



**GHENT
UNIVERSITY**

DATA ANALYTICS: MATHEMATICAL BACKGROUND AND ALGORITHMS

Tijl De Bie

TAKE-AWAY MESSAGE

**Mathematics is key
to data analytics**

**Search:
algorithms and optimization**

**Assessment:
statistical learning theory**

MACHINE LEARNING

Data
Credit card transaction properties
Patient symptoms
Movie properties
Customer properties
...

x_i



Label
Fraudulent or not
Diagnosis
Rating
Churn probability
...

y_i

Given n pairs (x_i, y_i) , infer a general rule

MACHINE LEARNING

- **Hypothesis space** \mathcal{H}
- Find $f \in \mathcal{H}$ such that $f(\mathbf{x}) \approx y$ most of the time
- Ideal approach:
 - Define **loss function**: $L(f(\mathbf{x}), y)$
 - **Risk** = expected loss given f :
$$R(f) = \mathbb{E}\{L(f(\mathbf{x}), y)\}$$
 - Find $f \in \mathcal{H}$ minimizing the **Risk**:

$$\min_{f \in \mathcal{H}} R(f)$$

EMPIRICAL RISK MINIMIZATION

- $R(f) = \mathbb{E}\{L(f(\mathbf{x}), y)\}$ is unknown
- Can be estimated!

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i)$$

- **Empirical Risk Minimization:**

$$\min_{f \in \mathcal{H}} \hat{R}(f)$$

$$\min_{f \in \mathcal{H}} \hat{R}(f)$$

Search

- How to find f efficiently?
- How does this depend on \mathcal{H} ?

