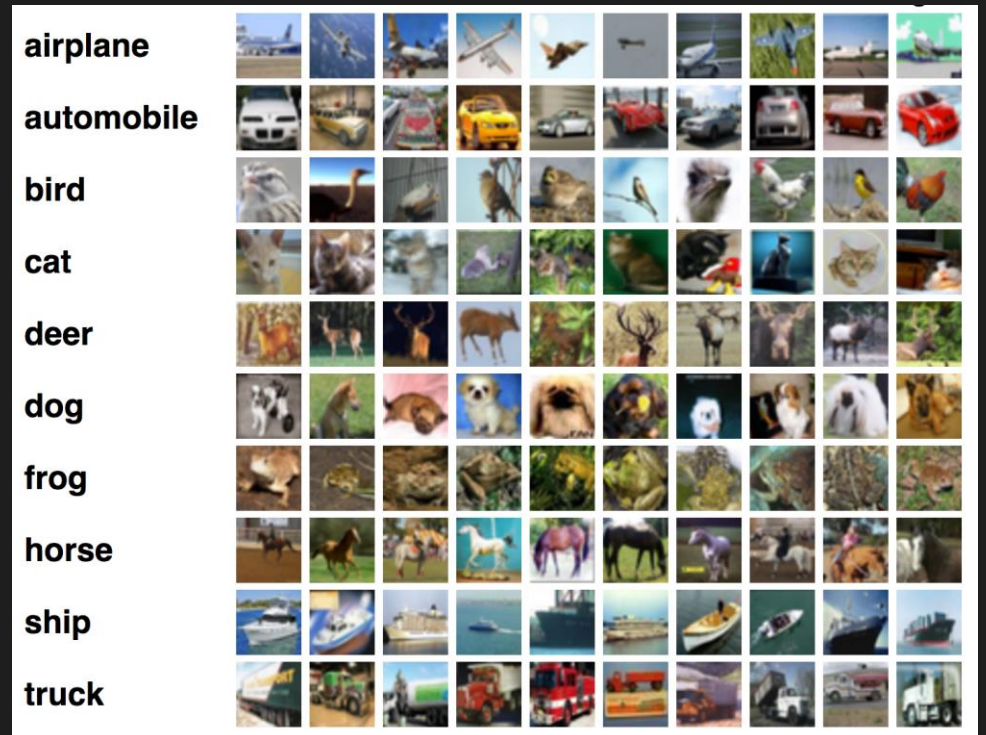# INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

