \mu_{\phi} = \sum_{i=1}^{n} V_i$$

Using the Chernoff's bounding method, we have:

$$P(\phi - \mu_{\phi} \ge t) \le \inf_{s>0} e^{-st} M_{\phi - \mu_{\pi}}(s)$$
$$= \inf_{s>0} e^{-st} \mathbb{E} \left[\exp \left(s \sum_{i=1}^{n} V_i \right) \right]$$

Claim 1: For all $1 \le i \le n$, $a_i \le V_i \le b_i$ and $b_i - a_i \le c_i$. Define the infimum and supremum of V_i as followed:

$$U_{i} = \sup_{x \in \mathcal{X}} \left\{ \mathbb{E} \left[\phi \middle| X_{1}, \dots, X_{i} = x \right] - \mathbb{E} \left[\phi \middle| X_{1}, \dots, X_{i-1} \right] \right\}$$
$$L_{i} = \inf_{x \in \mathcal{X}} \left\{ \mathbb{E} \left[\phi \middle| X_{1}, \dots, X_{i} = x \right] - \mathbb{E} \left[\phi \middle| X_{1}, \dots, X_{i-1} \right] \right\}$$

Clearly, $U_i \geq V_i \geq L_i$. We have:

$$\begin{split} U_i - L_i &= \sup_{x \in \mathcal{X}} \mathbb{E} \Big[\phi \Big| X_1, \dots, X_{i-1}, X_i = x \Big] - \inf_{x \in \mathcal{X}} \mathbb{E} \Big[\phi \Big| X_1, \dots, X_{i-1}, X_i = x \Big] \\ &= \sup_{x \in \mathcal{X}} \int \phi(X_1, \dots, X_{i-1}, x, x_{i+1}, \dots, x_n) dP(x_{i+1}, \dots, x_n | X_1, \dots, X_{i-1}, x) \\ &- \inf_{x \in \mathcal{X}} \int \phi(X_1, \dots, X_{i-1}, x, x_{i+1}, \dots, x_n) dP(x_{i+1}, \dots, x_n | X_1, \dots, X_{i-1}, x) \\ &= \sup_{x \in \mathcal{X}} \int \phi(X_1, \dots, X_{i-1}, x, x_{i+1}, \dots, x_n) dP(x_{i+1}, \dots, x_n) \\ &- \inf_{x \in \mathcal{X}} \int \phi(X_1, \dots, X_{i-1}, x, x_{i+1}, \dots, x_n) dP(x_{i+1}, \dots, x_n) \\ &= \sup_{x, y \in \mathcal{X}} \int \Big[\phi(X_1, \dots, X_{i-1}, x, x_{i+1}, \dots, x_n) - \phi(X_1, \dots, X_{i-1}, y, x_{i+1}, \dots, x_n) \Big] dP(x_{i+1}, \dots, x_n) \\ &\leq c_i \int dP(x_{i+1}, \dots, x_n) = c_i \end{split}$$

Claim 2: $\mathbb{E}\Big[V_i\Big|X_1,\ldots,X_{i-1}\Big]=0, \ \forall 1\leq i\leq n$. We have:

$$\mathbb{E}\left[V_{i}\middle|X_{1},\ldots,X_{i-1}\right] = \mathbb{E}\left[\mathbb{E}\left[\phi\middle|X_{1},\ldots,X_{i}\right]\middle|X_{1},\ldots,X_{i-1}\right] - \mathbb{E}\left[\mathbb{E}\left[\phi\middle|X_{1},\ldots,X_{i-1}\right]\middle|X_{1},\ldots,X_{i-1}\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\phi\middle|X_{1},\ldots,X_{i}\right]\middle|X_{1},\ldots,X_{i-1}\right] - \mathbb{E}\left[\phi\middle|X_{1},\ldots,X_{i-1}\right]$$

By the tower property, we have:

$$\mathbb{E}\left[\mathbb{E}\left[\phi\Big|X_1,\ldots,X_i\right]\bigg|X_1,\ldots,X_{i-1}\right] = \mathbb{E}\left[\phi\Big|X_1,\ldots,X_{i-1}\right]$$

