

BASICS of SUPERVISED LEARNING

Main goal (informally): Given a set of observations $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, $i = 1, 2, \dots, n$, predict the label $y \in \mathcal{Y}$ of a previously unseen input $x \in \mathcal{X}$.

Nomenclature: \mathcal{X} : input set, domain set
 $x \in \mathcal{X}$: input, feature, covariate
 $y \in \mathcal{Y}$: output, label, response

$S = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, 2, \dots, n\}$: training set, data set
previously unseen $(x, y) \in \mathcal{X} \times \mathcal{Y}$: test data

Some examples:

x	y
email	SPAM or not
image	CAT or DOG
image of a digit	0, 1, 2, ..., 9

Formal introduction to supervised learning:

First, let's describe the statistical nature of the problem.

$S = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, 2, \dots, n\}$ is given to the learner,
 $(x, y) \in \mathcal{X} \times \mathcal{Y}$ is the test data.

We assume that

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), (x, y)$ are random variables,

drawn independently from a distribution $P = P_X \otimes P_{Y|X}$ on $\mathcal{X} \times \mathcal{Y}$,
which is unknown to the learner.

The goal is (intuitively) to come up with a mapping $\hat{f}_S: \mathcal{X} \rightarrow \mathcal{Y}$
s.t. $\hat{f}_S(x) \approx y$ on the test data $(x, y) \in \mathcal{X} \times \mathcal{Y}$.

How to measure?

Loss function: $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. Given a predictor $f: \mathcal{X} \rightarrow \mathcal{Y}$,
 $l(f(x), y)$ is the loss incurred when the true label is y ,
and the prediction is $f(x)$.

For any prediction rule $f: \mathcal{X} \rightarrow \mathcal{Y}$, the performance criterion is

$$R(f) = \mathbb{E}_{(x, y) \sim P} [l(f(x), y)] \rightarrow \text{population risk}$$

The ultimate goal in supervised learning:

$$R^* = \inf_{\substack{f: \mathcal{X} \rightarrow \mathcal{Y} \\ f \text{ is measurable}}} R(f) \quad \text{is the Bayes risk}$$

Find a predictor $\hat{f} = \hat{f}_S$ based on S s.t. $R(\hat{f}) \approx R^*$ in expectation or with high probability.

Main Problem Classes in Supervised Learning

① Classification

\mathcal{X} an arbitrary set

$$\mathcal{Y} = \{0, 1\} \quad \text{or} \quad \mathcal{Y} = \{-1, 1\} \quad \left. \vphantom{\mathcal{Y}} \right\} \text{ binary classification}$$
$$\ell(\hat{y}, y) = \mathbb{1}\{\hat{y} \neq y\}$$

$$\mathcal{Y} = \{0, 1, \dots, m-1\} \quad \text{multi-class classification.}$$

② Regression

$$\mathcal{Y} = \mathbb{R}, \quad \ell(\hat{y}, y) = |\hat{y} - y|^2.$$

Remarks

① Why learning?

The underlying distribution of the nature, P , is unknown to the learner. Only partial knowledge via $S = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}; i \in [n]\}$ is available.

② Why supervised?

A supervisor pre-labels the training input data, the learner is trained under this supervision.

③ Why i.i.d.?

- (x_i, y_i) and (x, y) are identically distributed
 \Rightarrow all training samples provide useful info about the system and (x, y) .
- independence \Rightarrow each sample yields maximal info.

④ Realizable vs. agnostic: If $\exists f^*: \mathcal{X} \rightarrow \mathcal{Y}$ s.t. $y = f^*(x), \forall x \in \mathcal{X}$, the problem is called "realizable", and the task is to learn f^* .

Our setting is more general as it does not assume a deterministic mapping $x \mapsto f^*(x) = y$. The relation between $x \in \mathcal{X}$ and its label $y \in \mathcal{Y}$ can be random, governed by the conditional distribution $P_{Y|X}$.

Basics of Bayesian Decision - Making

Assume that P is known.

→ no learning

What is R^* ? How is $R(f)$ minimized?

Recall: X, Y, Z random variables.

$$E[E[X | Y, Z] | Y] = E[X | Y] \quad (\text{tower property})$$

By using the above property, for any $f: \mathcal{X} \rightarrow \mathcal{Y}$,

$$\begin{aligned} R(f) &= E[l(f(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} l(f(x), y) dP(x, y) \\ &= E[E[l(f(x), y) | x]] \\ &= \int_{\mathcal{X}} E[l(f(x'), y) | x = x'] dP_X(x') \end{aligned}$$

where P_X is the marginal distribution of X . Using this, we characterize R^* and the optimal prediction rule in the general case:

PROPOSITION For any $x' \in \mathcal{X}$, let

$$f^*(x') \in \arg \min_{y' \in \mathcal{Y}} E[l(y', y) | x = x'] = \int \mathbb{1}\{y' \neq y''\} dP_{Y|X}(y'' | x')$$

Then, $R(f^*) = R^*$, i.e., f^* is a Bayes optimal predictor.