## Lecture 13
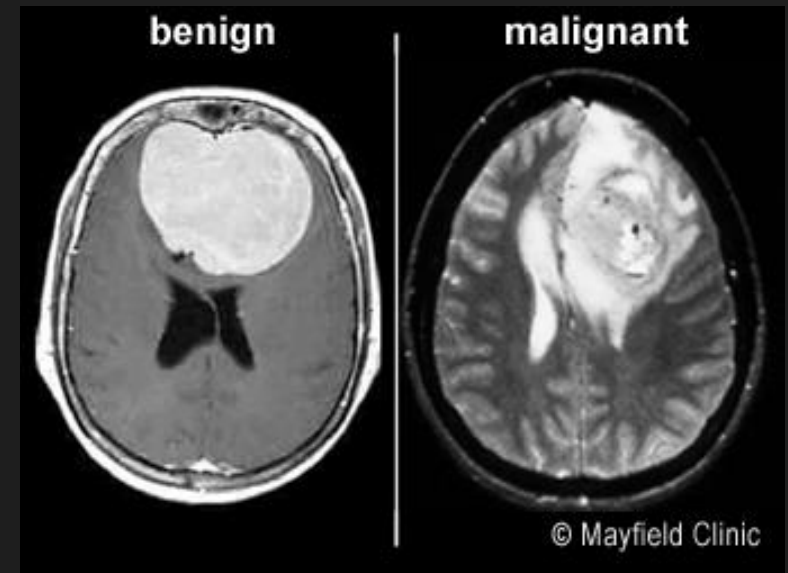
### Hassan Ashtiani

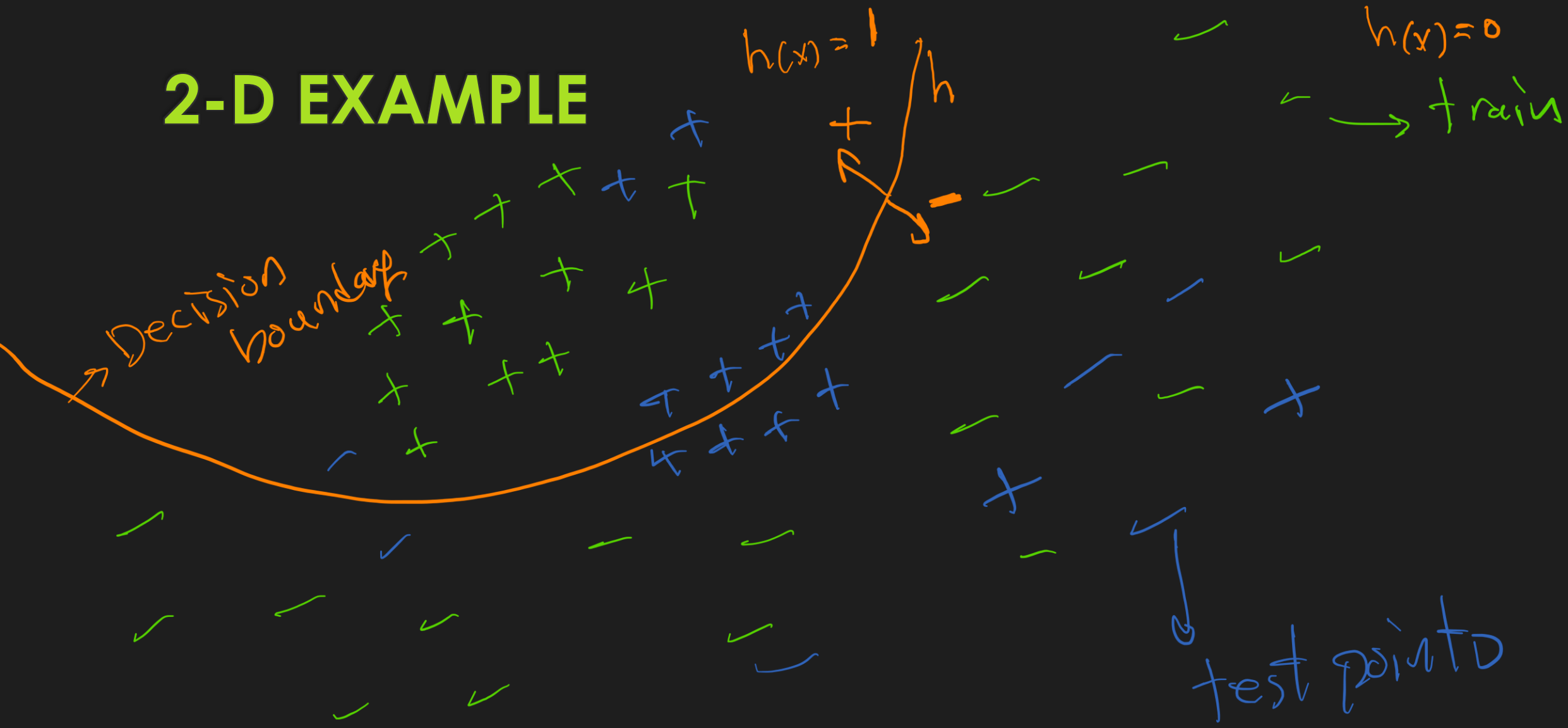# CLASSIFICATION

- PREDICT THE CATEGORY

- $k$-CLASS CLASSIFICATION
  - $y \in \{1,2,3,\ldots,k\}$

# 2-D EXAMPLE

$h(x) = 1$

$h(x) = 0$

$\rightarrow$ train
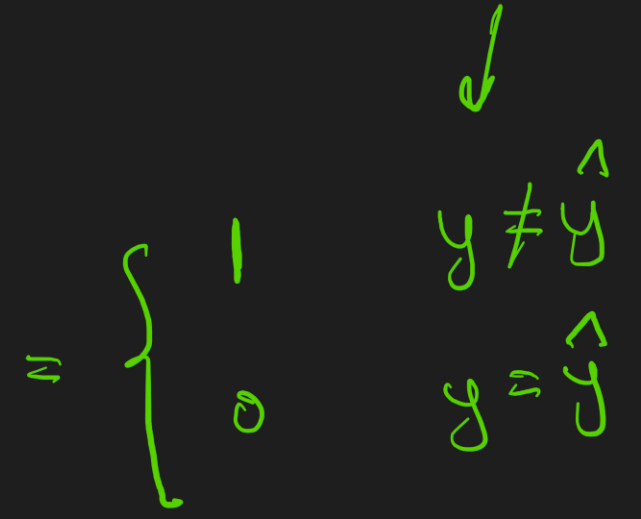
$h$

Decision boundary

test points

# ACCURACY, 0-1 LOSS

- FOR REGRESSION WE USED THE SQUARED LOSS
  - $l(y, \hat{y}) = (y - \hat{y})^2$.
- LOSS FOR CLASSIFICATION?