Hence,

$$\mathbb{E}\left[V_i\middle|X_1,\ldots,X_{i-1}\right] = \mathbb{E}\left[\phi\middle|X_1,\ldots,X_{i-1}\right] - \mathbb{E}\left[\phi\middle|X_1,\ldots,X_{i-1}\right]$$

$$= 0$$

From the above claims, we can make use of the Hoeffding's lemma 3.1 to bound the moment

generating functions. We have:

$$P(\phi - \mu_{\phi} \ge t) = \inf_{s>0} e^{-st} \mathbb{E} \left[\exp\left(s \sum_{i=1}^{n} V_{i}\right) \right]$$

$$= \inf_{s>0} e^{-st} \mathbb{E}_{X_{1},\dots,X_{n-1}} \mathbb{E}_{X_{n}|X_{1},\dots,X_{n-1}} \left[\exp\left(s \sum_{i=1}^{n} V_{i}\right) \right]$$

$$= \inf_{s>0} e^{-st} \mathbb{E}_{X_{n}|X_{1},\dots,X_{n-1}} \left[e^{sV_{n}} \right] \mathbb{E}_{X_{1},\dots,X_{n-1}} \left[\exp\left(s \sum_{i=1}^{n-1} V_{i}\right) \right]$$

$$\leq \inf_{s>0} \exp\left(-st + \frac{s^{2}c_{n}^{2}}{8}\right) \mathbb{E}_{X_{1},\dots,X_{n-1}} \left[\exp\left(s \sum_{i=1}^{n-1} V_{i}\right) \right] \quad (Lemma \ 3.1)$$

$$\vdots$$

$$\leq \inf_{s>0} \exp\left(-st + s^{2} \sum_{i=1}^{n} \frac{c_{i}^{2}}{8}\right)$$

Substituting $s = \frac{4t}{\sum_{i=1}^{n} c_i^2}$ to minimize the upperbound (just like the proof for Hoeffding's inequality 3.1). We have:

$$P(\phi - \mu_{\phi} \ge t) \le \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right)$$

 \Box .

6.2 Rademacher Complexity

Overview: Rademacher Complexity is a measure for the richness of a class of real-valued functions. In this sense, it is similar to VC dimension. However, unlike VC dimension, the Rademacher Complexity is not restricted to binary functions.

Definition 6.2 (Empirical Rademacher Complexity). Let $\mathcal{G} \subseteq [a,b]^{\mathcal{Z}}$ be a set of functions $\mathcal{Z} \to [a,b]$ where $a,b \in \mathbb{R}, a < b$. Let Z_1,\ldots,Z_n be an independently identically distributed random sample on \mathcal{Z} following some distribution P. Denote $S = (Z_1,\ldots,Z_n)$, we define the **Empirical Rademacher Complexity** as:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{G}) = \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(Z_{i}) \right]$$

Where $\sigma = (\sigma_1, \dots, \sigma_n)^T$, $\sigma_i \sim Uniform(-1, 1)$ are known as **Rademacher random variables**. Note that $\widehat{\mathfrak{R}}_S(\mathcal{G})$ is random due to randomness in S.

$$\mathfrak{R}_n(\mathcal{G}) = \mathbb{E}_S \left[\widehat{\mathfrak{R}}_S(\mathcal{G}) \right]$$

Theorem 6.2: One-sided Rademacher Complexity bound

Let Z be a random variable and $S=(Z_1,\ldots,Z_n)$ be an independently identically distributed sample over \mathcal{Z} . Consider a class of functions $\mathcal{G}\subseteq[a,b]^{\mathcal{Z}}$. $\forall \delta>0,g\in\mathcal{G}$, with at least probability $1-\delta$, we have:

(i)
$$\mathbb{E}\left[g(Z)\right] - \frac{1}{n} \sum_{i=1}^{n} g(Z_i) \le 2\mathfrak{R}_n(\mathcal{G}) + (b-a)\sqrt{\frac{\log 1/\delta}{2n}}$$

(ii)
$$\mathbb{E}\left[g(Z)\right] - \frac{1}{n} \sum_{i=1}^{n} g(Z_i) \le 2\widehat{\Re}_S(\mathcal{G}) + 3(b-a)\sqrt{\frac{\log 2/\delta}{2n}}$$

Proof (Theorem 6.2).