Note: If we knew the conditional dist'n $P_{Y|X}$, we would be able to find f^* . Without the a priori knowledge of P or $P_{Y|X}$, it is not possible.

Let us make the above proposition more explicit on our two broad problem classes.

Corollary (Bayes optimal predictor for binary classification)

Consider $\mathcal{Y} = \{0, 1\}$, and $l(y, y') = \mathbb{1}\{y \neq y'\}$. Given the knowledge of P ,

$$f^*(x) = \begin{cases} 1, & \text{if } P_{Y|X}(1|x) > 1/2, \\ 0, & \text{if } P_{Y|X}(1|x) < 1/2, \end{cases}$$

Pf:

$$\begin{aligned}
 R(f) &= P(f(x) \neq y) = P(f(x) \neq 1, y=1) + P(f(x) \neq 0, y=0) \\
 &= \int_{\mathcal{X}} \left(\mathbb{1}\{f(x') \neq 1\} P(y=1|x=x') + \mathbb{1}\{f(x') \neq 0\} P(y=0|x=x') \right) dP_X(x') \\
 &= \int_{\mathcal{X}} \left(P_{Y|X}(1|x') + \underbrace{\mathbb{1}\{f(x') \neq 0\} [P_{Y|X}(0|x') - P_{Y|X}(1|x')]}_{(+)} \right) dP_X(x')
 \end{aligned}$$

How to minimize $R(f)$?

For $x' \in \mathcal{X}$, • if

$$P_{Y|X}(0|x') - P_{Y|X}(1|x') > 0, \text{ then set } f(x') = 0$$

so that

$$\mathbb{1}\{f(x') \neq 0\} = 0, \text{ and } (+) = 0.$$

Otherwise, (+) would be strictly positive.

Recall that $P_{Y|X}(0|x') + P_{Y|X}(1|x') = 1, \forall x' \in \mathcal{X}$.

Thus, set $f(x') = 0$ if

$$1 - 2P_{Y|X}(1|x') > 0 \Rightarrow P_{Y|X}(1|x') < 1/2$$

$$\bullet \text{ if } P_{Y|X}(0|x') - P_{Y|X}(1|x') = 1 - 2P_{Y|X}(1|x') \leq 0,$$

then one must set $f(x') = 1$ so that

$$\mathbb{1}\{f(x') \neq 0\} = 1 \text{ and } (+) < 0. \quad \square$$

Corollary 2 (Bayes optimal predictor for the regression problem)

Given P , $\mathcal{Y} = \mathbb{R}$ and $R(f) = \mathbb{E}[(f(x) - y)^2]$, $f: \mathcal{X} \rightarrow \mathcal{Y}$, the Bayes optimal predictor is

$$f^*(x) = \mathbb{E}[y|x=x'], \quad \forall x' \in \mathcal{X}.$$

Pf: $R(f) = \mathbb{E}[(y - f(x))^2]$

$$= \mathbb{E}[(y - \mathbb{E}[y|x] + \mathbb{E}[y|x] - f(x))^2]$$

$$= \mathbb{E}[(y - \mathbb{E}[y|x])^2] + \mathbb{E}[(f(x) - \mathbb{E}[y|x])^2] \\ + 2 \cdot \underbrace{\mathbb{E}[(y - \mathbb{E}[y|x])(f(x) - \mathbb{E}[y|x])]}_{(*)}$$

Note that

$$(*) = \mathbb{E}[(y - \mathbb{E}[y|x])(f(x) - \mathbb{E}[y|x])] \\ = \mathbb{E}\left[\mathbb{E}[(y - \mathbb{E}[y|x]) \cdot \underbrace{(f(x) - \mathbb{E}[y|x])}_{\sigma(x)\text{-measurable}} \mid x]\right] \\ = \mathbb{E}\left[(f(x) - \mathbb{E}[y|x]) \cdot \mathbb{E}[y - \mathbb{E}[y|x] \mid x]\right] \\ = \mathbb{E}\left[(f(x) - \mathbb{E}[y|x]) \cdot \underbrace{(\mathbb{E}[y|x] - \mathbb{E}[\mathbb{E}[y|x] \mid x])}_{=0}\right] \\ = 0$$

Thus,

$$\mathcal{R}(f) = \underbrace{\mathbb{E}[(y - \mathbb{E}[y|x])^2]}_{\text{does not depend on } f} + \underbrace{\mathbb{E}[(f(x) - \mathbb{E}[y|x])^2]}_{\text{can be minimized if } f(x) = \mathbb{E}[y|x] \text{ a.s.}}$$

Thus,

$$\mathcal{R}^* = \mathbb{E}[(y - \mathbb{E}[y|x])^2] \\ = \mathbb{E}[\text{Var}(y|x)].$$

In summary, the optimal prediction rule heavily depends on $P_{Y|X}$:

- for binary classification, $f^*(x) = \begin{cases} 1, & \text{if } P_{Y|X}(1|x) \geq 1/2 \\ 0, & \text{otherwise} \end{cases}$
- for regression, $f^*(x) = \mathbb{E}[y|x]$.