  - THE ZERO-ONE LOSS (ERROR): $l^{0-1}(y, \hat{y}) = 1_{y \neq \hat{y}}$ $= \begin{cases} 1 & y \neq \hat{y} \\ 0 & y = \hat{y} \end{cases}$
- PREDICTOR/LABELING-FUNCTION: $h : X \to Y$
- TRAINING SET: $Z = \left((x^1, y^1), \dots, (x^m, y^m)\right)$
  - GENERATED FROM $D_{XY}$, OR $D$ FOR SHORT
- EMPIRICAL (TRAINING) ERROR OF $h$: $L_Z^{0-1}(h) = \frac{1}{m} \sum_i \ell^{0-1}(h(x^i), y^i)$
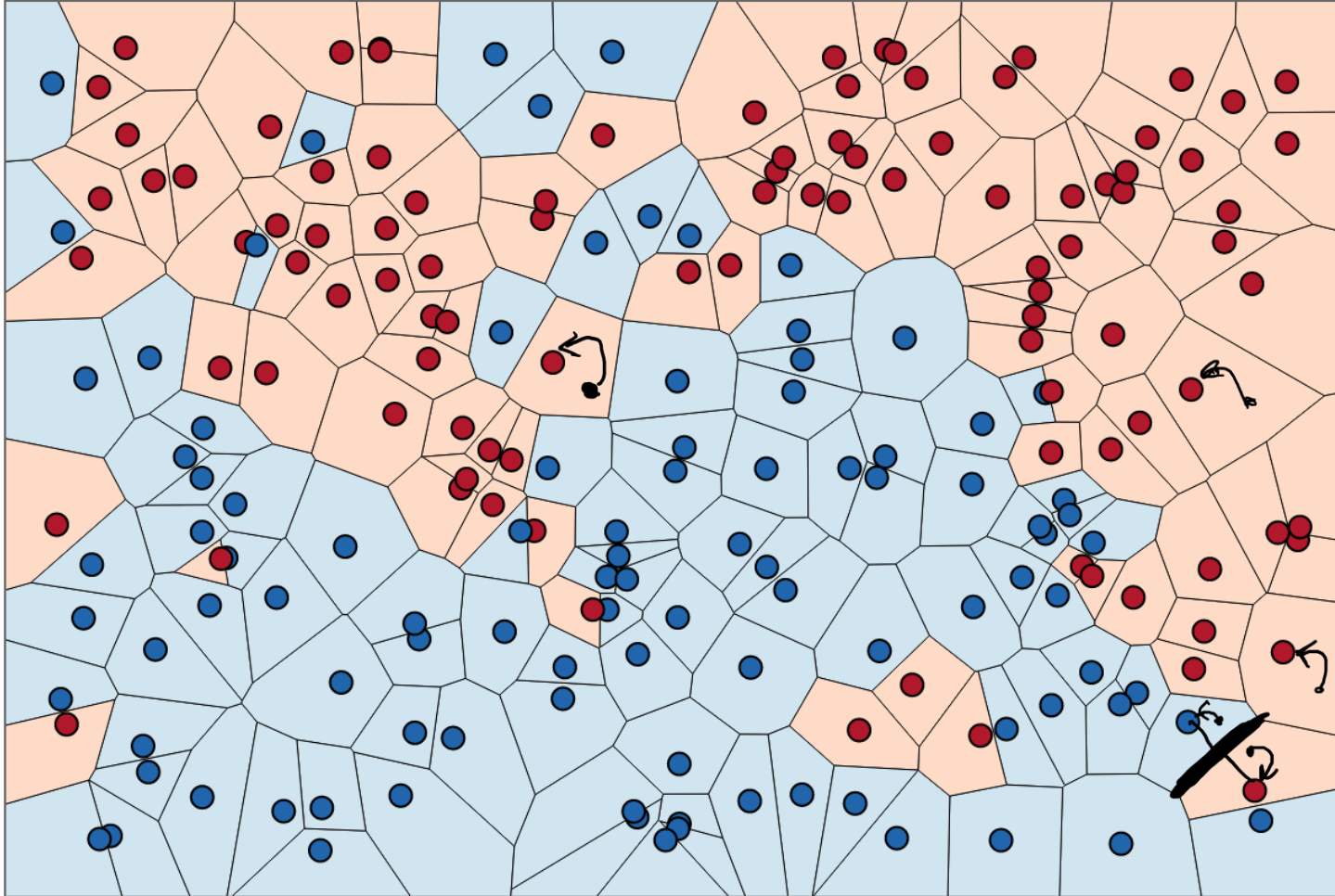- EXPECTED ERROR OF $h$: $L_D^{0-1} = \mathbb{E}_{(x,y) \sim D} \ell^{0-1}(h(x), y)$

