

Theory of classification:

①

Basic Model :

- $x \in \mathcal{X}$, \mathcal{X} : a measurable space equipped with a σ -algebra.
- $y \in \{-1, 1\}$ a class.
- x is an observation.

• $g: \mathcal{X} \rightarrow \{-1, 1\}$: a classifier.

• (X, Y) : a random pair.

• $\eta(x) = P\{Y=1 | X=x\}$: a posteriori probability.

• measure perf of classifier g by its probability of error :

$$L(g) = P\{g(X) \neq Y\}$$

• Given η , classifier with minimal prob. of error is:

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ -1 & \text{otherwise} \end{cases}$$

then $L(g^*) \leq L(g) \quad \forall g$.

• The minimal risk $L^* \triangleq L(g^*)$ is called Bayes risk.

• We have: $L(g) - L^* = E \left[\mathbb{1}_{\{g(X) \neq g^*(X)\}} |2\eta(X) - 1| \right] \geq 0$.

• g^* can be called the Bayes classifier.

• data is denoted $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ i.i.d.

• A classifier constructed using D_n is g_n .

• $g_n(X) = g_n(X; X_1, Y_1, \dots, X_n, Y_n)$.

• Performance of g_n : $L(g_n) = P[g_n(X) \neq Y | D_n]$.

3. Empirical risk minimization and Rademacher averages.

(2)

1. Consider a class C of classifiers $g: X \rightarrow \{-1, 1\}$ and use data-based estimates of the probabilities of error $L(g)$ to select a classifier from the class.

2. The error count is a natural estimate of $L(g)$.

$$L_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(X_i) \neq Y_i\}}$$

$L_n(g)$ is called the empirical error of the classifier g .

3. Denote by g_n^* the classifier that minimizes the estimated probability of error over the class:

$$L_n(g_n^*) \leq L_n(g) \quad \forall g \in C$$

4. We have:

$$L(g_n^*) - \inf_{g \in C} L(g) \leq 2 \sup_{g \in C} |L_n(g) - L(g)|$$

$$L(g_n^*) \leq L_n(g_n^*) + \sup_{g \in C} |L_n(g) - L(g)|$$

5. $nL_n(g)$ is a r.v. binomially distributed with parameters n & $L(g)$.

6. Thus, to obtain bounds for the success of empirical error minimization, we need to study uniform deviations of binomial r.v. from their means.

7. Let X_1, \dots, X_n be i.i.d r.v. $X_i \in X$.

\mathcal{F} be a class of bounded fnc $X \rightarrow [-1, 1]$.

$Pf = E[f(X_1)]$ denotes expectation.

$P_n f = \frac{1}{n} \sum_{i=1}^n f(X_i)$ denotes empirical average.

we are interested in upper bound for the maximal deviation:

$$\sup_{f \in \mathcal{F}} (Pf - P_n f)$$

③
 Theorem 3.1 (bounded differences inequality). Let $g: \mathcal{X}^n \rightarrow \mathbb{R}$ be a fun. of n vars $| \exists c_i$

$$\sup_{\substack{x_1, \dots, x_n \\ x'_i \in \mathcal{X}}} |g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i, \quad 1 \leq i \leq n.$$

Let X_1, \dots, X_n be n independent r.v.

r.v. $Z = g(X_1, \dots, X_n)$ satisfies

$$P[|Z - E[Z]| > t] \leq 2e^{-\frac{t^2}{C}}$$

$$\text{where } C = \sum_{i=1}^n c_i^2$$

An example of a fun. that satisfies the bounded differences assumption is:

$$Z = \sup_{f \in \mathcal{F}} |Pf - P_n f|$$

Z satisfies the bounded diff. ass. with $c_i = \frac{2}{n}$.

With prob at least $1 - \delta$:

$$\sup_{f \in \mathcal{F}} |Pf - P_n f| \leq E \left[\sup_{f \in \mathcal{F}} |Pf - P_n f| \right] + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}.$$

Introduce a "ghost sample" X'_1, \dots, X'_n , independent of the X_i and inequality. distributed identically. If $P'_n f = \frac{1}{n} \sum_{i=1}^n f(X'_i)$, then by Jensen's ~~ineq~~

$$E \left[\sup_{f \in \mathcal{F}} |Pf - P_n f| \right] = E \left[\sup_{f \in \mathcal{F}} (E[|P'_n f - P_n f| | X_1, \dots, X_n]) \right] \\ \leq E \left[\sup_{f \in \mathcal{F}} |P'_n f - P_n f| \right].$$

It can be shown that:

$$E \left[\sup_f |P'_n f - P_n f| \right] \leq 2 E \left[\sup_f \frac{1}{n} \left| \sum_{i=1}^n \sigma_i f(X_i) \right| \right].$$

Let $A \in \mathbb{R}^n$ be a bounded set of vectors $a = (a_1, \dots, a_n)$ and

$$R_n(A) = E \left[\sup_{a \in A} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i a_i \right| \right]$$

• $P_n(A)$ is called the Rademacher average associated with A .

