# Classification

Machine Learning Course - CS-433
Oct 18, 2022
Nicolas Flammarion
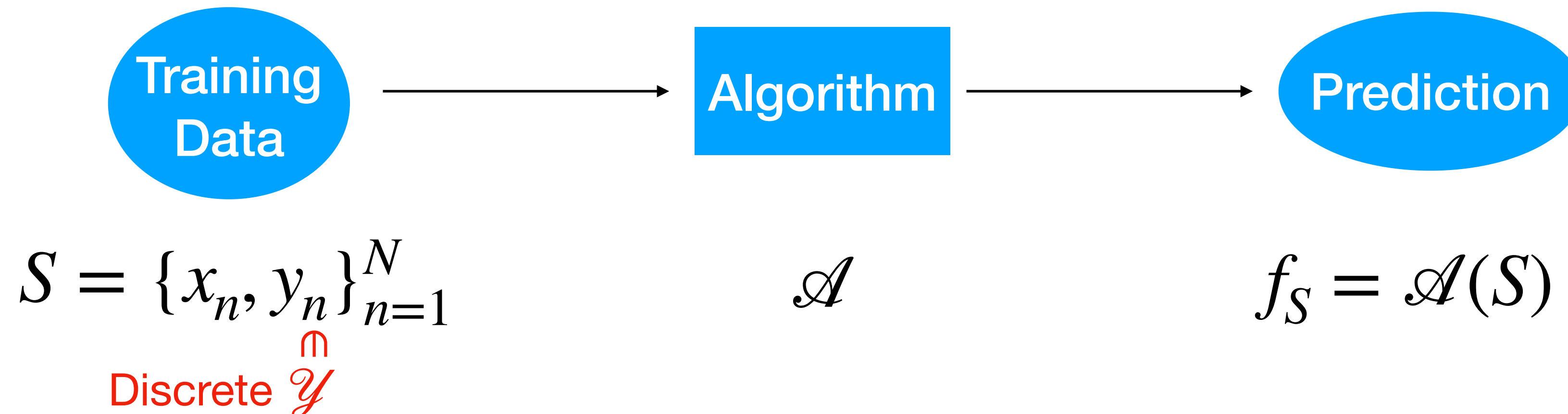
# Definition of classification

We observe some data $S = \{x_n, y_n\}_{n=1}^N \in \mathcal{X} \times \underbrace{\mathcal{Y}}$

**Discrete Set**

<u>Goal</u>: given a new observation $x$, we want to predict its label $y$

<u>How:</u>

Training Data → Algorithm → Prediction

$$S = \{x_n, y_n\}_{n=1}^N \qquad\qquad \mathcal{A} \qquad\qquad f_S = \mathcal{A}(S)$$

Discrete $\mathcal{Y}$

# Classification: relates input to a categorical variable

$$(x, y) \in \mathcal{X} \times \underbrace{\mathcal{Y}}_{\text{Discrete Set}}$$

**Binary Classification**: $y$ can take two values

$y \in \{c_1, c_2\}$ where $c_i$ are the class labels. We often use $\{0,1\}$ or $\{-1,1\}$

**Multi-class classification**: $y$ can take more than two values

$y \in \{c_1, \cdots, c_{K-1}\}$ for a $K$ classes problem. We often use $\{0, \cdots, K-1\}$

⚠ no ordering between classes

# Spam Detection

# Image classification

# Credit Card Default



+   individual who defaulted on their credit card payments

○   individual who did not

# Classifier

A classifier $f : \mathcal{X} \to \mathcal{Y}$ divides the input space into a collection of regions belonging to each class



**Decision boundary**

# Classifier

A classifier $f : \mathcal{X} \to \mathcal{Y}$ divides the input space into a collection of regions belonging to each class

It can be linear



**Linear Decision boundary**

# Classifier

A classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ divides the input space into a collection of regions belonging to each class

It can also be nonlinear



**Nonlinear Decision boundary**

# Classification: a special case of regression?

Classification is a **regression problem** with discrete labels:

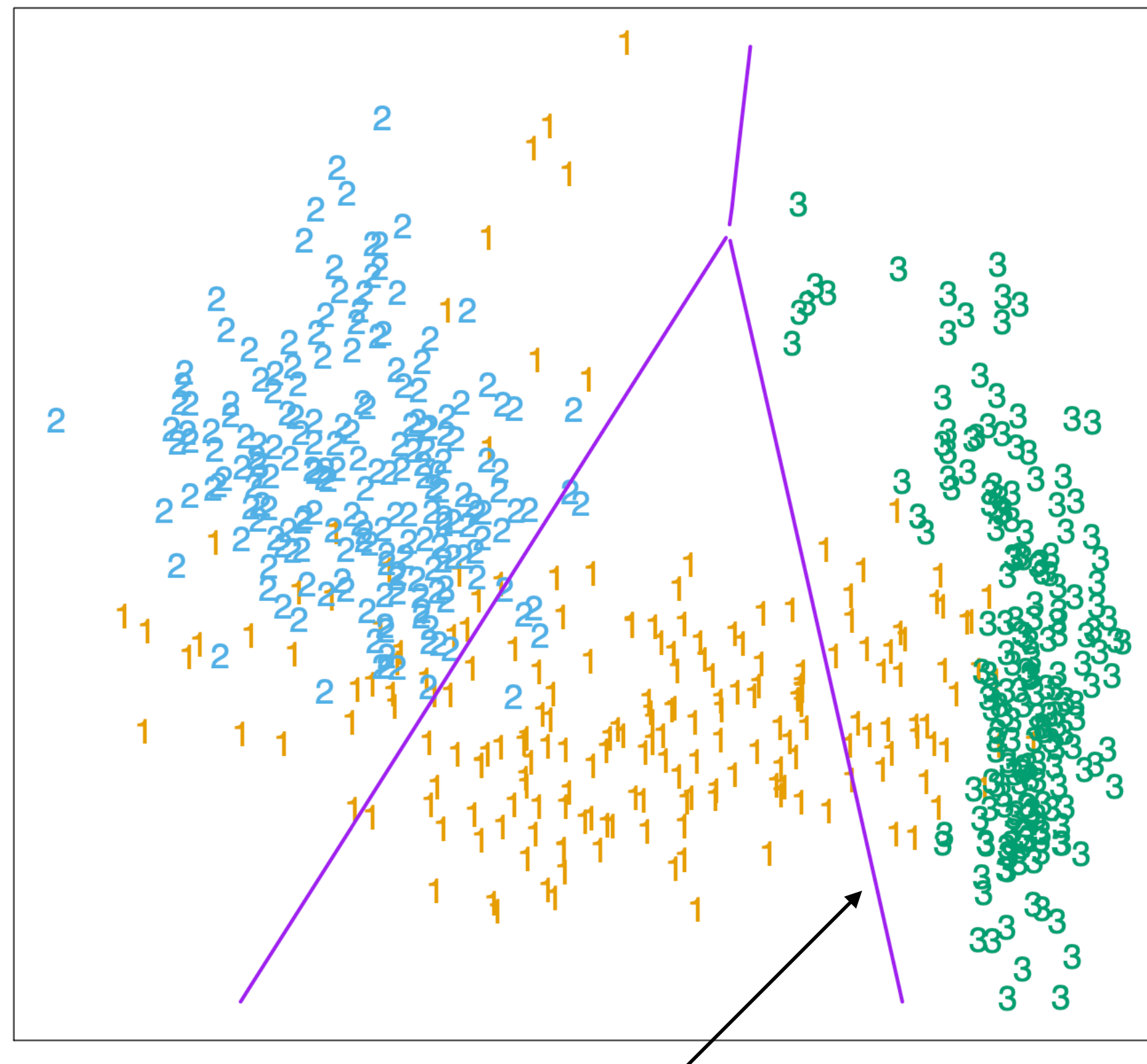$$(x, y) \in \mathcal{X} \times \{0, 1\} \subset \mathcal{X} \times \mathbb{R}$$

Could we use previously seen regression methods to solve it?



$$S_{\text{train}} \colon y = \begin{cases} 0 & \text{if} \quad C_1 \\ 1 & \text{if} \quad C_2 \end{cases}$$

$$f_{S_{\text{train}}}$$

$$\begin{cases} C_1 & \text{if} \quad f_{S_{\text{train}}}(x) < 0.5 \\ C_2 & \text{if} \quad f_{S_{\text{train}}}(x) \geq 0.5 \end{cases}$$

# Is it a good idea?

Credit-card default problem:

# Classification is not just a special form of regression

A. The predicted values are not probabilities (not in [0,1])

# Classification is not just a special form of regression

B. Sensitivity to unbalanced data



The position of the line depends crucially on how many points are in each class

# Classification is not just a special form of regression

C. Sensitivity to extreme values:



The position of the line depends crucially on where the points lie

**Why**: the square loss we used for regression is not suitable for classification

# How to perform classification?

- A lot of different approaches have been developed
- We will not detail them exhaustively today
- Rather we will provide quick introductions

- Fundamental task of classification:

**separate the space into various decision regions**

# Nearest Neighbor

Assume that nearby points are likely to have similar label



You assign the label of the closest point in your training set.

