

# Machine Learning I

Supervised learning framework - Optimal predictions

---

Souhaib Ben Taieb

March 7, 2022

University of Mons

# Table of contents

Regression with squared error loss

Classification with zero-one loss

Bias and variance

# Optimal prediction function

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h)$$

$$g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{in}}(h)$$

Recall that the **optimal prediction function** is given by

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \underbrace{\mathbb{E}_x [E_{\text{out}}(h, x)]}_{E_{\text{out}}(h)}, \quad (1)$$

where

$$E_{\text{out}}(h, x) = \mathbb{E}_{y|x}[L(y, h(x))|x].$$

and  $L(\cdot, \cdot)$  is the loss function.

It sufficed to minimize the error pointwise, i.e. compute

$$f(x) = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h, x), \quad (2)$$

for all  $x \in \mathcal{X}$ .

# Table of contents

Regression with squared error loss

Classification with zero-one loss

Bias and variance

## Optimal predictions in regression (squared error loss)

With the squared error loss function  $L(y, \hat{y}) = (y - \hat{y})^2$ , the optimal prediction function is given by

$$f(x) = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{y|x}[(y - h(x))^2|x] \quad (3)$$

$$= \mathbb{E}_{y|x}[y|x], \quad (4)$$

i.e. the conditional expectation, also known as the **regression function**.

In other words, when *best is measured by expected squared error*, the best prediction for  $y$  at any point  $x$  is the conditional expectation at  $x$ .

## Optimal predictions in regression (squared error loss)

$$E_{\text{out}}(h, x) \tag{5}$$

$$= \mathbb{E}_{y|x}[(y - h(x))^2|x] \tag{6}$$

$$= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \tag{7}$$

$$= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \tag{8}$$

$$= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \tag{9}$$

$$= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2 \tag{10}$$

- The second term is non-negative, and will be equal to zero if  $h(x) = \mathbb{E}[y|x]$ .
- The first term corresponds to the inherent unpredictability, or noise, of the output, and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.

# Table of contents

Regression with squared error loss

Classification with zero-one loss

Bias and variance

## Optimal predictions in regression (zero-one loss)

For a multi-class classification problem with  $K$  categories, i.e.

$y \in \mathcal{C} = \{C_1, \dots, C_K\}$  and the zero-one loss function

$L(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\}$ , the optimal prediction function is given by

$$f(x) = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{C}} \mathbb{E}_{y|x}[\mathbb{1}\{y \neq h(x)\}|x] \quad (11)$$

$$= \operatorname{argmax}_{h: \mathcal{X} \rightarrow \mathcal{C}} \mathbb{P}(y = h(x)|x). \quad (12)$$

The optimal classifier is called the **Bayes classifier**, which has the following error rate at  $x$ :

$$1 - \max_{k=1, \dots, K} \mathbb{P}(y = C_k|x),$$

also called the **Bayes error rate**, which gives the lowest possible error rate that could be achieved if we knew  $\mathbb{P}(y|x)$ .



## Optimal predictions in regression (zero-one loss)

$$\begin{aligned}E_{\text{out}}(h, x) &= \mathbb{E}_{y|x}[\mathbb{1}\{y \neq h(x)\}|x] \\&= \sum_{k=1}^K \mathbb{1}\{C_k \neq h(x)\} \mathbb{P}(y = C_k|x) \\&= \sum_{k: C_k \neq h(x)} 1 \times \mathbb{P}(y = C_k|x) + 0 \times \mathbb{P}(y = h(x)|x) \\&= \sum_{k: C_k \neq h(x)} \mathbb{P}(y = C_k|x) \\&= \sum_{k: C_k \neq h(x)} \mathbb{P}(y = C_k|x) + \mathbb{P}(y = h(x)|x) - \mathbb{P}(y = h(x)|x) \\&= \sum_{k=1}^K \mathbb{P}(y = C_k|x) - \mathbb{P}(y = h(x)|x) \\&= 1 - \mathbb{P}(y = h(x)|x).\end{aligned}$$

# Optimal predictions in classification

Using the fundamental bridge, we can directly write

$$\begin{aligned}\mathbb{E}_{y|x}[\mathbb{1}\{y \neq h(x)\}|x] \\ &= \mathbb{P}(y \neq h(x)|x) \\ &= 1 - \mathbb{P}(y = h(x)|x).\end{aligned}$$

In conclusion, we have

$$f(x) = \operatorname{argmin}_{h:\mathcal{X} \rightarrow \mathcal{C}} \mathbb{E}_{y|x}[\mathbb{1}\{y \neq h(x)\}|x] \quad (13)$$

$$= \operatorname{argmin}_{h:\mathcal{X} \rightarrow \mathcal{C}} 1 - \mathbb{P}(y = h(x)|x) \quad (14)$$

$$= \operatorname{argmax}_{h:\mathcal{X} \rightarrow \mathcal{C}} \mathbb{P}(y = h(x)|x). \quad (15)$$

# Table of contents

Regression with squared error loss

Classification with zero-one loss

Bias and variance

# Quantifying the approximation-generalization tradeoff

The difference between the out-of-sample error of  $g$  and  $f$  can be decomposed as follows

$$E_{\text{out}}(g) - E_{\text{out}}(f) = \underbrace{[E_{\text{out}}(g^*) - E_{\text{out}}(f)]}_{\text{Approximation error}} + \underbrace{[E_{\text{out}}(g) - E_{\text{out}}(g^*)]}_{\text{Estimation error}}$$

Recall that simple models underfit the data and complex models overfit. There is an **approximation-generalization** tradeoff.

The **bias-variance** decomposition allows to quantify the approximation-generalization tradeoff for the **squared error** loss function.

# The bias-variance tradeoff

See board.

# The bias-variance tradeoff