# THE NEAREST NEIGHBOR CLASSIFIER

- $\hat{y}(x; Z) =$
  - FIND THE CLOSEST $x'$ TO $x$ IN THE DATA SET
    - $\min\limits_{x'} \|x - x'\|_2$
  - OUTPUT THE LABEL OF $x'$
- DECISION BOUNDARY?

# VORONOI DIAGRAM

# NEAREST NEIGHBOR: PROS AND CONST

- Pros
  - Basically no training is needed.
  - No parameter or hyper-parameter
  - Easy to implement
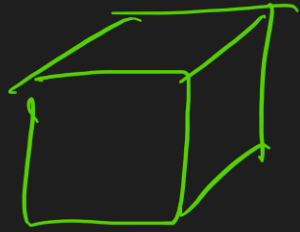  - Powerful and flexible
    - non-parametric!
- Cons
  - High test-time computational complexity, memory intensive
  - Curse-of-dimensionality!
- Can we use nearest neighbor for regression?
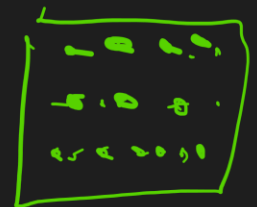
# CURSE OF DIMENSIONALITY FOR NN

- Assume points are in the $d$-dimensional unit cube

    - How many training points do I need to "cover" this cube?

        - E.g., for any $x \in [0,1]^d$, I want to have at least one training point $x^i$ such that $\left\| x - x^i \right\|_2 < 0.1$

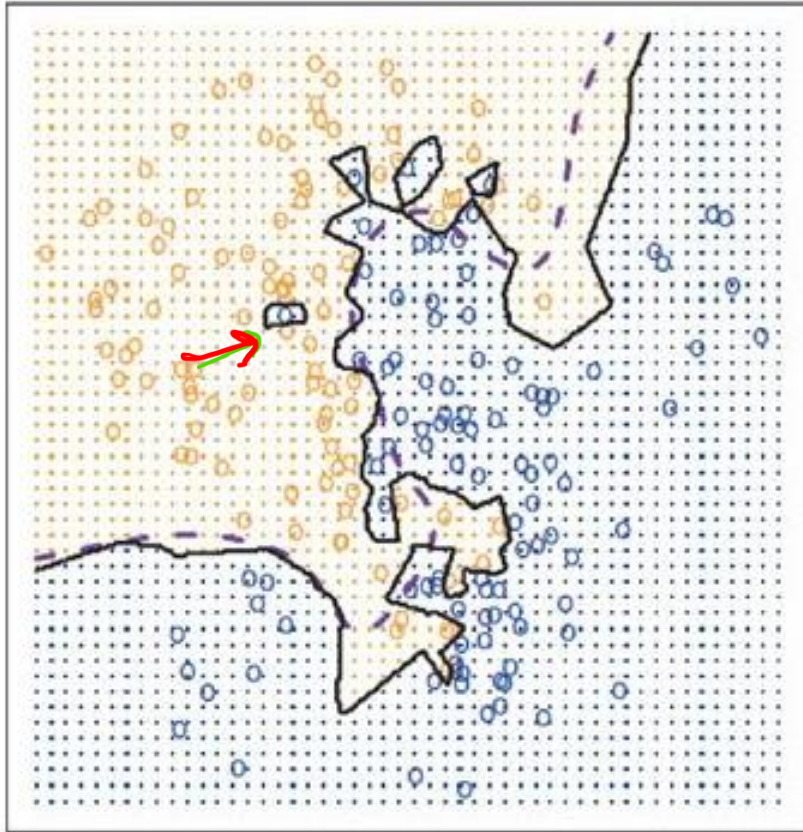    - We need $10^d$ training points!