# Nearest Neighbor

Pros:

- **No optimization** or training
- **Easy** to implement
- Works well in **low dimensions** where you can get very complex decision boundaries

Cons:

- **Slow** at query time
- Bad for high dimensional data
- Choice of **local distance** is crucial

# k-Nearest Neighbor

A new point $x$ is classified by a **majority vote among the k-nearest neighbor of** $x$



Generalization: smoothing kernels; weighted linear combination of elements

# Linear Decision boundaries

Assume we restrict ourself to linear decision boundaries (hyperplane):

➡ Prediction: $g(x) = \text{sign}(x^\top w)$

# Separating hyperplane

Assume we restrict ourself to linear decision boundaries (hyperplane):

Assume the data are linearly separable, i.e, a separating hyperplane exists



Which separating hyperplane would you pick?

# Margin

Key concept: The margin is the distance between the hyperplane and the closest point



➡ Take the one with the largest margin!

# Max-margin separating hyperplane

Pick the hyperplane which maximizes the margin



Why: If we slightly change the training set, the number of misclassifications will stay low
➡ It will lead us to support vector machine (SVM) and logistic regression

# Nonlinear classifier

- Linear decision boundaries will not always work.

- Feature augmentation $(x, x^2, x^3, x^4)$

- Kernel Method

# Do we still have time?

# A little bit of theory

$$(X, Y) \sim \mathcal{D} \text{ with ranges } \mathcal{X}, \mathcal{Y} = \{0,1\}$$

Loss function:

$$\ell(y, y') = 1_{y \neq y'} = \begin{cases} 1 & \text{if } y \neq y' \\ 0 & \text{if } y = y' \end{cases}$$

True risk for the classification:

$$L_{\mathcal{D}}(g) = \mathbb{E}_{\mathcal{D}}[1_{Y \neq g(X)}] = \mathbb{P}_{\mathcal{D}}[Y \neq g(X)]$$

classification error          probability of making an error

# Bayes classifier

What is the **optimal performance**, regardless of the finiteness of the training data?

Def: The classifier $g_* = \arg\min\limits_{g} L_{\mathscr{D}}(g)$ is called the **Bayes classifier**

Claim:

$$g_*(x) = \arg\max\limits_{y \in \{0,1\}} \mathbb{P}(Y = y \,|\, X = x)$$

# Proof of the Bayes classifier

**Claim 1**: $\forall x \in \mathcal{X}, h_*(x) \in \arg\min_{y \in \mathcal{Y}} \mathbb{P}(Y \neq y \,|\, X = x) \implies h_* \in \arg\min_{h:\mathcal{X} \to \mathcal{Y}} L_{\mathcal{D}}(h)$

$$
\begin{aligned}
L_{\mathcal{D}}(h) = \mathbb{E}_{X,Y}[1_{Y \neq h(X)}] &= \mathbb{E}_X[\mathbb{E}_{Y|X}[1_{Y \neq h(X)} \,|\, X]] \\
&= \mathbb{E}_X[\mathbb{P}(Y \neq h(X) \,|\, X)] \\
&\geq \mathbb{E}_X[\min_{y \in \mathcal{Y}} \mathbb{P}(Y \neq y \,|\, X)] \\
&= \mathbb{E}_X[\mathbb{P}(Y \neq h_*(X) \,|\, X)] = \mathbb{E}_{X,Y}[1_{Y \neq h_*(X)}] = L_{\mathcal{D}}(h_*)
\end{aligned}
$$

**Claim 2:** $g_*(x) = \arg\min_{y \in \mathcal{Y}} \mathbb{P}(Y \neq y \,|\, X = x)$

$$
g_*(x) = \arg\max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \,|\, X = x) = \arg\min_{y \in \mathcal{Y}} \mathbb{P}(Y \neq y \,|\, X = x)
$$

# Bonus: a good regressor implies a good classifier

For all regression functions $\eta : \mathscr{X} \to \mathbb{R}$ we can define a classifier as

$$\mathscr{X} \to \{0,1\}$$

$$g_\eta : \ x \mapsto 1_{\eta(x) \geq 1/2}$$

**<u>Claim:</u>**

$$L_{\mathscr{D}}^{\mathsf{classif}}(g_\eta) - L_{\mathscr{D}}^{\mathsf{classif}}(g*) \leq 2\sqrt{L_{\mathscr{D}}^{\ell_2}(\eta) - L_{\mathscr{D}}^{\ell_2}(\eta*)}$$

Where $L_{\mathscr{D}}^{\mathsf{classif}}(g_\eta) = \mathbb{E}_{\mathscr{D}}[1_{g(X) \neq Y}]$, $L_{\mathscr{D}}^{\ell_2}(f) = \mathbb{E}_{\mathscr{D}}[(Y - f(X))^2]$ and $\eta_* = \arg\min_\eta L_{\mathscr{D}}^{\ell_2}(\eta)$

➡ If $\eta$ is good for regression then $g_\eta$ is good for classification too (converse is not true)