 \Box .

Theorem 6.3: Two-sided Rademacher Complexity bound

Consider a set of classifiers $\mathcal{G} \subseteq [a, b]^{\mathcal{Z}}$. Then, $\forall \delta > 0$, with probability of at least $1 - \delta$, we have:

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E} \left[g(Z) \right] - \frac{1}{n} \sum_{i=1}^{n} g(Z_i) \right| \le 2\mathfrak{R}_n(\mathcal{G}) + (b - a) \sqrt{\frac{\log 2/\delta}{2n}}$$

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E} \left[g(Z) \right] - \frac{1}{n} \sum_{i=1}^{n} g(Z_i) \right| \le 2 \widehat{\mathfrak{R}}_n(\mathcal{G}) + 3(b-a) \sqrt{\frac{\log 4/\delta}{2n}}$$

Proof (Theorem 6.3).

(The proof of this theorem is left as an exercise).

 \Box .

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 $\overline{\mathcal{R}}$ 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:

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

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:

maximize:
$$P_{H_1}(\mathcal{R}) = \int_{\mathcal{R}} f_1(x) dx$$

subjected to: $P_{H_0}(\mathcal{R}) = \int_{\mathcal{R}} f_0(x) dx \leq \alpha$

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 α , 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:

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

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 α , we have:

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

$$\implies P_{H_1}(\mathcal{R}) - P_{H_1}(\mathcal{R}') \ge c \left(\int_{\mathcal{R}} f_0(x) dx - \int_{\mathcal{R}'} f_0(x) dx \right)$$

$$= c \left(\alpha - \int_{\mathcal{R}'} f_0(x) dx \right)$$

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

$$P_{H_1}(\mathcal{R}) - P_{H_1}(\mathcal{R}') \ge 0 \implies P_{H_1}(\mathcal{R}') \le P_{H_1}(\mathcal{R})$$

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

A.2 Rademacher Complexity bound for linear function classes

In the following section, we will explore a simple exercise for bounding the Rademacher Complexity for linear function classes.

A.2.1 Problem Statement

Problem: Let \mathcal{F} be a linear function class defined as followed:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \to \mathbb{R} \middle| f(x) = wx, ||w||_2 \le R \right\}$$

Our objective is to prove the following Rademacher Complexity bound:

$$\mathfrak{R}_n(\mathcal{F}) \leq \tilde{O}\left(\frac{R}{\sqrt{n}}\right)$$

Before solving the above problem, we have to get familiar with the definition of **covering number** and some related lemmas.

A.2.2 Covering Number

Definition A.3 (ϵ -Cover). $_$

Let Q be a set. A subset $C \subset Q$ is called an ϵ -cover of Q with respect to a metric ρ if:

$$\forall v \in Q, \exists v' \in \mathcal{C} : \rho(v, v') \le \epsilon$$

Basically, C can be thought of as a collection of centers of ϵ -balls overlapping Q (Figure 4).

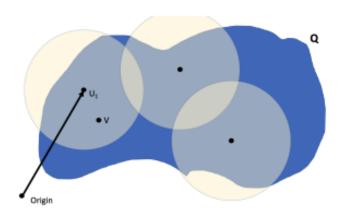


Figure 4: Examples of ϵ -balls. The centers of the balls would form an ϵ -cover if they completely contain Q.

Definition A.4 (Covering Number of sets $(\mathcal{N}(Q, \epsilon, \rho))$.

The covering number of a set Q is defined as the size of the smallest ϵ -cover needed to completely contain Q. In other words, it is the minimum number of ϵ -balls needed to completely contain Q:

$$\mathcal{N}(Q,\epsilon,\rho) = minimum \ size \ of \ \epsilon\text{-covers} \ of \ Q \ w.r.t \ \rho$$

Visual illustration of covering number is included in figure 5.



