

Machine Learning I

Supervised Learning: Optimal Predictions

Souhaib Ben Taieb

University of Mons



Table of contents

Optimal predictions with the squared error loss

Optimal predictions with the zero-one loss

Optimal predictions with a general binary classification loss

Optimal prediction function

$$\boxed{f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h)} \quad g^* = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{out}}(h), \quad g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{in}}(h)$$

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

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{x,y} [L(y, h(x))], \quad (1)$$

where $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is the loss function.

Using the **law of iterated expectation**, we can write

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

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

$$f(x) = \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}_{y|x} [L(y, h(x))|x], \quad (2)$$

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

Optimal prediction function

$$\boxed{f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h)} \quad g^* = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{out}}(h), \quad g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{in}}(h)$$

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

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{x,y}[L(y, h(x))], \quad (1)$$

where $L: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is the loss function.

Using the **law of iterated expectation**, we can write

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

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

$$f(x) = \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}_{y|x}[L(y, h(x))|x], \quad (2)$$

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

Optimal prediction function

$$\boxed{f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h)} \quad g^* = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{out}}(h), \quad g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{in}}(h)$$

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

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{x,y}[L(y, h(x))], \quad (1)$$

where $L: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is the loss function.

Using the **law of iterated expectation**, we can write

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

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

$$f(x) = \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}_{y|x}[L(y, h(x))|x], \quad (2)$$

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

Table of contents

Optimal predictions with the squared error loss

Optimal predictions with the zero-one loss

Optimal predictions with a general binary classification loss

Optimal predictions with the squared error loss

Let $\mathcal{Y} \subseteq \mathbb{R}$. With the squared error loss function $L(y, \hat{y}) = (y - \hat{y})^2$, the expected error at x can be decomposed as follows:

$$\begin{aligned}\mathbb{E}[L(y, h(x))|x] &= \mathbb{E}[(y - h(x))^2|x] \\ &= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \\ &= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\ &= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\ &= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2\end{aligned}$$

- ▶ The first term corresponds to the inherent unpredictability, or noise of y , and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.
- ▶ The second term is non-negative, and will be equal to zero if

$$h(x) = \mathbb{E}[y|x].$$

Optimal predictions with the squared error loss

Let $\mathcal{Y} \subseteq \mathbb{R}$. With the squared error loss function $L(y, \hat{y}) = (y - \hat{y})^2$, the expected error at x can be decomposed as follows:

$$\begin{aligned}\mathbb{E}[L(y, h(x))|x] &= \mathbb{E}[(y - h(x))^2|x] \\&= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \\&= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2\end{aligned}$$

- ▶ The first term corresponds to the inherent unpredictability, or noise of y , and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.
- ▶ The second term is non-negative, and will be equal to zero if

$$h(x) = \mathbb{E}[y|x].$$

Optimal predictions with the squared error loss

Let $\mathcal{Y} \subseteq \mathbb{R}$. With the squared error loss function $L(y, \hat{y}) = (y - \hat{y})^2$, the expected error at x can be decomposed as follows:

$$\begin{aligned}\mathbb{E}[L(y, h(x))|x] &= \mathbb{E}[(y - h(x))^2|x] \\&= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \\&= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2\end{aligned}$$

- ▶ The first term corresponds to the inherent unpredictability, or noise of y , and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.
- ▶ The second term is non-negative, and will be equal to zero if

$$h(x) = \mathbb{E}[y|x].$$

Optimal predictions with the squared error loss

Let $\mathcal{Y} \subseteq \mathbb{R}$. With the squared error loss function $L(y, \hat{y}) = (y - \hat{y})^2$, the expected error at x can be decomposed as follows:

$$\begin{aligned}\mathbb{E}[L(y, h(x))|x] &= \mathbb{E}[(y - h(x))^2|x] \\&= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \\&= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\&= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2\end{aligned}$$

- ▶ The first term corresponds to the inherent unpredictability, or noise of y , and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.
- ▶ The second term is non-negative, and will be equal to zero if

$$h(x) = \mathbb{E}[y|x].$$

Optimal predictions with the squared error loss

Let $\mathcal{Y} \subseteq \mathbb{R}$. With the squared error loss function $L(y, \hat{y}) = (y - \hat{y})^2$, the expected error at x can be decomposed as follows:

$$\begin{aligned}\mathbb{E}[L(y, h(x))|x] &= \mathbb{E}[(y - h(x))^2|x] \\ &= \mathbb{E}[y^2 - 2yh(x) + h(x)^2|x] \\ &= \mathbb{E}[y^2|x] - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\ &= \text{Var}(y|x) + (\mathbb{E}[y|x])^2 - 2h(x)\mathbb{E}[y|x] + h(x)^2 \\ &= \text{Var}(y|x) + (\mathbb{E}[y|x] - h(x))^2\end{aligned}$$

