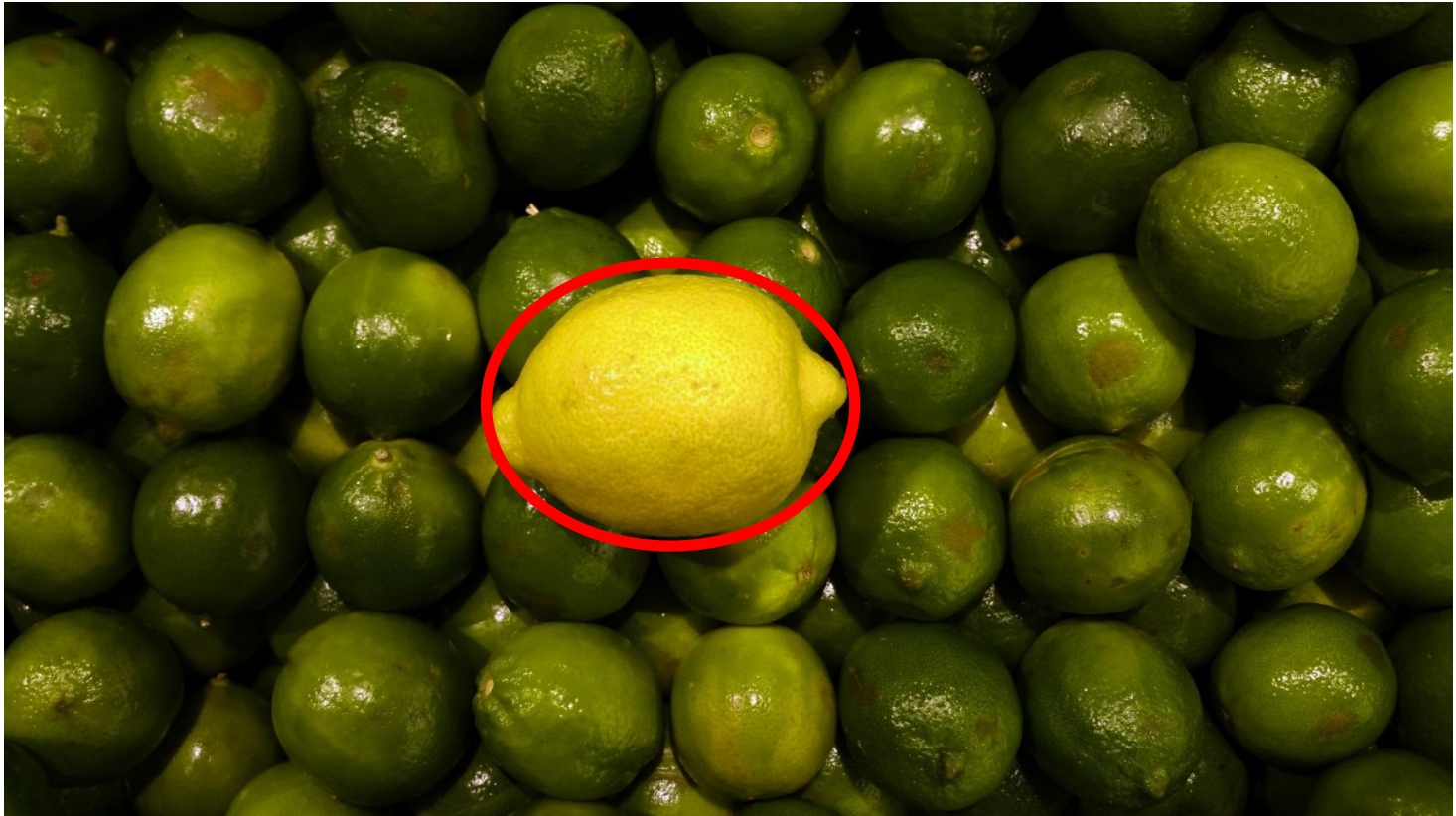# Finite Sample Complexity of Rare Pattern Anomaly Detection