Figure 5: Example of covering number. For \mathcal{F} , the covering number is 5. For \mathcal{F}' , the covering number is 10.

Definition A.5 (Covering number of function class $(\mathcal{N}(\mathcal{F}, \epsilon, \rho))$). Let \mathcal{F} be a function class. Define the restriction of \mathcal{F} to the observations $\{x_1, \ldots, x_n\}$ as:

$$\mathcal{F}_{x_1,\ldots,x_n} = \left\{ (f(x_1),\ldots,f(x_n)) : f \in \mathcal{F} \right\}$$

Then, we can define the **Covering Number** for \mathcal{F} as followed:

$$\mathcal{N}(\mathcal{F}, \epsilon, \rho) = \sup_{x_1, \dots, x_n} \mathcal{N}(\mathcal{F}_{x_1, \dots, x_n}, \epsilon, \rho)$$

We can call $\mathcal{N}(\mathcal{F}_{x_1,...,x_n},\epsilon,\rho)$ the **Empirical Covering Number**.

A.2.3 Massart's Lemma

A.2.4 Dudley's Theorem

Theorem A.2: Dudley's Theorem

If \mathcal{F} is a function class from $\mathcal{Z} \to \mathbb{R}$ (where \mathcal{Z} is a vector space and the norm of $f \in \mathcal{F}$ is not necessarily bounded), then:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) \leq 12 \int_{0}^{\sup_{f \in \mathcal{F}} \|f\|_{2}} \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})}{n}} d\epsilon$$

Where for any $f \in \mathcal{F}$, given a sample x_1, \ldots, x_n , we define the $\|.\|_2$ norm as followed:

$$||f||_2 = \sqrt{\frac{1}{n} \sum_{i=1}^n f(x_i)^2}$$

Proof (Theorem A.2).

The main idea of the proof is via chaining. Define the following sequence $\{\epsilon_j\}_{j=0}^n$:

$$\begin{cases} \epsilon_0 &= \sup_{f \in \mathcal{F}} \|f\|_2 \\ \epsilon_j &= 2^{-j} \epsilon_0 \end{cases}$$

Next, we define the sequence of ϵ -covers \mathcal{N}_{ϵ_i} corresponding to each ϵ_i . By definition, we have:

$$\forall f \in \mathcal{F}, \exists g_j \in \mathcal{N}_{\epsilon_j} : ||f - g_j||_2 \le \epsilon_j$$

We can write any $f \in \mathcal{F}$ as the following telescoping sum:

$$f = f - g_n + \sum_{j=1}^{n} (g_j - g_{j-1})$$

Now, define the Empirical Rademacher Complexity as followed:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} f(x_{i}) \right] = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left\langle \sigma, f \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

We have:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left\langle \sigma, f \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

$$= \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left\langle \sigma, f - g_{n} + \sum_{j=1}^{n} \left(g_{j} - g_{j-1} \right) \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

$$\leq \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left\langle \sigma, f - g_{n} \right\rangle + \sup_{f \in \mathcal{F}} \sum_{j=1}^{n} \left\langle \sigma, g_{j} - g_{j-1} \right\rangle \middle| x_{1}, \dots, x_{n} \right] \quad (\sup \Sigma \leq \Sigma \sup)$$

$$\leq \|\sigma\|_{2} \cdot \|f - g_{n}\|_{2} + \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \sum_{j=1}^{n} \left\langle \sigma, g_{j} - g_{j-1} \right\rangle \middle| x_{1}, \dots, x_{n} \right] \quad (Cauchy-Schwarz)$$

$$\leq \epsilon_{n} + \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \sum_{j=1}^{n} \left\langle \sigma, g_{j} - g_{j-1} \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

$$\leq \epsilon_{n} + \sum_{j=1}^{n} \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left\langle \sigma, g_{j} - g_{j-1} \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

From here, note that:

$$||g_j - g_{j-1}||_2 \le ||g_j - f||_2 + ||g_{j-1} - f||_2 \le \epsilon_j + \epsilon_{j-1} = 3\epsilon_j$$

