# Classification Problems

<u>Examples</u>:
- ▷ Fraud detection in credit card payments
- ▷ Categorize digital images: dog, cat, human

. <u>Supervised Learning</u>:

- ▷ $X \subset \mathbb{R}^k$ : domain set
- ▷ $Y \subset N$ : <u>finite</u> target set, $|N| = q$

# 1. (Multiclass) Logistic Regression

$$X \subset \mathbb{R}^k, \quad Y = \{0,1\}^q \qquad F = \left\{ h: \mathbb{R}^k \to \mathbb{R}^q : x \longmapsto \text{softmax}(Wx), \; W \in \mathbb{R}^{q \times k} \right\}$$

with $\text{softmax}(z) = \underset{\mathbb{R}^q}{\underbrace{\left[ \dfrac{\exp(z_1)}{\sum\limits_{i=1}^{q} \exp(z_i)}, \; \dfrac{\exp(z_2)}{\sum\limits_{i=1}^{q} \exp(z_i)}, \; \cdots, \; \dfrac{\exp(z_q)}{\sum\limits_{i=1}^{q} \exp(z_i)} \right]}}$

▷ <u>**loss function**</u> $\quad \ell(y, z) := - \sum\limits_{i=1}^{q} y_i \log(z_i) = - \langle y, \log(z) \rangle$

> LinReg:
> $\ell(y,z) = (y-z)^2$

▷ <u>**empirical risk**</u> $\quad L_S(W) = \dfrac{1}{n} \sum\limits_{j=1}^{n} \ell\left( y^j, \text{softmax}(Wx^j) \right)$

if $S = \left\{ x^j, y^j \right\}_{j=1}^{n}$

with **hot encoding** s.t. $\quad y^j = (0, 0, \ldots, 0, 1, 0, \ldots, 0)$

⤷ if $j$-th data point is in class $i$

ith position

- Add ridge/lasso type regularization term is an option

$\left[ W \longmapsto L_S(W) \text{ is convex} \right]$

- Optimization is non-trivial, but *convex* optimization can be used.

# 2. K - Nearest Neighbors

Let $X \subset \mathbb{R}^k$, $Y \subset \{0,1\}^q$. Let $d : X \times X \to \mathbb{R}$ be a <span style="color:yellow">metric</span>,

e.g. $d(x, x') := \|x - x'\|_2$.

If $S_x = \{x^1, \dots, x^n\}$ is a set, define $\pi_i(x)$ as the
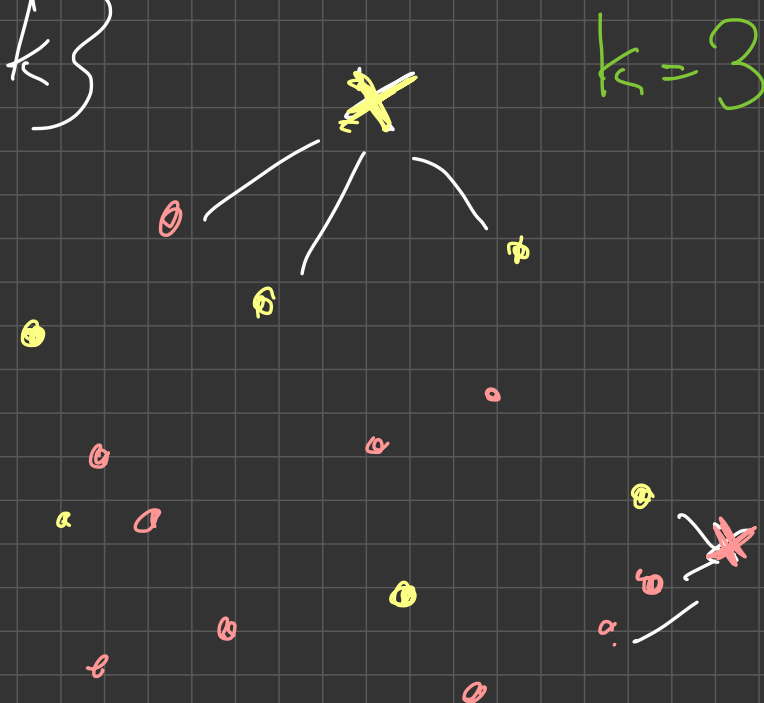
$i$-th closest member of $S_x$ to $x$ (w.r.t. metric $d$)

Algorithm: Input: Training set $S = \{x^j, y^j\}_{j=1}^n$, <span style="color:green">parameter $k$</span>.