Md Amran Siddiqui, Alan Fern, Thomas G. Dietterich and Shubhomoy Das

School of EECS

Oregon State University

Oregon State
UNIVERSITY
OSU

# Anomaly Detection

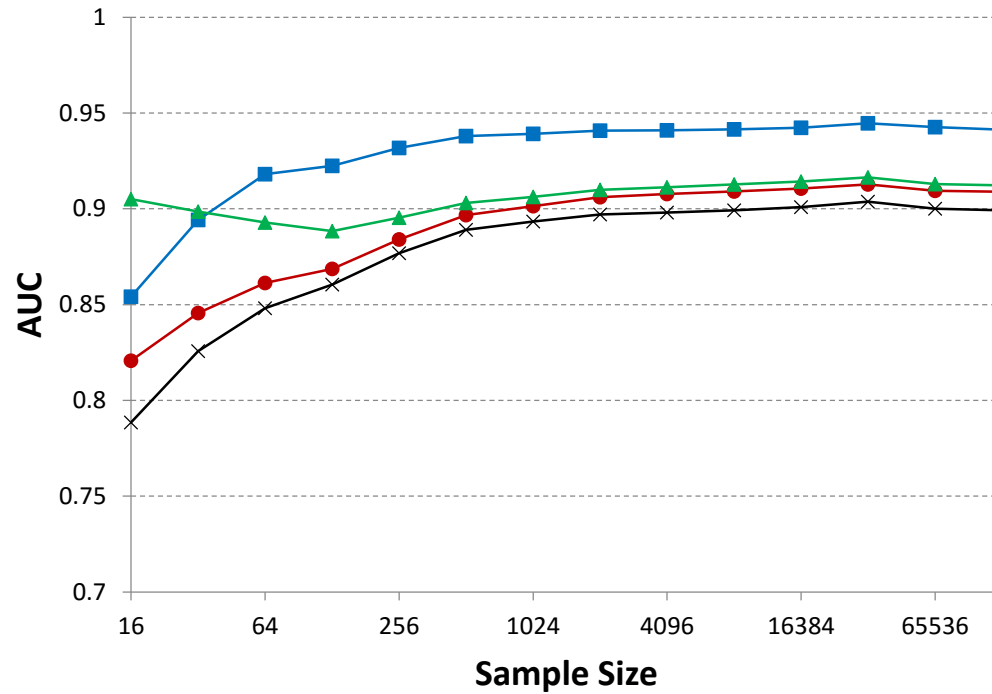- **Goal:** Identify rare or strange objects

# Challenges

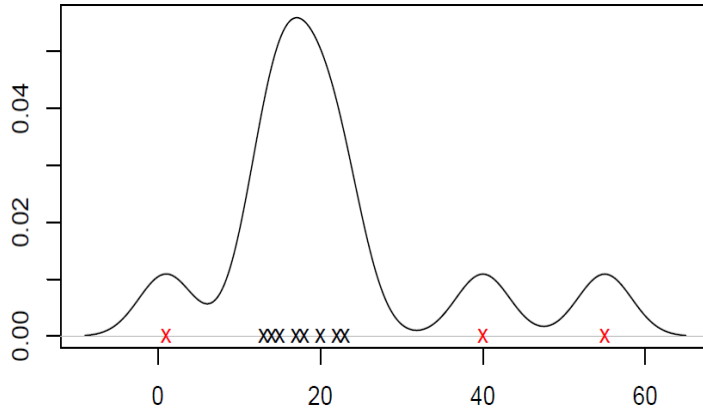- Every object is unusual in some ways!



- Anomaly detection in high-dimension seems impossible ☹
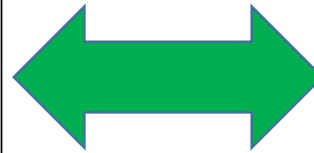
# State-of-the-art