Also, define the following classes of functions:

$$\mathcal{H}_{j} = \left\{ g_{j} - g_{j-1} \middle| g_{j} \in \mathcal{N}_{\epsilon_{j}}, g_{j-1} \in \mathcal{N}_{\epsilon_{j-1}} \right\}$$

Continuing the above, we have:

$$\widehat{\Re}_{S}(\mathcal{F}) \leq \epsilon_{n} + \sum_{j=1}^{n} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}_{j}} \left\langle \sigma, h \right\rangle \middle| x_{1}, \dots, x_{n} \right]$$

$$\leq \epsilon_{n} + \sum_{j=1}^{n} \sup_{h \in \mathcal{H}_{j}} \|h\|_{2} \cdot \sqrt{\frac{2 \log |\mathcal{H}_{j}|}{n}} \quad (Massart's \ lemma)$$

$$\leq \epsilon_{n} + \sum_{j=1}^{n} 6(\epsilon_{j} - \epsilon_{j-1}) \cdot \sqrt{\frac{2 \log |\mathcal{H}_{j}|}{n}} \quad \left(\sup_{h \in \mathcal{H}_{j}} \|h\|_{2} \leq 3\epsilon_{j} \leq 6(\epsilon_{j} - \epsilon_{j+1}) \right)$$

Now we have:

$$|\mathcal{H}_j| \le |\mathcal{N}_{\epsilon_j}| \cdot |\mathcal{N}_{\epsilon_{j-1}}| \le |\mathcal{N}_{\epsilon_j}|^2$$

Hence, we have:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) \leq \epsilon_{n} + 6 \sum_{j=1}^{n} (\epsilon_{j} - \epsilon_{j-1}) \cdot \sqrt{\frac{2 \log |\mathcal{N}_{\epsilon_{j}}|^{2}}{n}}$$

$$= \epsilon_{n} + 12 \sum_{j=1}^{n} (\epsilon_{j} - \epsilon_{j-1}) \cdot \sqrt{\frac{\log |\mathcal{N}_{\epsilon_{j}}|}{n}}$$

$$\leq \epsilon_{n} + 12 \sum_{j=1}^{n} \int_{\epsilon_{j+1}}^{\epsilon_{j}} \sqrt{\frac{\log |\mathcal{N}_{\epsilon_{j}}|}{n}} dt$$

$$\leq \epsilon_{n} + 12 \sum_{j=1}^{n} \int_{\epsilon_{j+1}}^{\epsilon_{j}} \sqrt{\frac{\log |\mathcal{N}_{t}|}{n}} dt$$

$$\leq \epsilon_{n} + 12 \int_{\epsilon_{n+1}}^{\epsilon_{0}} \sqrt{\frac{\log |\mathcal{N}_{t}|}{n}} dt$$

Now take $n \to \infty$, we notice that $\epsilon_n \to 0$ and the above inequality holds for every ϵ -cover for $0 < \epsilon < \epsilon_0 = \sup_{f \in \mathcal{F}} \|f\|_2$. Therefore,

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) \leq 12 \int_{0}^{\sup_{f \in \mathcal{F}} \|f\|_{2}} \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})}{n}} d\epsilon$$

 \Box .

Remark: Using theorem A.2, we can translate the covering number to the Empirical Rademacher Complexity given that the covering number has some special formulation and the function class \mathcal{F} is bounded in [-1,1]. For example:

• (i) $\mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2) \approx (1/\epsilon)^R$.

Then we have $\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2) \approx R \log(1/\epsilon)$. Therefore,

$$\int_{0}^{1} \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})}{n}} d\epsilon = \int_{0}^{1} \sqrt{\frac{R \log(1/\epsilon)}{n}} d\epsilon \approx \sqrt{\frac{R}{n}}$$