Output: function $h_S : X \to Y$ such that $h_S(x)$ is
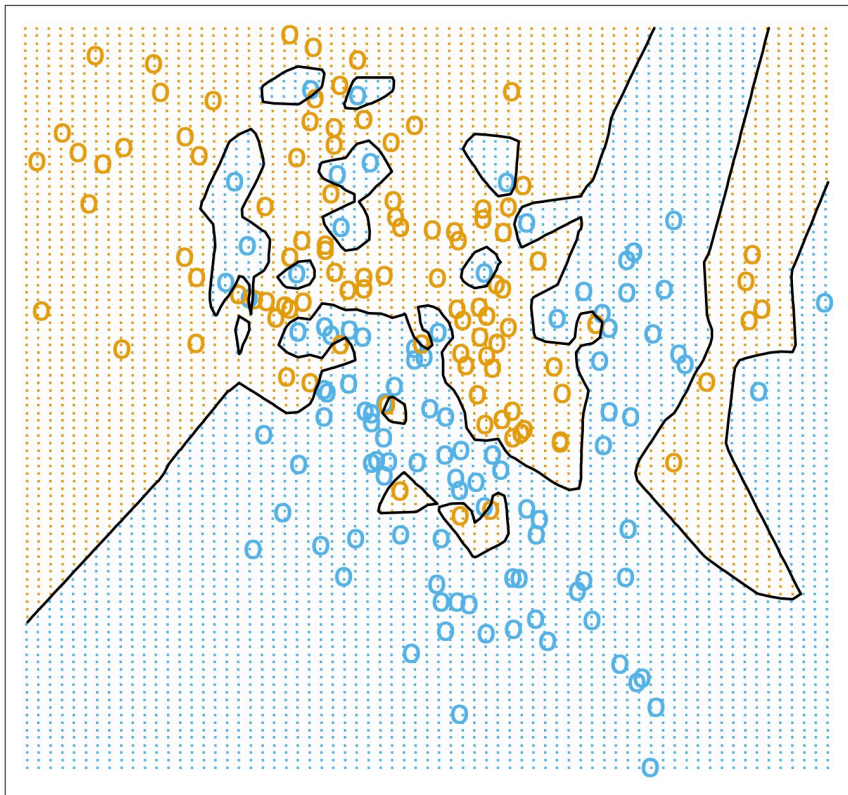
the majority label among $\{y^{\pi_i(x)}, i \in k\}$

$\oplus$ "local" method, simple

$\ominus$ Needs all pairwise comparison

( $\hat{O}(kn)$ computations )

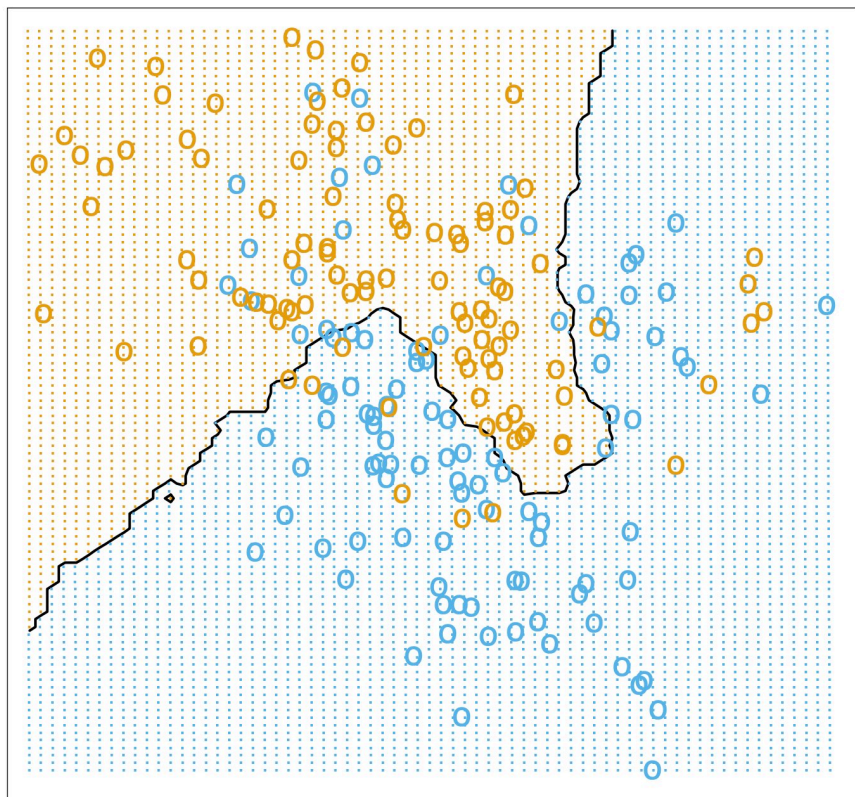$\ominus$ Suffer from a "course of dimensionality"

if $k$ is large

<span style="color:green">$k = 3$</span>

1-Nearest Neighbor Classifier

15-Nearest Neighbor Classifier

# Feature selection: Natural language Processing

$Q:$ How to work with **text** data?

→ spam/no detector

categorizing text items according to topics

"document"

$X =$ " Michael likes walking his dog in his neighborhood. "

## 1. Tokenize:

" Michael ", " likes " " walking ", " dog ", " in ", " his ", " neighborhood "

## 2. Build dictionary : Do this for all documents $x$ of interest.

" Aaron ", " Amsterdam ", " am "..... , " his ", ..? "Michael" ,... "walking" ,..., " zebra "

List of $d$ words

## 3. Encoding : a) Count how often each word/token in $x$ occurs, create

sparse vector in $\mathbb{R}^d$ :

$$\tilde{\Phi} : \mathcal{D} \to \mathbb{R}^d \quad , \quad \Phi(x) = \frac{\tilde{\Phi}(x)}{\|\tilde{\Phi}(x)\|} , \text{ where } \tilde{\Phi}(x) = ( 0, 0, 0, ..., 2, ...?, ..., 1..., 0)$$

$\tilde{\Phi}(x)_w$

Modifications can include:

▽ Removal of very common words such as "in", "the"

▽ "n-grams": Use "Michael likes", "likes walking", etc. as tokens

("2-gram")

⊕ Better semantic understand

⊖ Computationally more challenging as dictionary dimension d is larger.

b) Term-Frequency – Inverse Document Frequency (TF-IDF)

Choose $\overline{\Phi}(x)_w = freq_w \cdot \left(\log\left(\frac{d}{N_w}\right) + 1\right)$

$freq_w$: frequency of word $w$ in document $x$

$d$: total number of words

$N_w$: nr. of documents containing word $w$.

⊕ scales down importance of words that are common across documents.