- Often perform very well with a surprisingly small number of examples ☺

- Performance depends on:
  - ✓Sample Complexity
  - ✓Notion of Anomaly

# Notion of Anomaly



## Statistical

## Semantic
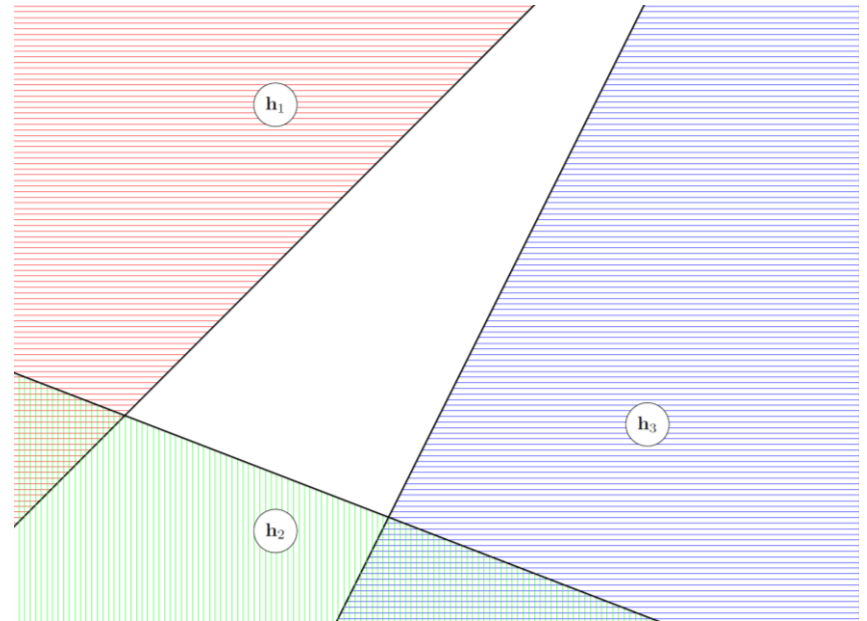
- Algorithm/representation specific

- Application specific

- Example: density of a point

- Example: A threat in security

Oregon State UNIVERSITY OSU

# Motivation

Many state-of-the-art algorithms [Chen et al. 2015, Liu et al. 2008, Wu et al. 2014, Tomas Pevny 2016] exhibit the following steps:

1. Choose a "pattern space"
   (analogous to hypothesis space)

Oregon State UNIVERSITY OSU

# Motivation

Many state-of-the-art algorithms [Chen et al. 2015, Liu et al. 2008, Wu et al. 2014, Tomas Pevny 2016] exhibit the following steps:
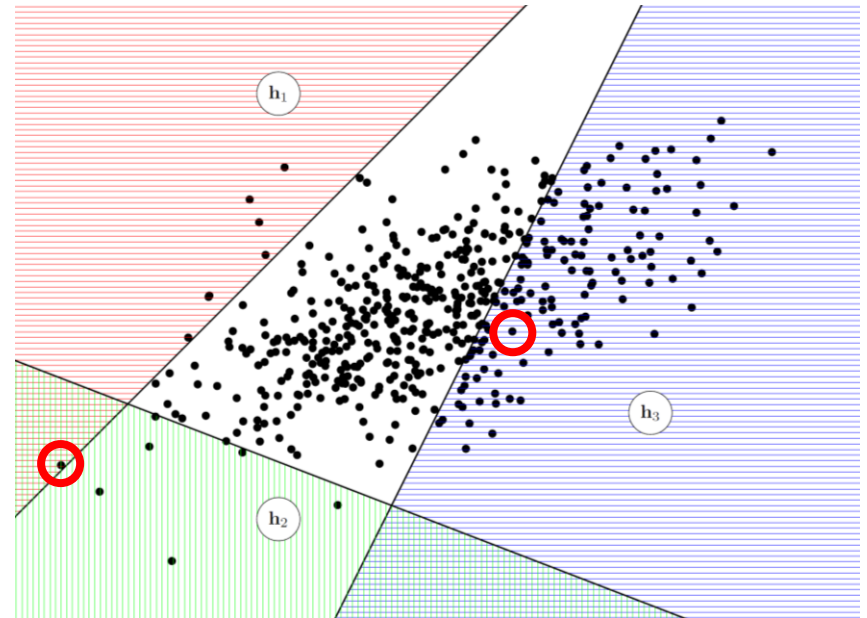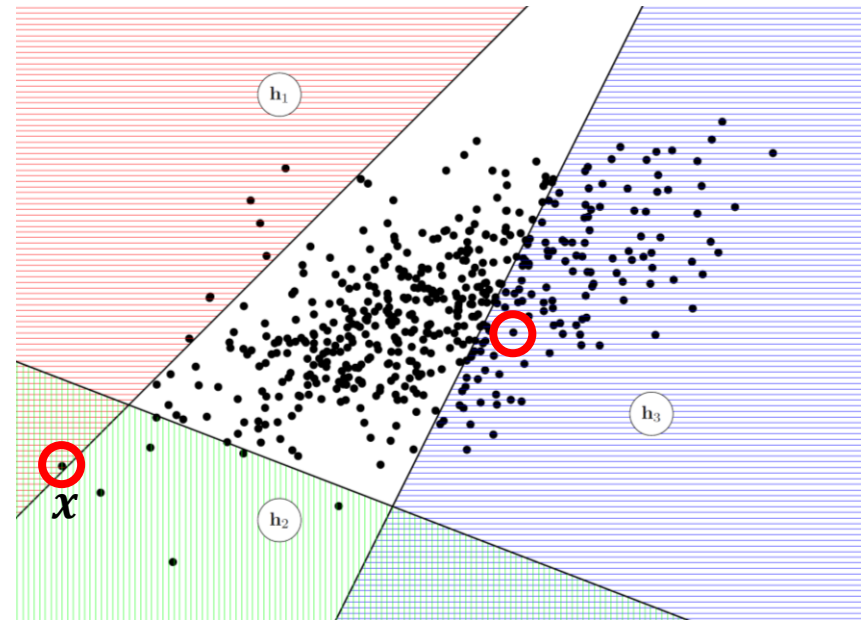
1. Choose a "pattern space" (analogous to hypothesis space)