- ▶ The first term corresponds to the inherent unpredictability, or noise of y , and is called the **Bayes error**. It is the smallest error any learning algorithm can achieve.
- ▶ The second term is non-negative, and will be equal to zero if

$$h(x) = \mathbb{E}[y|x].$$

Optimal predictions with the squared error loss

In summary, the optimal prediction at x is given by

$$\begin{aligned} f(x) &= \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}[(y - h(x))^2 | x] \\ &= \mathbb{E}[y | x], \end{aligned}$$

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 .

Table of contents

Optimal predictions with the squared error loss

Optimal predictions with the zero-one loss

Optimal predictions with a general binary classification loss

Optimal predictions with the zero-one loss

Consider a multi-class classification problem with K categories where $\mathcal{Y} = \{C_1, \dots, C_K\}$. With the zero-one loss $L(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\}$, the expected error (error rate) at x can be decomposed as follows:

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

Hence, we have

$$\begin{aligned}f(x) &= \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}[\mathbb{1}\{y \neq h(x)\}|x] \\ &= \operatorname{argmin}_{h(x) \in \mathcal{Y}} 1 - \mathbb{P}(y = h(x)|x) \\ &= \operatorname{argmax}_{h(x) \in \mathcal{Y}} \mathbb{P}(y = h(x)|x).\end{aligned}$$

Optimal predictions with the zero-one loss

Consider a multi-class classification problem with K categories where $\mathcal{Y} = \{C_1, \dots, C_K\}$. With the zero-one loss $L(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\}$, the expected error (error rate) at x can be decomposed as follows:

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

Hence, we have

$$\begin{aligned}f(x) &= \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}[\mathbb{1}\{y \neq h(x)\}|x] \\ &= \operatorname{argmin}_{h(x) \in \mathcal{Y}} 1 - \mathbb{P}(y = h(x)|x) \\ &= \operatorname{argmax}_{h(x) \in \mathcal{Y}} \mathbb{P}(y = h(x)|x).\end{aligned}$$

Optimal predictions with the zero-one loss

In summary, the **optimal prediction** at x is given by

$$f(x) = \operatorname{argmin}_{h(x) \in \mathcal{Y}} \mathbb{E}[\mathbb{1}\{y \neq h(x)\} | x] \quad (3)$$

$$= \operatorname{argmax}_{h(x) \in \mathcal{Y}} \mathbb{P}(y = h(x) | x). \quad (4)$$

This optimal classifier is called the **Bayes classifier**, and has the following expected error (error rate) at x :

$$\begin{aligned} \mathbb{E}[\mathbb{1}\{y \neq h(x)\} | x] &= 1 - \mathbb{P}(y = \operatorname{argmax}_{h(x) \in \mathcal{Y}} \mathbb{P}(y = h(x) | x) | x) \\ &= 1 - \max_{k=1, \dots, K} \mathbb{P}(y = C_k | x), \end{aligned}$$

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

Table of contents

Optimal predictions with the squared error loss

Optimal predictions with the zero-one loss

Optimal predictions with a general binary classification loss

Optimal predictions

Consider a binary classification problem where $\mathcal{Y} = \{0, 1\}$ and the general binary classification loss function:

$$L(y, \hat{y}) = \begin{cases} L(0, 0) & \text{if } y = 0 \text{ and } \hat{y} = 0; \\ L(0, 1) & \text{if } y = 0 \text{ and } \hat{y} = 1; \\ L(1, 0) & \text{if } y = 1 \text{ and } \hat{y} = 0; \\ L(1, 1) & \text{if } y = 1 \text{ and } \hat{y} = 1, \end{cases}$$

where we assume $L(1, 0) > L(1, 1) \geq 0$ and $L(0, 1) > L(0, 0) \geq 0$.

The **zero-one loss** is a particular case where $L(0, 0) = L(1, 1) = 0$ and $L(1, 0) = L(0, 1) = 1$.