• for a given seq. $x_1, \dots, x_n \in X$, we write $F(x_i^n)$ for the class of n -vectors $(f(x_1), \dots, f(x_n))$ with $f \in F$.

$$\text{so: } F(x_i^n) = \{ (f(x_1), \dots, f(x_n)) : f \in F \}.$$

• With prob. at least $1 - \delta$:

$$\sup_f |Pf - P_n f| \leq 2E[R_n(F(x_i^n))] + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}.$$

We also have:

$$\sup_f |Pf - P_n f| \leq 2 \uparrow R_n(F(x_i^n)) + \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}.$$

this is a data-dependent bound.

• Properties of Rademacher averages. Let A, B be bounded subsets of \mathbb{R}^n &

let $c \in \mathbb{R}$ be a constant. Then:

$$R_n(A \cup B) \leq R_n(A) + R_n(B).$$

$$R_n(c \cdot A) = |c| R_n(A)$$

$$c \cdot A = \{ca : a \in A\}$$

$$R_n(A \oplus B) \leq R_n(A) + R_n(B). \quad A \oplus B = \{a+b : a \in A, b \in B\}$$

• Moreover, if $A = \{a^{(1)}, \dots, a^{(N)}\} \subset \mathbb{R}^n$ is a finite set, then

$$R_n(A) \leq \max_{j=1, \dots, N} \|a^{(j)}\| \frac{\sqrt{2 \log N}}{n}$$

$\|\cdot\|$ is Euclidean norm.

$$\text{If } \text{absconv}(A) = \left\{ \sum_{j=1}^N c_j a^{(j)} : N \in \mathbb{N}, \sum_{j=1}^N |c_j| \leq 1, a^{(j)} \in A \right\}$$

is the absolute convex hull of A , then:

$$R_n(A) = R_n(\text{absconv}(A)).$$

• The contraction principle states that if $\phi: \mathbb{R} \rightarrow \mathbb{R}$ is a fun. with $\phi(0) = 0$ and Lipschitz constant L_ϕ and $\phi \circ A$ is the set of vectors of form $(\phi(a_1), \dots, \phi(a_n)) \in \mathbb{R}^n$ with $a_i \in A$.

$$R_n(\phi \circ A) \leq L_\phi R_n(A)$$

• Consider the case when F is a class of indicator functions. (5)

~~Recall~~

For any collections of points $x_i^n = (x_1, \dots, x_n)$, $F(x_i^n)$ is a finite subset of R^n whose cardinality is denoted by $S_F(x_i^n)$ and is called the VC shatter coefficient.

• Obviously, $S_F(x_i^n) \leq 2^n$.

• We have, $\forall x_i^n$, $R_n(F(x_i^n)) \leq \sqrt{\frac{2 \log S_F(x_i^n)}{n}}$ (we used the fact that $\sum_i t(x_i)^2 \leq n$)

• In particular,

$$E \left[\sup_f |Pf - P_n f| \right] \leq 2 E \left[\sqrt{\frac{2 \log S_F(x_i^n)}{n}} \right]$$

• The log. of the VC shatter coefficient may be upper bounded in terms of a combinatorial quantity, called the VC dimension.

• If $A \subseteq \{-1, 1\}^n$, then the VC dimension of A is the size V of the largest set of indices $\{i_1, \dots, i_V\} \subset \{1, \dots, n\}$ | \forall for each binary V -vector $b = (b_1, \dots, b_V) \in \{-1, 1\}^V$, \exists an $a = (a_1, \dots, a_n) \in A$ | $(a_{i_1}, \dots, a_{i_V}) = b$.

• A key inequality establishing a rel. bit shatter coefficient & VC dimension is known as Sauer's lemma which states that the cardinality of any set $A \subseteq \{-1, 1\}^n$ may be upper bounded as

$$|A| \leq \sum_{i=0}^V \binom{n}{i} \leq (n+1)^V \quad \text{where } V \text{ is the VC dimension of } A.$$

• In particular, $\log S_F(x_i^n) \leq V(x_i^n) \log(n+1)$, where we denote by $V(x_i^n)$ the VC dimension of $F(x_i^n)$.

• Thus: $E \left[\sup_f [Pf - P_n f] \right] \leq 2 E \left[\sqrt{2 V(x_i^n) \frac{\log(n+1)}{n}} \right]$

- To obtain distribution-free upper bounds, introduce the VC dim. of a class of binary fns. \mathcal{F} , defined by:

$$V = \sup_{n, x_1^n} V(x_1^n).$$

- VC inequality. For all distributions, one has
- $$E \left[\sup_f (Pf - P_n f) \right] \leq 2 \sqrt{\frac{2V \log(n+1)}{n}}$$

also,

$$E \left[\sup_f (Pf - P_n f) \right] \leq C \sqrt{\frac{V}{n}}$$

for a universal constant C .

- One useful property: let G be an m -dimensional vector space of real-valued functions defined on X .

The class of indicator functions:

$$\mathcal{F} = \{f(x) = \mathbb{1}_{g(x) \geq 0} : g \in G\}$$

has VC dimension $V \leq m$.