2. Monitor the empirical frequency of the patterns

3. Compute anomaly score based on the frequencies



**Rare Pattern Anomaly Detection (RPAD)**

Oregon State UNIVERSITY OSU

# Rare Pattern Anomaly Detection (RPAD)

| | |
|---|---|
| $\mathcal{H}$ | Pattern space, $\{h_1, h_2, h_3\}$ |
| $\mathcal{H}[x]$ | Set of patterns that contain $x$, $\{h_1, h_2\}$ |
| $f(h)$ | Frequency of a pattern $h$, $f(h_1) < f(h_3)$ |
| $\tau$ | Detection threshold |

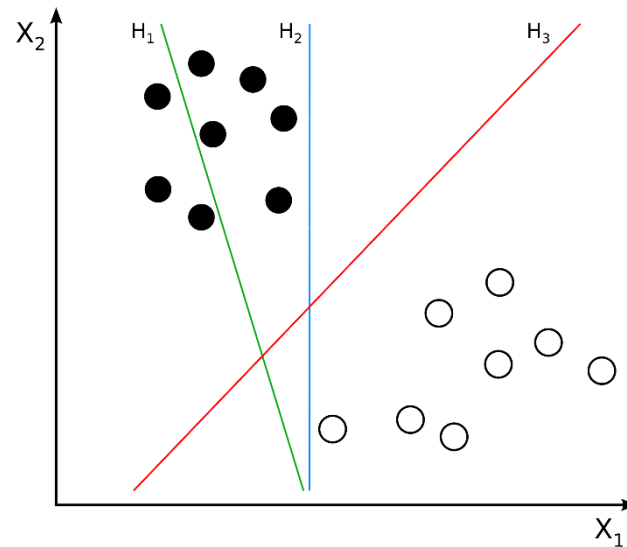| A point $x$ is | $\tau$-outlier : If $\mathcal{H}[x]$ contains an $h$ with $f(h) \leq \tau$ |
|---|---|
| | $\tau$-common : Otherwise |

# Learning Protocol

- Assumption: Input is generated from a distribution $\mathcal{P}$ i.e. $x \sim \mathcal{P}$

- Let, $\mathcal{A}$ be an anomaly detection algorithm

- $\mathcal{A}$ can draw a training set $\mathcal{D}$ of any size $\mathcal{N}$ from $\mathcal{P}$

- Given a new point $x$ : $\mathcal{A}$ has to either "detect" or "reject"

- Ideally, $\mathcal{A}$ is "correct":
    if $\mathcal{A}$ "detects" all $\tau$-outliers and
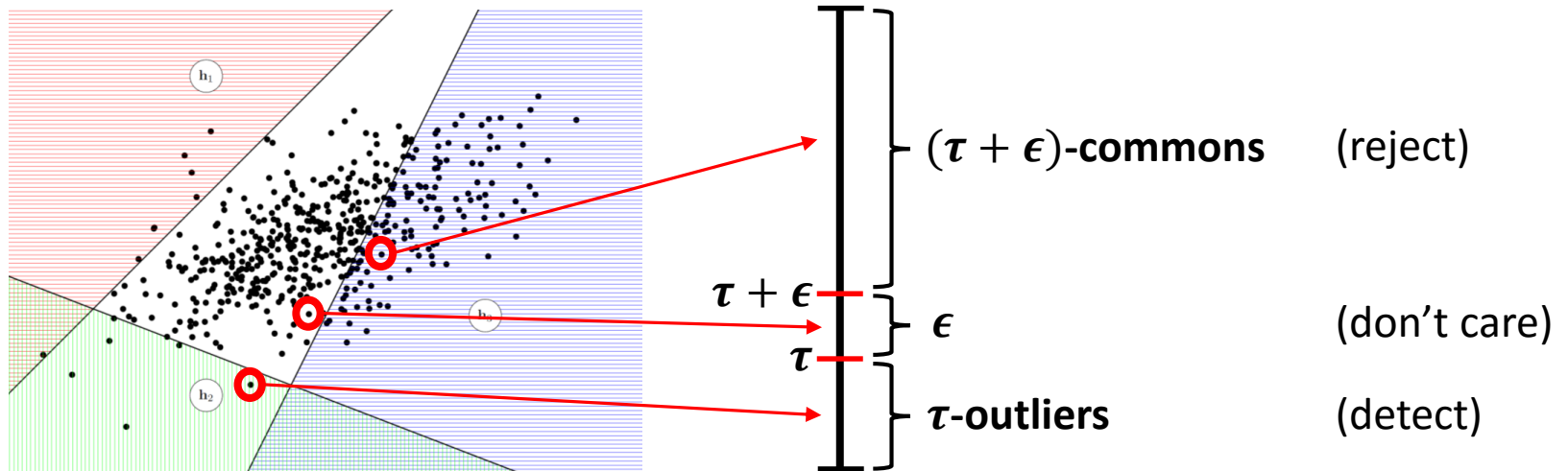        "rejects" all $\tau$-commons

