

Binary Classification, Part 1

$$(X, Y) \quad \left. \begin{array}{l} X \in \mathcal{X} \text{ (feature)} \\ Y \in \{0, 1\} \text{ (label)} \end{array} \right\}$$

classifiers: $\mathcal{X} \rightarrow \{0, 1\}$
 $\mathbb{1}_{\{X \in C\}}$ - prediction

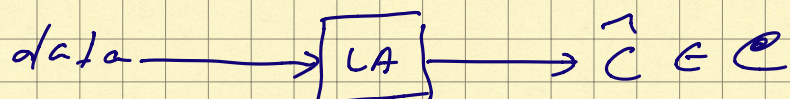
$$L(C) := \mathbb{P}[Y \neq \mathbb{1}_{\{X \in C\}}]$$

$P = \mathcal{L}((X, Y))$ known: Bayes optimal classifier

$$f^*(x) = \mathbb{1}_{\{\eta(x) \geq \frac{1}{2}\}} \quad \Leftrightarrow \quad C^* = \{x \in \mathcal{X} : \eta(x) \geq \frac{1}{2}\}$$
$$\eta(x) = \mathbb{P}[Y=1 | X=x]$$
$$= \mathbb{E}[Y | X=x]$$

P unknown, iid samples given:

try to 'learn' the 'best' C in some \mathcal{C}



$$L(\hat{C}) \approx \inf_{C \in \mathcal{C}} L(C) \quad \text{whp}$$

Preview: \mathcal{C} has to be a VC class

Ex.: $\mathcal{X} = \mathbb{R}^d$

$$C = \{x \in \mathbb{R}^d : (w, x) + b \geq 0\}$$

$$w \in \mathbb{R}^d \setminus \{0\}$$
$$b \in \mathbb{R}$$

- linear discriminant rules

$$\text{VC-dim} = d+1$$

ERM, binary classification version

\mathcal{X} : feature space

\mathcal{C} : collection of subsets of \mathcal{X}

$(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} P$ on $\mathcal{X} \times \{0, 1\}$

ERM: $C \in \mathcal{C} \rightarrow \ell_C(x, y) := \frac{1}{2} |y \neq \mathbb{1}_C(x)|$

$$\hat{C}_n = \operatorname{argmin}_{C \in \mathcal{C}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell_C(x_i, y_i)}_{\text{fraction of mistakes of } C \text{ on the data}}$$

Abstract ERM (reminder):

$z_1, \dots, z_n \stackrel{iid}{\sim} P$ on \mathcal{Z}

\mathcal{F} : a class of fns $f: \mathcal{Z} \rightarrow [0, 1]$

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(z_i)$$

Here: $\mathcal{Z} = \mathcal{X} \times \{0, 1\}$ $z = (x, y)$

$$\mathcal{F} = \mathcal{F}_{\mathcal{C}} := \left\{ (x, y) \mapsto \ell_C(x, y) : C \in \mathcal{C} \right\}$$

Thm With prob. $\geq 1 - \delta$,

$$L(\hat{C}_n) \leq \inf_{C \in \mathcal{C}} L(C) + 4 \mathbb{E}_n(\mathcal{F}_{\mathcal{C}}(z^n)) + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

Fix $z^n = (z_1, \dots, z_n)$ $z_i = (x_i, y_i)$

$$\begin{aligned} \mathcal{F}_{\mathcal{C}}(z^n) &= \{ (f(z_1), \dots, f(z_n)) : f \in \mathcal{F}_{\mathcal{C}} \} \\ &= \left\{ \left(\mathbb{1}_{\{y_1 \neq 1, x_1 \in \mathcal{C}\}}, \dots, \mathbb{1}_{\{y_n \neq 1, x_n \in \mathcal{C}\}} \right) : c \in \mathcal{C} \right\} \\ &\subseteq \{0, 1\}^n \end{aligned}$$

FCL: $R_n(\mathcal{F}_{\mathcal{C}}(z^n)) \leq 2 \sqrt{\frac{\log |\mathcal{F}_{\mathcal{C}}(z^n)|}{n}}$