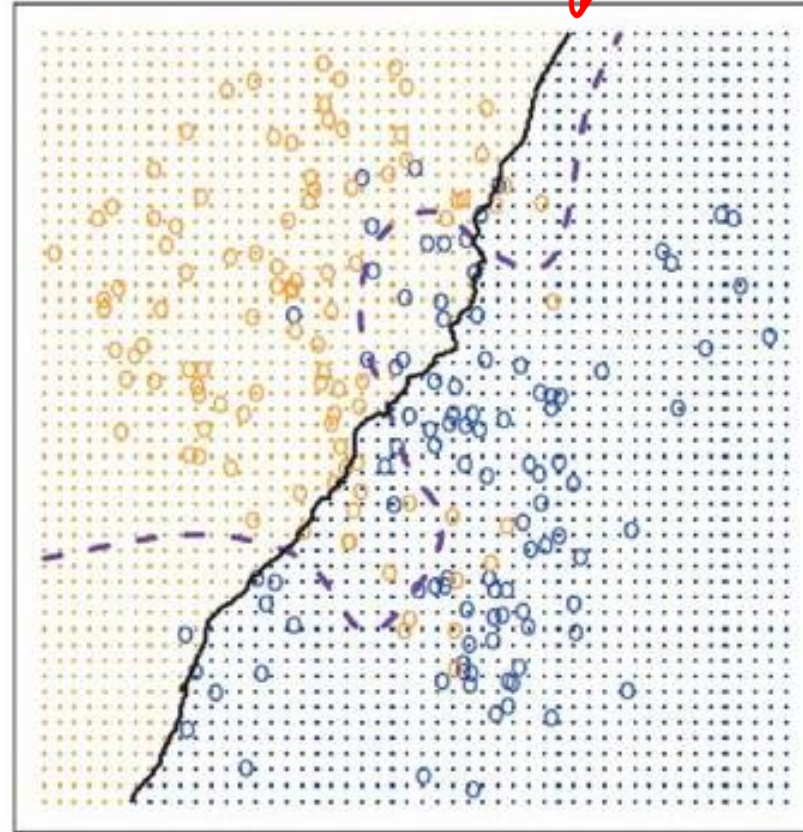- Nearest neighbor does not work well for high-dimensional data