# Supervised PAC Learning Framework

- Consider a hypothesis space $\mathcal{H}$ i.e. set of linear separators
- **Goal:** Learn a hypothesis that will make small error with high probability



- Sample complexity is related to the complexity of $\mathcal{H}$ : VC-dimension

- What is analogous for Anomaly Detection?

# PAC-RPAD Framework



**Definition 1.** (**PAC-RPAD**) Detection algorithm $\mathcal{A}$ is PAC-RPAD if for any $\mathcal{P}$ and any $\tau$, with probability at least $1 - \delta$ (over draws of $\mathcal{D}$), $\mathcal{A}$ detects all $\tau$-outliers and rejects all $(\tau + \epsilon)$-commons.

**Sample efficient :** if $\mathcal{A}$ draws polynomial (in $d$, $\frac{1}{\delta}$ and $\frac{1}{\epsilon}$) number of training examples from $\mathcal{P}$

# RAREPATTERNDETECT Algorithm

Input:

$\delta$: Probability tolerance

$\epsilon$: Error tolerance

$\tau$: Detection threshold

1. Draw a training set $\mathcal{D}$ of $\mathcal{N}(\delta, \epsilon)$ instances from $\mathcal{P}$
2. Decision Rule for any $x$:
   **"detect"**: If $x$ has a $\tau$-rare pattern
   **"reject"**: Otherwise

Is RAREPATTERNDETECT Sample efficient?

# Sample Complexity of RAREPATTERNDETECT

- For finite pattern space $\mathcal{H}$:

$$\mathcal{N}(\delta, \epsilon) = O\left(\frac{1}{\epsilon^2}\left(\log|\mathcal{H}| + \log\frac{1}{\delta}\right)\right)$$

- For infinite pattern space $\mathcal{H}$, but bounded VC-dimension $\mathcal{V}_{\mathcal{H}}$:

$$\mathcal{N}(\delta, \epsilon) = O\left(\frac{1}{\epsilon^2}\left(\mathcal{V}_{\mathcal{H}} \log\frac{1}{\epsilon^2} + \log\frac{1}{\delta}\right)\right)$$

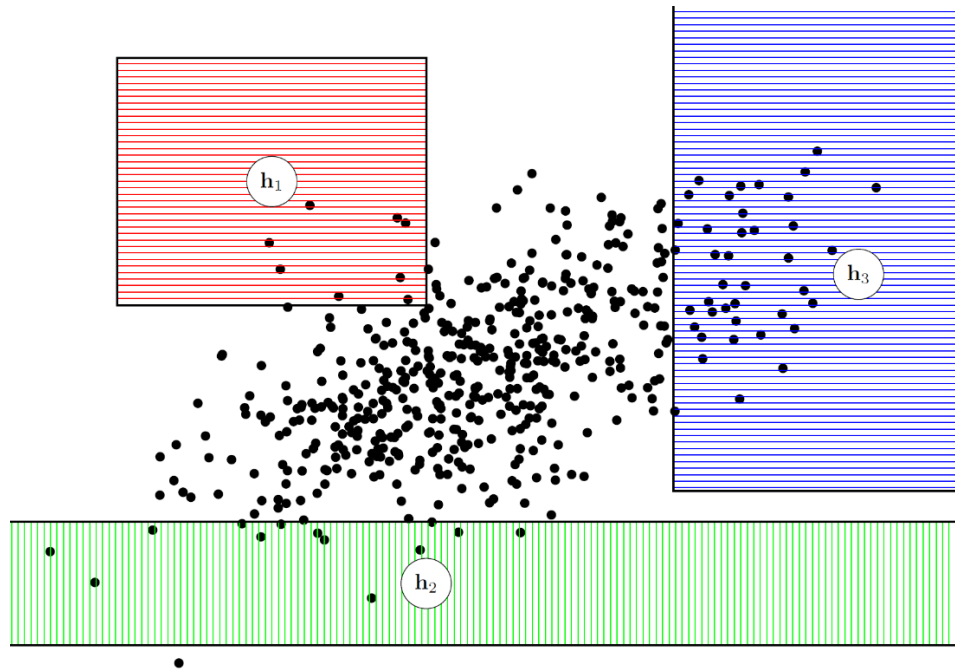- Polynomial in $\mathcal{V}_{\mathcal{H}}, \frac{1}{\delta}$ and $\frac{1}{\epsilon}$