- if $\mathcal{F}_{\mathcal{C}}$ is a VC class [$V(\mathcal{F}_{\mathcal{C}}) < \infty$],
 $\log |\mathcal{F}_{\mathcal{C}}(z^n)| \leq V(\mathcal{F}_{\mathcal{C}}) \log(n+1)$ [Sauer-Shelah]

Key risk bound for ERM: w.p. $\geq 1 - \delta$,

$$L(\hat{C}_n) \leq \inf_{C \in \mathcal{C}} L(C) + 8 \sqrt{\frac{V(\mathcal{F}_{\mathcal{C}}) \log(n+1)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

NB: using Dudley's chaining,

$$\sqrt{\frac{V(\mathcal{F}_{\mathcal{C}}) \log(n+1)}{n}} \rightarrow \text{const} \sqrt{\frac{V(\mathcal{F}_{\mathcal{C}})}{n}}$$

$$\mathcal{F}_{\mathcal{C}} = \left\{ (x, y) \mapsto \mathbb{1}_{\{y \neq 1, x \in \mathcal{C}\}} : c \in \mathcal{C} \right\}$$

\mathcal{C} - fundamental object

Lemma

$$V(\mathcal{F}_{\mathcal{C}}) = V(\mathcal{C})$$

↓
 class of subsets of $\mathcal{X} \times \{0, 1\}$ ↪ class of subsets of \mathcal{X}

E.g.: \mathcal{C} are indicators of half-spaces in \mathbb{R}^d ,

$$V(\mathcal{F}_{\mathcal{C}}) = V(\mathcal{C}) = d+1$$

— excess risk $O\left(\sqrt{\frac{d + \log(1/\delta)}{n}}\right)$.

Proof (of lemma)

0) Assume $V(\mathcal{C}), V(\mathcal{F}_{\mathcal{C}}) < \infty$

1) A set $\{x_1, \dots, x_n\} \subset \mathcal{X}$ is shattered by \mathcal{C}
iff $\{(x_1, 0), \dots, (x_n, 0)\} \subset \mathcal{X} \times \{0, 1\}$ is
shattered by $\mathcal{F}_{\mathcal{C}}$

$$\forall c \in \mathcal{C} : l_c(x, y) = \mathbb{1}_{\{y \neq 1_c(x)\}}$$

$$l_c(x, 0) = \mathbb{1}_{\{0 \neq 1_c(x)\}}$$

$$= \mathbb{1}_{\{x \in c\}}$$

$$\begin{aligned} (\mathbb{1}_{\{x_1 \in c\}}, \dots, \mathbb{1}_{\{x_n \in c\}}) &= (b_1, \dots, b_n) \text{ iff} \\ (l_c(x_1, 0), \dots, l_c(x_n, 0)) &= (b_1, \dots, b_n) \end{aligned}$$

2) $V(\mathcal{C}) \leq V(\mathcal{F}_{\mathcal{C}})$

$$n = V(\mathcal{C})$$

$\exists \{x_1, \dots, x_n\} \subset \mathcal{X}$ shattered by \mathcal{C}

$\Rightarrow \{(x_1, 0), \dots, (x_n, 0)\} \subset \mathcal{X} \times \{0, 1\}$ shattered by $\mathcal{F}_{\mathcal{C}}$

3) $V(\mathcal{F}_{\mathcal{C}}) \leq V(\mathcal{C})$

$$n = V(\mathcal{F}_{\mathcal{C}})$$

$\exists \{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathcal{X} \times \{0, 1\}$ shattered by $\mathcal{F}_{\mathcal{C}}$

↓

$\forall b = (b_1, \dots, b_n) \in \{0, 1\}^n \quad \exists c \in \mathcal{C} \text{ s.t.}$

$$l_c(x_i, y_i) = \mathbb{1}_{\{y_i \neq l_c(x_i)\}} = b_i \quad i \in [n]$$

Goal: show that $\{x_1, \dots, x_n\} \subset \mathcal{X}$ shattered by \mathcal{C}

— need x_i 's to be distinct — can assume this:

$\{(x, 0), (x, 1)\}$ can't be shattered by $\mathcal{F}_{\mathcal{C}}$

$$l_c(x, 0) \neq l_c(x, 1) \quad \forall x \in \mathcal{X}$$

$$\begin{aligned} l_c(x_i, y_i) &= \mathbb{1}_{\{y_i \neq l_c(x_i)\}} & l_c(x, y) \\ &= y_i \oplus \mathbb{1}_{\{x_i \in C\}} &= \begin{cases} 0, & y = l_c(x) \\ 1, & y \neq l_c(x) \end{cases} \\ &= y_i \oplus l_c(x_i, 0) \end{aligned}$$

Know that $\{(x_1, 0), \dots, (x_n, 0)\}$ shattered by $\mathcal{F}_{\mathcal{C}}$
 $\Rightarrow \{x_1, \dots, x_n\}$ shattered by \mathcal{C}

Let $b = (b_1, \dots, b_n) \in \{0, 1\}^n$ be given;
take $b_i' := y_i \oplus b_i \quad \forall i \in [n]$

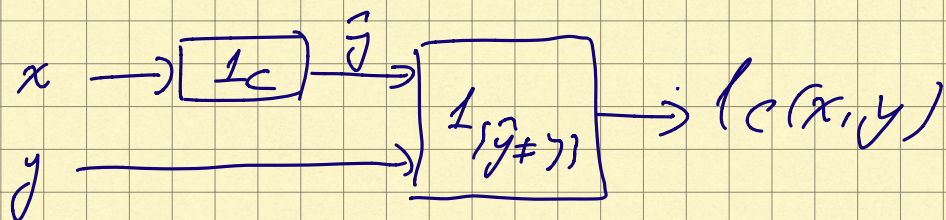
$\exists c \in \mathcal{C}$ s.t. $l_c(x_i, y_i) = b_i' \quad \forall i$

$$y_i \oplus \mathbb{1}_{\{x_i \in C\}} = y_i \oplus b_i$$

$$\mathbb{1}_{\{x_i \in C\}} = b_i \quad \square$$

Some intuition:

① $l_c(x, y)$ depends on x only through $\mathbb{1}_{\{x \in C\}}$



② $(X, Y) = (X, 0)$ w.p. 1

$$L(C) = \mathbb{P}[L_C(X) = 1] = \mathbb{P}(X \in C)$$

1) If \mathcal{C} is a VC-class, then the corresp. learning problem is PAC-learnable (by ERM).

2) The converse is also true: if binary classif. problem (for a given \mathcal{C}) is PAC-learnable, then $V(\mathcal{C}) < \infty$

$$\text{PAC: } \sup_P \{L_P(\hat{C}_n) - \inf_{C \in \mathcal{C}} L_P(C)\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

Examples

1) linear discriminant rules

$$\mathcal{X} = \mathbb{R}^d$$

\mathcal{C} - indicators of half-spaces

$$V(\mathcal{C}) = d+1$$

$$L(\hat{C}_n) - \inf_{C \in \mathcal{C}} L(C) = O(\sqrt{d/n}) \quad (\hat{C}_n: \text{ERM})$$

2) Dudley classifiers

\mathcal{X} arbitrary

\mathcal{G} : linear space of fns $\mathcal{X} \rightarrow \mathbb{R}$
spanned by lin. ind. $\{\psi_1, \dots, \psi_m\}$

$$g(x) = \sum_{j=1}^m c_j \psi_j(x)$$

$$\mathcal{X} \xrightarrow{(\psi_1, \dots, \psi_m)} \mathbb{R}^m \xrightarrow{(c_1, \dots, c_m)} g(x) = \sum_{j=1}^m c_j \psi_j(x)$$

(feature space)

$$C_g = \{x \in \mathcal{X} : g(x) \geq 0\} = \text{pos}(g)$$

$$V(\text{pos}(C_g)) = m, \quad \text{excess risk} = O(\sqrt{m/n})$$

Limitations

1) Expressivity: how small can $\inf_{C \in \mathcal{C}} L(C)$ be?

2) Computational tractability:

finding ERM classifier may be expensive

e.g. : complexity of finding a half-space in \mathbb{R}^d to minimize classification error on n points is $O(n^d)$ — prohibitive when $d \geq 5$.

Next lecture: surrogate losses via penalty fns