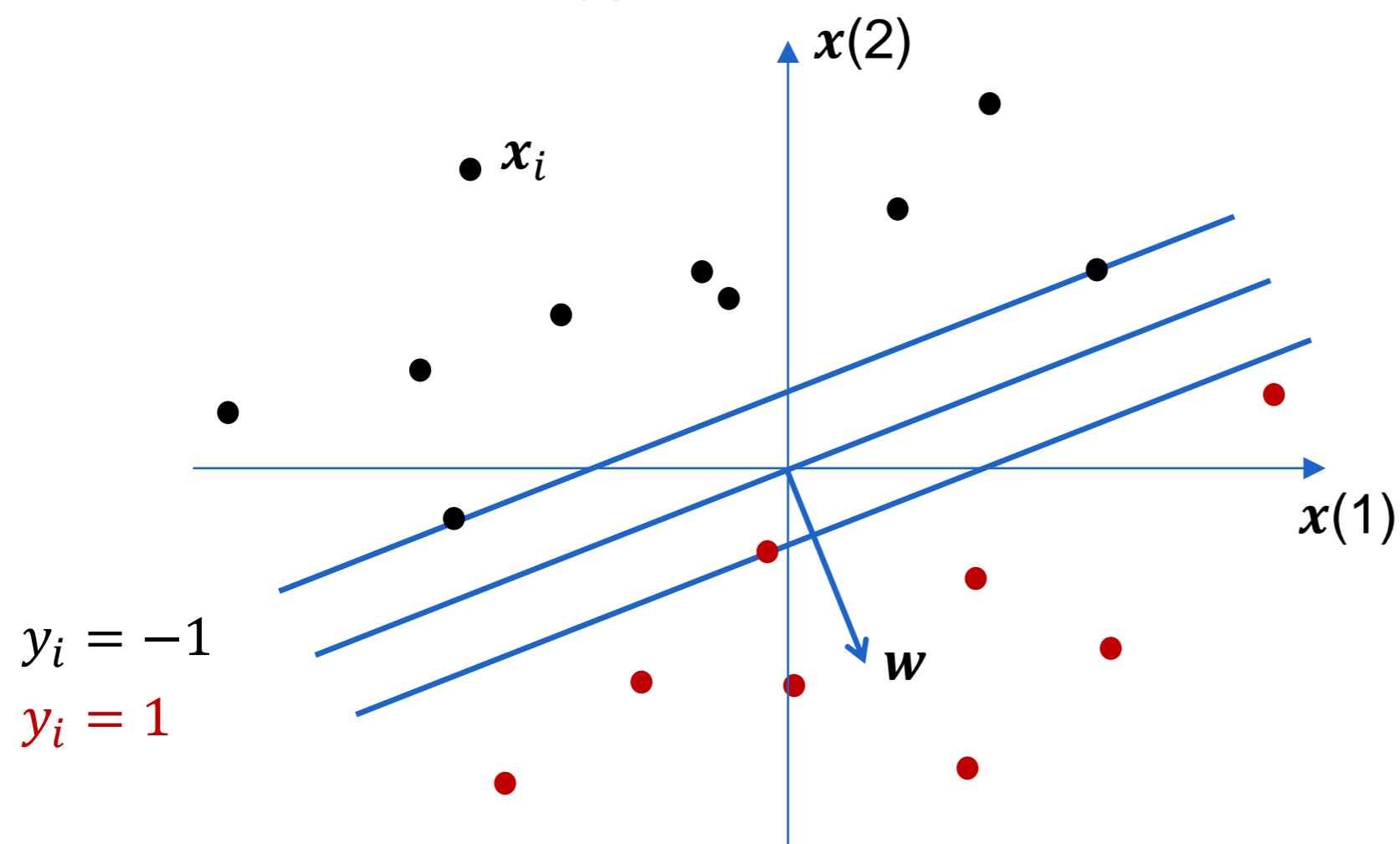
Assessment

- Given $\hat{R}(f)$, what about $R(f)$?
- How does this depend on \mathcal{H} ?