- For the example spaces, $\mathcal{V}_{\mathcal{H}}$ are polynomial in data dimension $d$

- Hence, $\mathcal{H}$ can be learned efficiently

Oregon State OSU

# Pattern Spaces for Anomaly Detectors

- Half-spaces
  - ✓The half-space mass algorithm [Chen et al. 2015]

- Axis aligned hyper rectangle
  - ✓Isolation Forest [Liu et al. 2008] and RS-Forest [Wu et al. 2014]

- Stripes
  - ✓Light weight online detectors of anomaly (LODA) [Tomas Pevny 2016]

- Ellipsoids and shells
  - ✓Density based detectors, for example, multivariate Guassians

Oregon State **OSU**

# Axis Aligned Hyper Rectangles

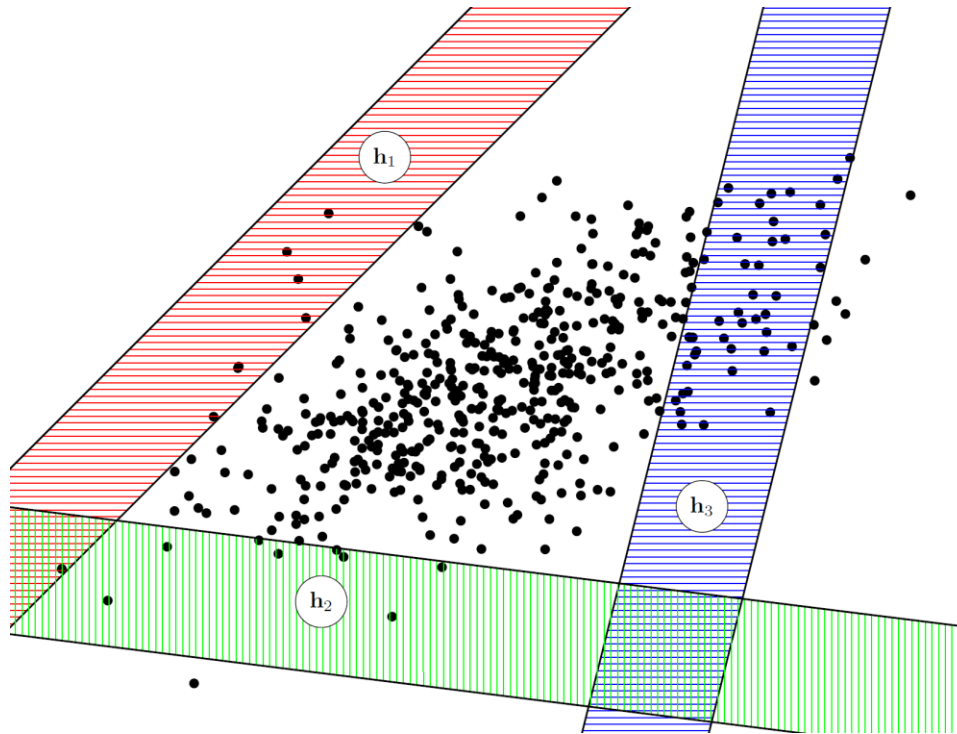- An axis aligned hyper rectangle (bounded or unbounded) is defined by $k$ boundaries in $d$-dimensional space

- Isolation Forest [Liu et al. 2008] and RS-Forest [Wu et al. 2014]



- VC-dimension = $O(d)$

# Stripes

- A stripe pattern is an intersection of two parallel half-spaces with opposite orientations

