

The classification challenge in machine learning focuses on separating data into distinct categories. We have an input domain  $X$  (all potential examples) and an output domain  $Y$  (the corresponding categories). The primary objective is to develop an algorithm that accurately assigns examples to their categories while reducing errors. Classification problems can be categorized into two main types: Supervised Learning: The algorithm is trained on labeled data, meaning the correct classifications are already known (e.g., identifying objects).

Unsupervised Learning: Here, the algorithm seeks patterns in unlabeled data (e.g., grouping customers based on their buying habits).

Classification falls under the supervised learning category. In this framework, we consider two spaces:  $X$  (the input data) and  $Y$  (the corresponding output labels). In the case of binary classification, we categorize objects into two classes:  $-1$  and  $+1$ . The aim is to determine a function  $f: X \rightarrow Y$  that accurately matches objects to their labels. This is achieved using a training dataset comprising pairs  $(X_1, Y_1), \dots, (X_n, Y_n)$ , where each pair is drawn independently from a broader distribution. Given this dataset, the classification algorithm devises a function  $f$  that seeks to minimize classification errors. This process allows the machine to learn how to classify new instances.

Ultimately, the challenge lies in identifying the optimal function

$f$  that effectively maps inputs  $X$  to output labels  $Y$ . The ideal classifier, often referred to as the Bayesian classifier, operates on the principle that if the probability of an object belonging to class  $+1$  is at least  $0.5$ , it should be labeled as  $+1$ ; otherwise, it receives the label  $-1$ . This approach is statistically optimal for minimizing errors, but in practice, accurately calculating it is complicated due to the unknown nature of the data distribution.

Thus, the binary classification task can be summarized as: given a training dataset  $(X_1, Y_1), \dots, (X_n, Y_n)$  sourced from an unknown distribution  $P$  and a defined loss function, the goal is to construct a function  $f: X \rightarrow Y$  that minimizes the risk  $R(f)$ , thereby approaching the risk associated with the Bayesian classifier.