The **optimal prediction** at x is given by

$$f(x) = \operatorname{argmin}_{h(x) \in \{0, 1\}} \mathbb{E}[L(y, h(x)) | x].$$

Let us consider the two cases: $h(x) = 0$ and $h(x) = 1$.

Optimal predictions

Consider a binary classification problem where $\mathcal{Y} = \{0, 1\}$ and the general binary classification loss function:

$$L(y, \hat{y}) = \begin{cases} L(0, 0) & \text{if } y = 0 \text{ and } \hat{y} = 0; \\ L(0, 1) & \text{if } y = 0 \text{ and } \hat{y} = 1; \\ L(1, 0) & \text{if } y = 1 \text{ and } \hat{y} = 0; \\ L(1, 1) & \text{if } y = 1 \text{ and } \hat{y} = 1, \end{cases}$$

where we assume $L(1, 0) > L(1, 1) \geq 0$ and $L(0, 1) > L(0, 0) \geq 0$.

The **zero-one loss** is a particular case where $L(0, 0) = L(1, 1) = 0$ and $L(1, 0) = L(0, 1) = 1$.

The **optimal prediction** at x is given by

$$f(x) = \operatorname{argmin}_{h(x) \in \{0, 1\}} \mathbb{E}[L(y, h(x)) | x].$$

Let us consider the two cases: $h(x) = 0$ and $h(x) = 1$.

Optimal predictions

We can expand the expected loss for each of the two possible predictions.

$h(x) = 0$:

$$\mathbb{E}[L(y, 0)|x] = L(0, 0)\mathbb{P}(y = 0|x) + L(1, 0)\mathbb{P}(y = 1|x).$$

$h(x) = 1$:

$$\mathbb{E}[L(y, 1)|x] = L(0, 1)\mathbb{P}(y = 0|x) + L(1, 1)\mathbb{P}(y = 1|x).$$

Since we want to minimize the expected loss, the optimal prediction is $h(x) = 1$ ($h(x) = 0$) whenever the second expression is smaller (larger) than the first.

Optimal predictions

In other words, the optimal prediction at x is given by

$$\begin{aligned} f(x) &= \mathbb{1}\{\mathbb{E}[L(y, 1)|x] \leq \mathbb{E}[L(y, 0)|x]\} \\ &= \mathbb{1}\left\{\mathbb{P}(y = 1|x) \geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)}\mathbb{P}(y = 0|x)\right\} \\ &= \mathbb{1}\left\{\mathbb{P}(y = 1|x) \geq \frac{L(0, 1) - L(0, 0)}{L(0, 1) - L(0, 0) + L(1, 0) - L(1, 1)}\right\} \end{aligned}$$

Examples

Consider the following three cases:

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 1, L(1, 1) = 0$ (zero-one loss)

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq \mathbb{P}(y = 0|x) \} = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{2} = 0.5 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 10, L(1, 1) = 0$

$$f(x) = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{10} \times \mathbb{P}(y = 0|x) \right\} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{11} \approx 0.09 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1000, L(1, 0) = 1, L(1, 1) = 0$

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq 1000 \times \mathbb{P}(y = 0|x) \} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1000}{1001} \approx 0.99 \right\}$$

Examples

Consider the following three cases:

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 1, L(1, 1) = 0$ (zero-one loss)

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq \mathbb{P}(y = 0|x) \} = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{2} = 0.5 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 10, L(1, 1) = 0$

$$f(x) = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{10} \times \mathbb{P}(y = 0|x) \right\} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{11} \approx 0.09 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1000, L(1, 0) = 1, L(1, 1) = 0$

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq 1000 \times \mathbb{P}(y = 0|x) \} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1000}{1001} \approx 0.99 \right\}$$

Examples

Consider the following three cases:

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 1, L(1, 1) = 0$ (zero-one loss)

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq \mathbb{P}(y = 0|x) \} = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{2} = 0.5 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1, L(1, 0) = 10, L(1, 1) = 0$

$$f(x) = \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{10} \times \mathbb{P}(y = 0|x) \right\} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1}{11} \approx 0.09 \right\}$$

- $L(0, 0) = 0, L(0, 1) = 1000, L(1, 0) = 1, L(1, 1) = 0$

$$f(x) = \mathbb{1} \{ \mathbb{P}(y = 1|x) \geq 1000 \times \mathbb{P}(y = 0|x) \} \mathbb{1} \left\{ \mathbb{P}(y = 1|x) \geq \frac{1000}{1001} \approx 0.99 \right\}$$