- Light weight online detectors of anomaly (LODA) [Tomas Pevny 2016]
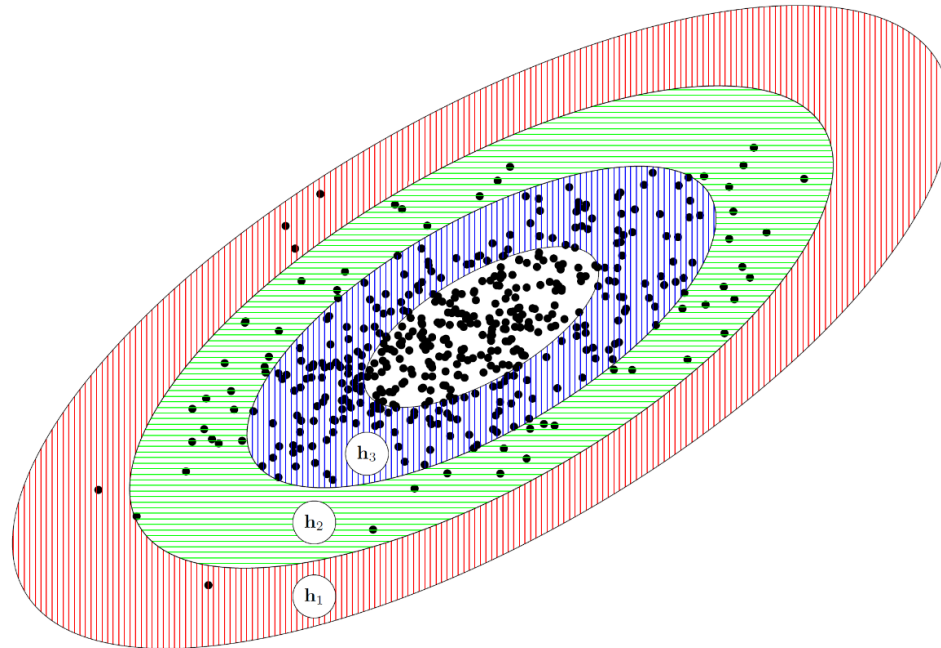


- VC-dimension = $O(d)$

# Ellipsoidal Shells

- An Ellipsoidal shell is a subtraction between two ellipsoids with same center and shape but different volumes

- Density based detectors, for example, multivariate Gaussians
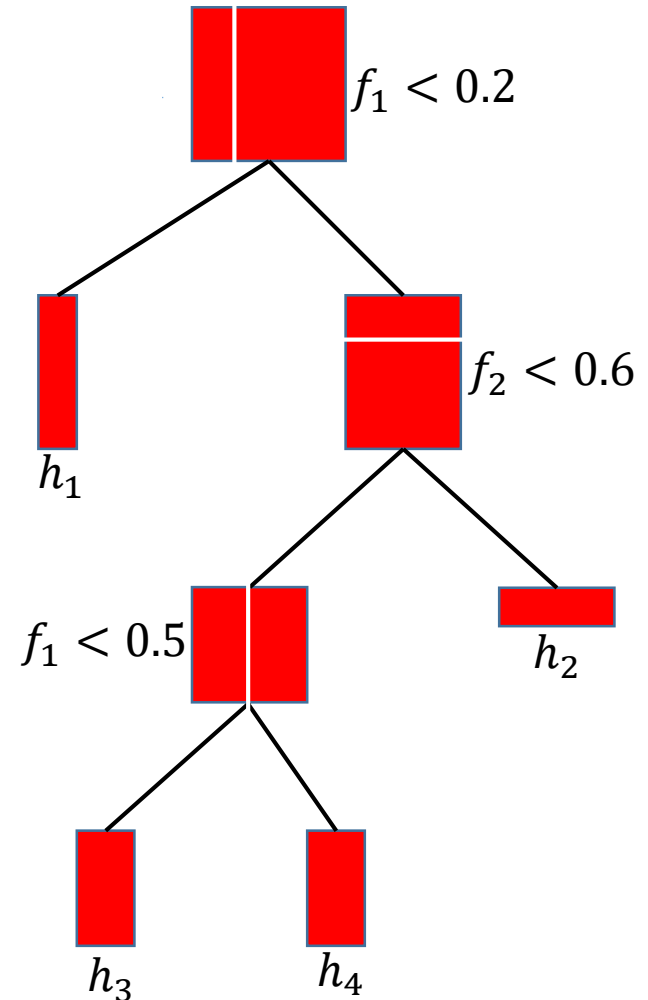


- VC-dimension = $O(d^2)$

# Experiments

- What are the qualitative properties of the learning curves of RarePatternDetect?