# K-NEAREST NEIGHBOR

- $\hat{y}(x; Z) =$
  - FIND THE $k$ CLOSEST POINTS TO $x$ IN THE DATA SET
  - LET $y^1, y^2, \ldots, y^k$ BE THE LABELS OF THESE NEIGHBOR POINTS
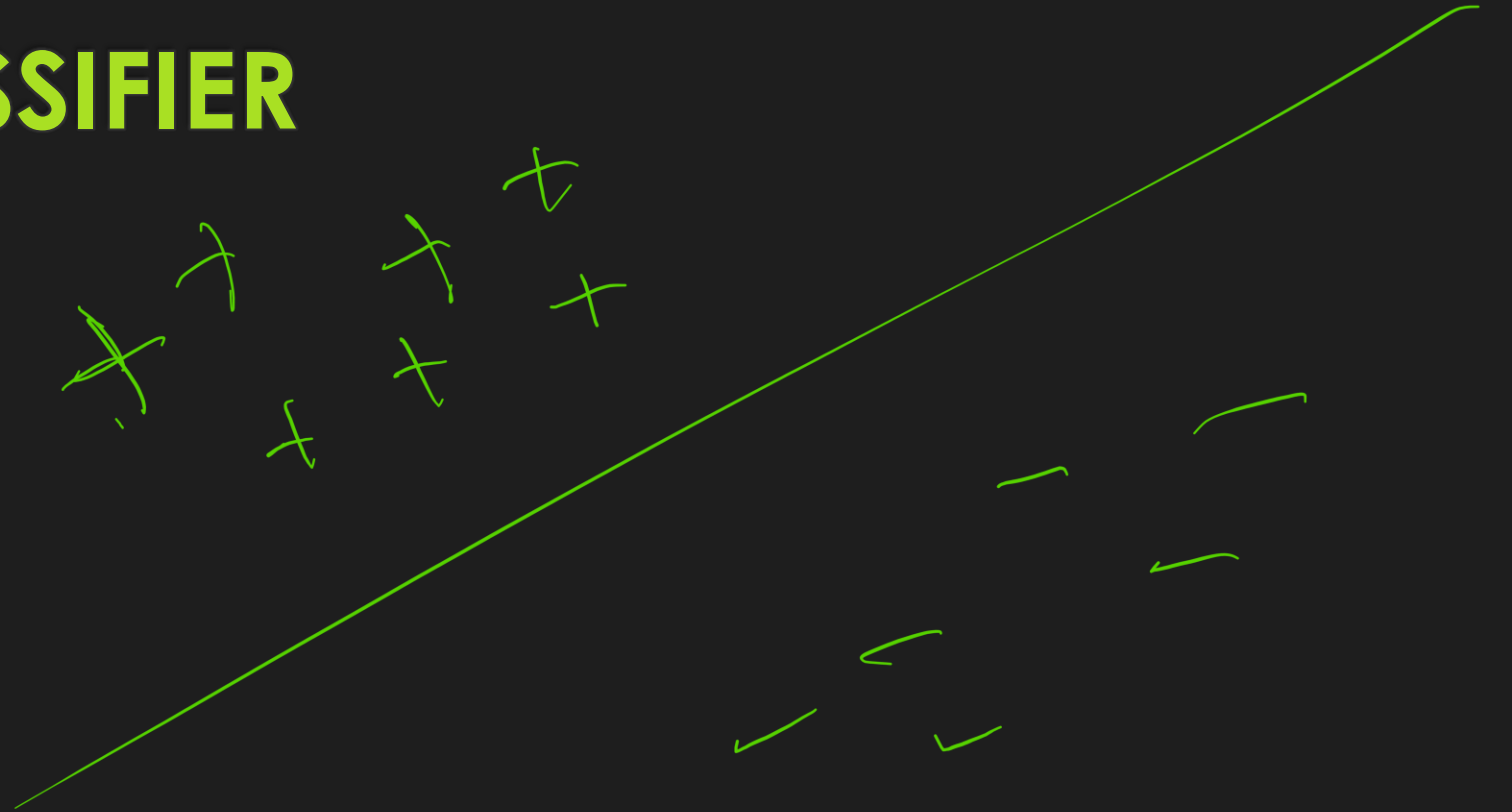  - OUTPUT THE **MAJORITY VOTE** AMONG $y^i$S

- MORE ROBUST, SMOOTHER DECISION BOUNDARY, HIGHER-TRAINING ERROR BUT LESS PRONE TO OVERFITTING
  - STILL CURSE OF DIMENSIONALITY

# LINEAR CLASSIFIER

- $\hat{y}_W(x) = sign(W^T x) = 1_{W^T x \geq 0}$
- HOW TO FIND $W$ GIVEN DATA SET $Z$?

$Y = \{0, 1\}$

$Y = \{-1, +1\}$

# LINEAR CLASSIFIERS

- EMPIRICAL RISK MINIMIZATION

  - $\widehat{W} = arg\min\limits_{w} L_Z^{0-1}(\hat{y}_w) = arg\min\limits_{W} \frac{1}{m} \sum \ell\left(\hat{y}_w(x^i), y^i\right)$

- COMPUTATIONAL COMPLEXITY FOR $d = 2$?

$O(m^2)$ time

- HIGHER $d$?

$m^d$

# HARDNESS OF LINEAR CLASSIFICATION

- In general, finding the linear separator with minimum classification error is NP-hard (with respect to $d$)

- Unless….

  - Data is linearly separable!

  - ..or optimize another loss function instead!

- Can't we just go into higher dimensions to make the data linearly separable?

# LINEARLY SEPARABLE DATA

# SURROGATE LOSS FUNCTIONS

- So far, we assumed the classifier returns a discrete value
  - E.g., $\hat{y}_w = sign(W^T x) \in \{0,1\}$
- What if the classifier's output is continuous
  - E.g., $\hat{y}_w = W^T x$
  - $\hat{y}$ will capture the "confidence" of the classifier too
- Other (continuous) loss functions
  - Margin loss, cross-entropy/negative-log-likelihood loss, ...
- Potential benefits?
  - Easier to optimize
  - Mitigate Overfitting