• (ii) $\mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2) \approx a^{R/\epsilon}$. Then we have $\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2) \approx \frac{R}{\epsilon} \log a$. Therefore,

$$\int_{0}^{1} \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})}{n}} d\epsilon \approx \int_{0}^{1} \sqrt{\frac{R}{n\epsilon}} \log a d\epsilon$$

$$= \sqrt{\frac{R}{n}} \log a \int_{0}^{1} \sqrt{\frac{1}{\epsilon}} d\epsilon$$

$$= \tilde{O}\left(\sqrt{\frac{R}{n}}\right)$$

• (iii) $\mathcal{N}(\mathcal{F}, \epsilon, \|.\|^2) \approx a^{R/\epsilon^2}$. Then, we have $\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|^2) \approx \frac{R}{\epsilon^2} \log a$. Therefore,

$$\int_0^1 \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|^2)}{n}} d\epsilon \approx \sqrt{\frac{R}{n} \log a} \int_0^1 \frac{1}{\epsilon} d\epsilon = \infty$$

A.2.5 Bound on covering number of linear function class

Lemma A.1: Bound on $\mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2)$ for linear functions

Let $\mathcal F$ be a linear function class defined as followed:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \to \mathbb{R} \middle| f(x) = wx, ||w||_q \le a, ||x||_p \le b \right\}$$

Where p,q are Holder conjugates and $2 \le p \le \infty$. Then, we have:

$$\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_2) \le \left\lceil \frac{a^2 b^2}{\epsilon^2} \right\rceil \log(2d+1)$$

The proof of this lemma is discussed in Theorem 3 in [9].

Proof (Lemma A.1).

A.2.6 Rademacher Complexity bound for linear functions class

In this section, we finally solve the problem stated in section A.2.1. First, consider the following lemma, then prove the proposition A.1:

 \Box .

Theorem A.3: Dudley's Entropy Integral bound

Let \mathcal{F} be a real-valued function class and assume that $\mathbf{0} \in \mathcal{F}$. Then,

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) \leq \inf_{\alpha > 0} \left(4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^{\sup_{f \in \mathcal{F}} \|f\|_{2}} \sqrt{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})} d\epsilon \right)$$

Where we define the $\|.\|_2$ norm as followed:

$$||f||_2 = \sqrt{\frac{1}{n} \sum_{i=1}^n f(x_i)^2}, \ f \in \mathcal{F}$$

Proof (Theorem A.3).

Proposition A.1: Rademacher Complexity bound for linear function class

Given the following function class \mathcal{F} whose range is bounded within [-1,1]:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \to \mathbb{R} \middle| f(x) = wx, ||w||_2 \le a, ||x||_2 \le b \right\}$$

Then, we have the following bound for the Rademacher Complexity:

$$\mathfrak{R}_n(\mathcal{F}) \leq \tilde{O}\left(\frac{R}{\sqrt{n}}\right), \quad R = ab$$

Proof (Proposition A.1).

From lemma A.1, we have the following bound on the covering number of \mathcal{F} :

$$\log \mathcal{N}(\mathcal{F}, \epsilon, ||.||_2) \le \left\lceil \frac{R^2}{\epsilon^2} \right\rceil \log(2d+1) < 2\frac{R^2}{\epsilon^2} \log(2d+1) = \frac{R^2}{\epsilon^2} \log(4d^2 + 4d + 1)$$

The second inequality holds under the assumption that $R^2 > \epsilon^2$. Let $D = 4d^2 + 4d + 1$, we have:

$$\int_{\alpha}^{\sqrt{n}} \sqrt{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})} < R\sqrt{\log D} \int_{\alpha}^{\sqrt{n}} \frac{1}{\epsilon} d\epsilon$$
$$= R\sqrt{\log D} \left(\log \sqrt{n} - \log \alpha\right)$$

Using theorem A.3, we have:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) \leq \inf_{\alpha > 0} \left(4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^{\sqrt{n}} \sqrt{\log \mathcal{N}(\mathcal{F}, \epsilon, \|.\|_{2})} d\epsilon \right)$$
$$< \inf_{\alpha > 0} \left(4\alpha + \frac{12}{\sqrt{n}} R \sqrt{\log D} \left(\log \sqrt{n} - \log \alpha \right) \right)$$

Letting $\alpha = \frac{3R}{\sqrt{n}}$, we have:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{F}) < \frac{12R}{\sqrt{n}} + \frac{12R}{\sqrt{n}} \sqrt{\log D} \left(\log \sqrt{n} - \log \frac{3R}{\sqrt{n}} \right)$$

$$= \frac{12R}{\sqrt{n}} \left[1 + \sqrt{\log D} \log \frac{n}{3R} \right]$$

$$= \widetilde{O} \left(\frac{R}{\sqrt{n}} \right)$$

Since the right-hand-side does not depend on the sample S, we can just take expectation over the samples for both sides and we have:

$$\mathfrak{R}_n(\mathcal{F}) < \tilde{O}\left(\frac{R}{\sqrt{n}}\right)$$

 \Box .