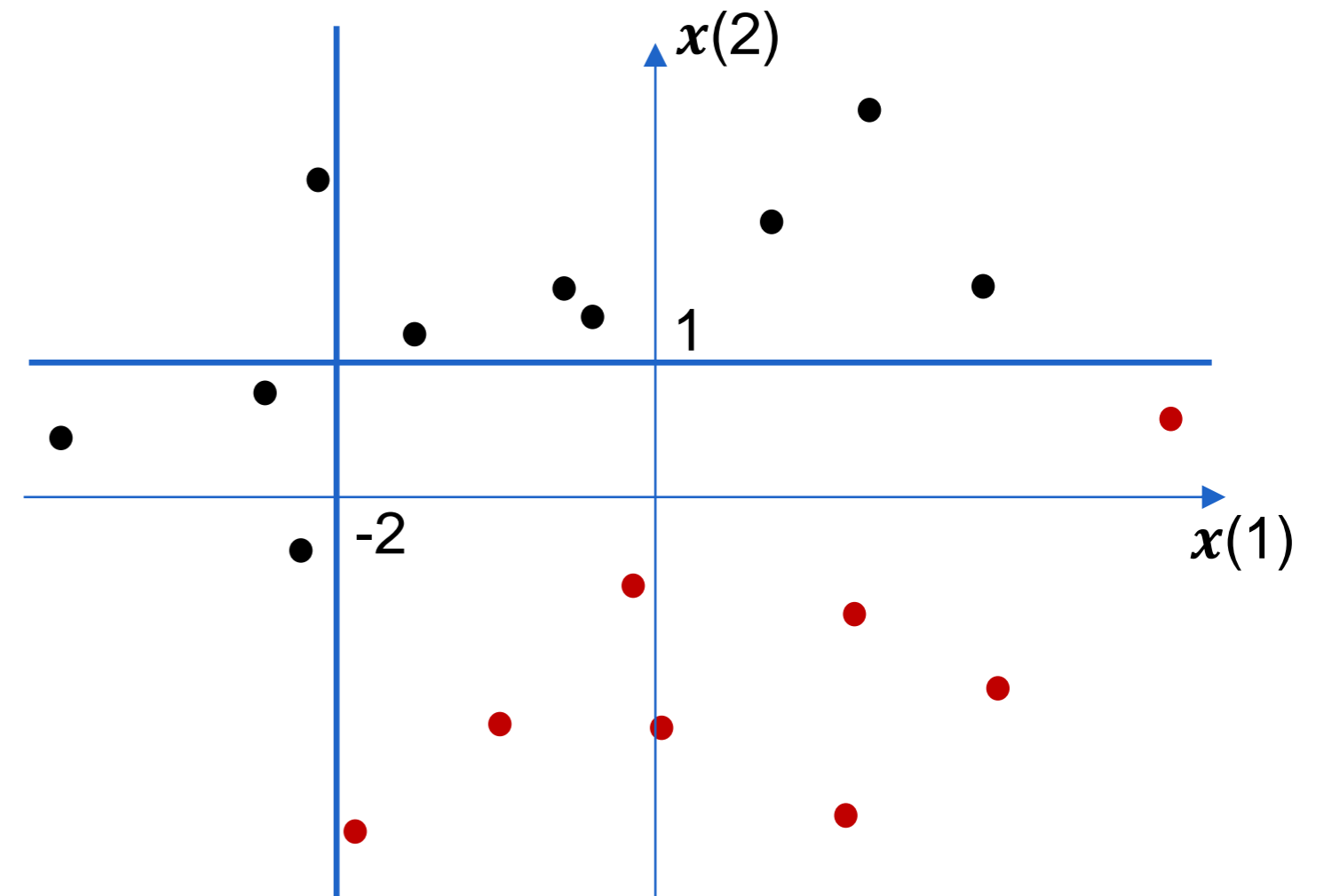
SEARCH

ALGORITHMS & OPTIMIZATION

Support Vector Machines



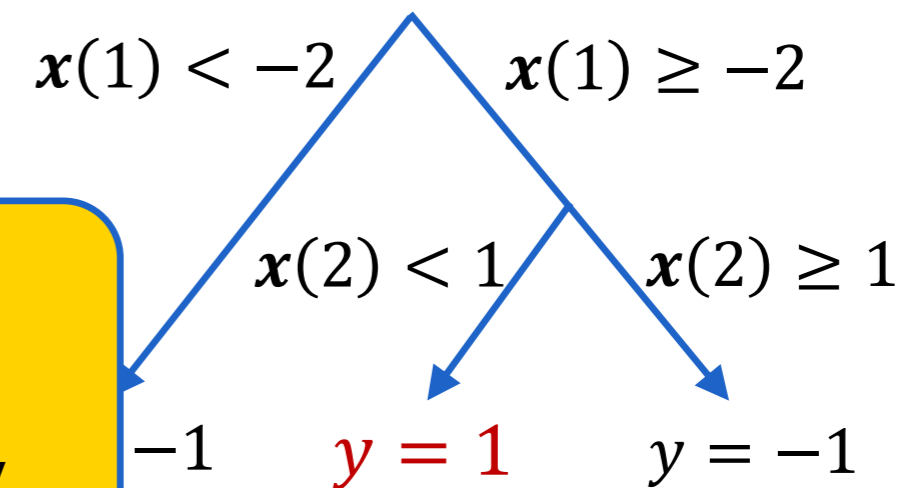
Decision Trees



$$\max_w \frac{2}{\mathbf{w}^T \mathbf{w}}$$

subject to: $y_i (\mathbf{x}_i^T \mathbf{w}) \geq 1$

Optimization theory
Algorithms
Computational complexity theory



ASSESSMENT

STATISTICAL LEARNING THEORY

- Given f , statements of the following type:

$$P(R(f) > \hat{R}(f) + \varepsilon) < \delta(\varepsilon)$$

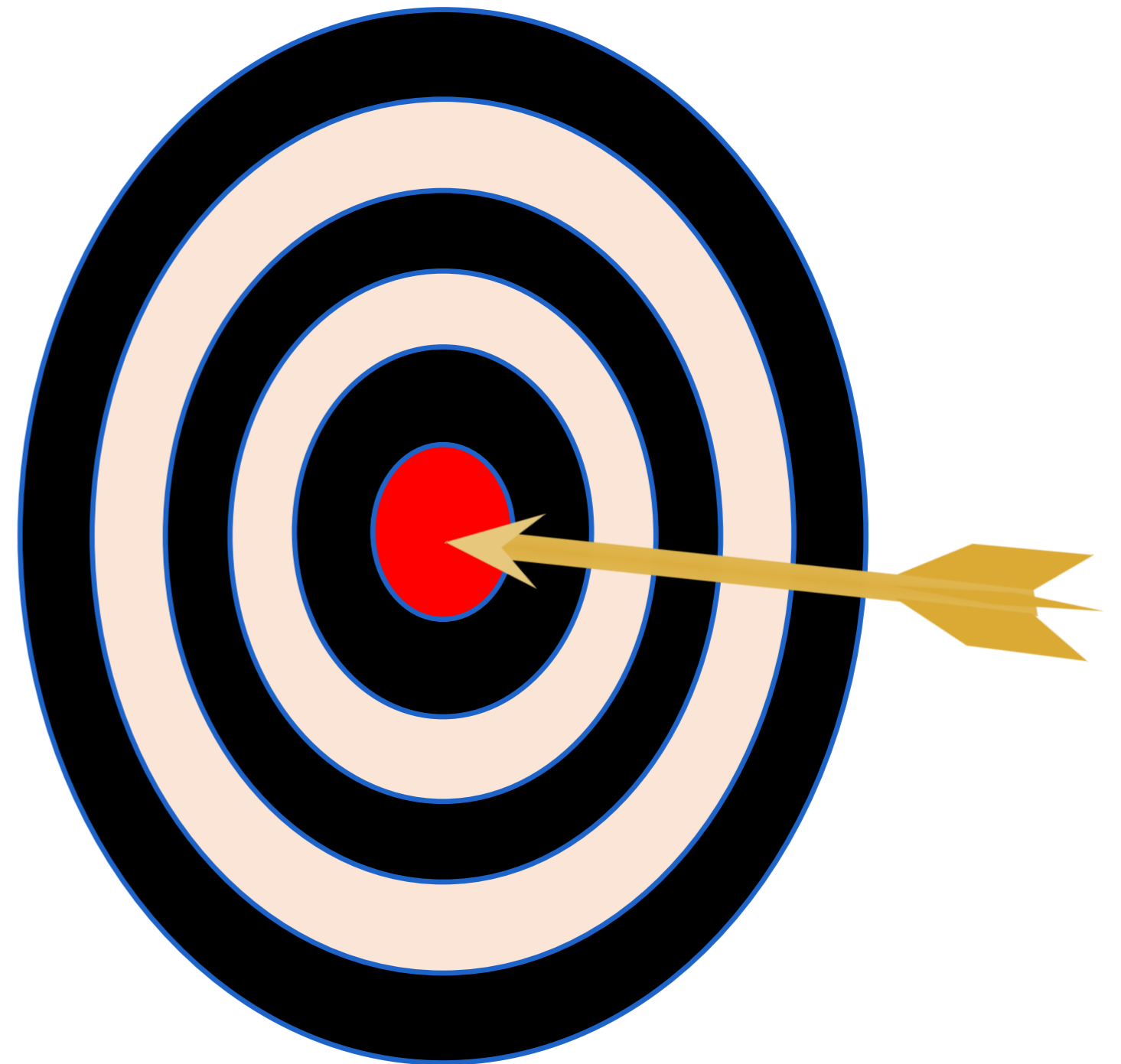
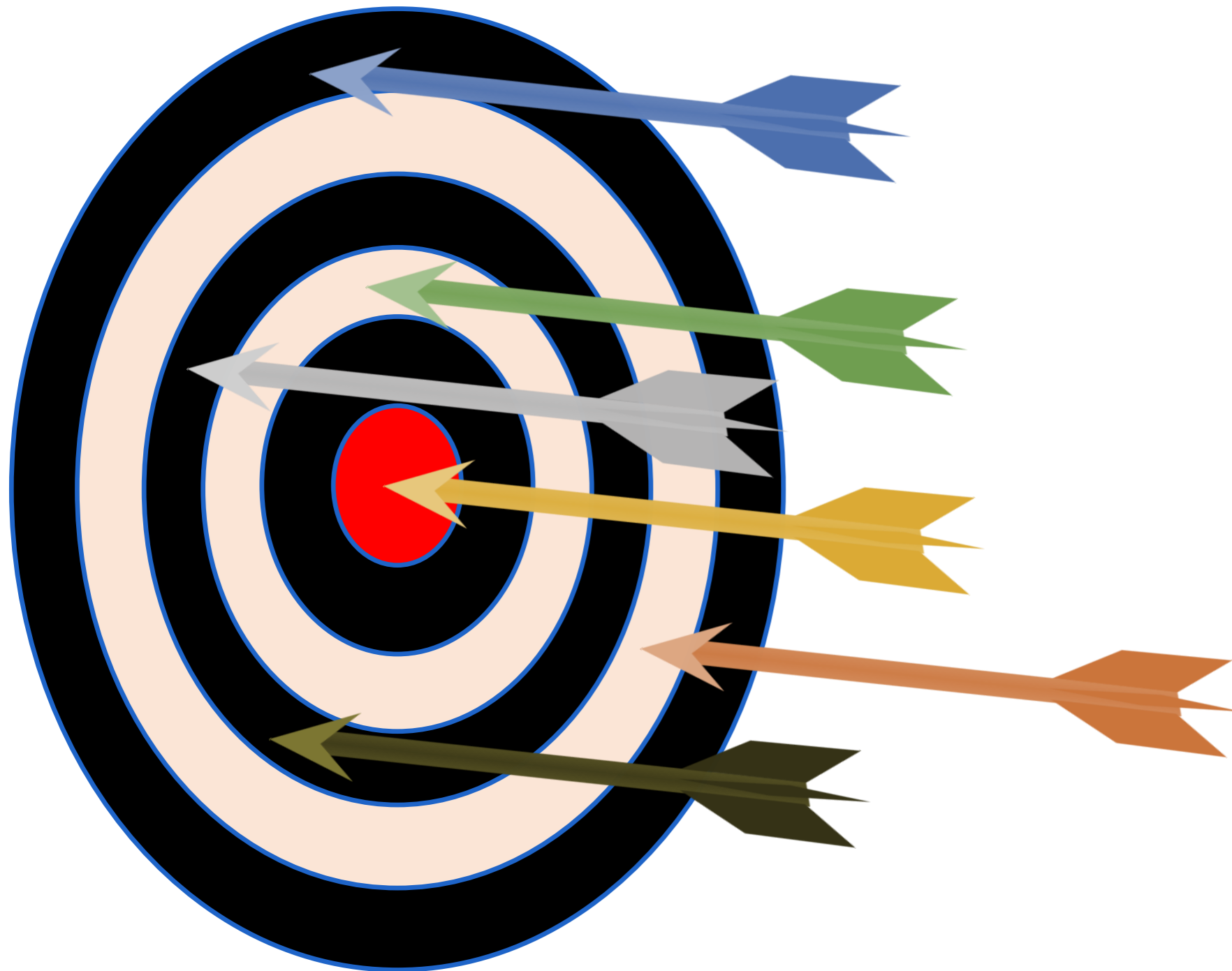
- Typically, $\delta(\varepsilon)$ decreases rapidly with increasing ε
 - For a given f , easy to do!
-
- However, there is a problem...

THE PROBLEM

- **Friend A:** “Last night I spotted a European roller (scharrelaar) in my garden!”
- **Friend B:** “I’m not impressed, I’ll check tonight and I’m sure I’ll spot one in my garden too.”
- Lo and behold, friend B does spot a European roller...
- **Friend C** would like to see a European roller too. Which garden should she try?



THE PROBLEM



FIXING THE PROBLEM

- Quick fix:

$$P(R(f) > \hat{R}(f) + \delta) < |\mathcal{H}| \varepsilon(\delta)$$

- Not very tight bound
- Fails when $|\mathcal{H}| = \infty$

- **Statistical Learning Theory:**

$$P(R(f) > \hat{R}(f) + \delta) < C_{\mathcal{H}} \varepsilon(\delta)$$

TAKE-AWAY MESSAGE

**Mathematics is key
to data analytics**

**Search:
algorithms and optimization**

**Assessment:
statistical learning theory**

**Both theories inform
how to choose \mathcal{H}**

Tijl De Bie

Professor

IDLAB – IMEC
ELECTRONICS AND INFORMATION SYSTEMS

E tijl.debie@ugent.be

T +32 9 264 33 66

M +32 468 29 69 74

www.ugent.be

 Ghent University

 @ugent

 Ghent University