 \Box .

B List of Definitions

| 1.1 1.2 1.3 1.4 1.5 1.6 1.7 2.1 3.1 4.1 4.2 4.3 5.1 | Definition (Classifier (h)) Definition (Decomposition of P_{XY}) Definition (Hypothesis space (\mathcal{H})) Definition (Learning algorithm (\mathcal{L}_n)) Definition (Risk $(R(h))$) Definition (Bayes Risk (R^*)) Definition (Consistency of learning algorithms) Definition (Plug-in classifier) Definition (Empirical Risk $(\widehat{R_n})$) Definition (Empirical Risk Minimization $(\widehat{h_n})$) Definition (Uniform Deviation Bounds (UDB)) Definition (PAC & Sample Complexity $(N(\epsilon, \delta))$) Definition (Restriction $(N_{\mathcal{H}})$) | 3 3 4 4 4 5 5 5 9 16 22 22 25 31 |
|---|--|--|
| 5.2 5.3 5.4 5.5 6.1 6.2 | Definition (Shattering Coefficient $(S_{\mathcal{H}})$) Definition (VC-dimension $(V_{\mathcal{H}})$) Definition (VC Class) Definition (VC Theory for sets) Definition (Bounded difference property) Definition (Empirical Rademacher Complexity) | 31 36 38 44 46 |
| A.4 | Definition (Rademacher Complexity) | 46 48 49 50 51 |
| | | |
| C I | Important Theorems | |
| 2.1 2.2 3.1 4.1 5.1 5.2 5.3 6.1 6.2 6.3 A.1 A.2 | Properties of Bayes classifier . Properties of Bayes classifier (Multi-class) Hoeffding's Inequality Uniform Deviation Bounds for finite #. Sauer's Lemma VC Theorem (for classifiers) VC Theorem (for sets) Bounded Difference (McDiarmid's) Inequality One-sided Rademacher Complexity bound Two-sided Rademacher Complexity bound Neyman-Pearson Lemma | 6 8 15 22 32 34 39 44 47 47 49 52 54 |
| 2.1 2.2 3.1 4.1 5.1 5.2 5.3 6.1 6.2 6.3 A.1 A.2 | Properties of Bayes classifier . Properties of Bayes classifier (Multi-class) Hoeffding's Inequality Uniform Deviation Bounds for finite # Sauer's Lemma VC Theorem (for classifiers) VC Theorem (for sets) Bounded Difference (McDiarmid's) Inequality One-sided Rademacher Complexity bound Two-sided Rademacher Complexity bound Neyman-Pearson Lemma Dudley's Theorem | 8 15 22 32 34 39 44 47 47 49 52 |

| 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 $\widehat{h_n}$ - $\delta \to \epsilon$ relation (VC-Theorem) | 35 |
| 5.6 | Linear classifiers have finite $V_{\mathcal{H}}$ | |
| 5.7 | Dvoretzky-Kiefer-Wolfowitz Inequality | 39 |
| | | |
| \mathbf{E} I | mportant Propositions | |
| | inportant repositions | |
| 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 |
| A.1 | Rademacher Complexity bound for linear function class | 55 |
| | | |

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.
- [4] Erhan undefinedinlar. Probability and Stochastics. Springer New York, 2011. ISBN: 9780387878591.
 DOI: 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.
- [9] Tong Zhang. "Covering Number Bounds of Certain Regularized Linear Function Classes". In: Journal of Machine Learning Research 2 (2002). URL: https://www.jmlr.org/papers/volume2/zhang02b/zhang02b.pdf.