- Is RarePatternDetect competitive?
  - ✓ State-of-the-art anomaly detector Isolation Forest (IF)
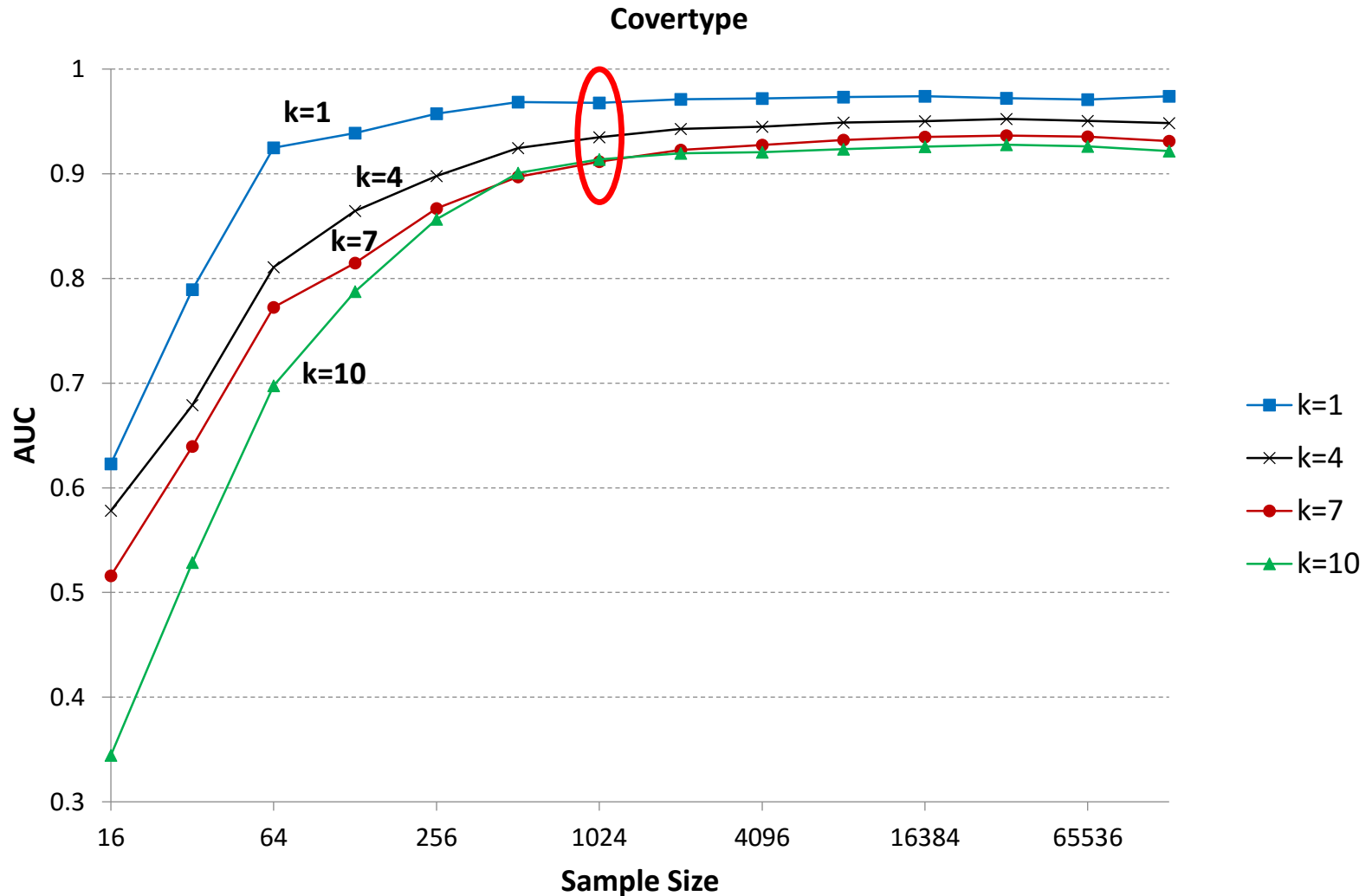  - ✓ Pattern space: axis aligned hyper rectangles

| Dataset | Dimension | # Instances | % Anomaly |
|---------|-----------|-------------|-----------|
| Covertype | 10 | 286K | 0.9% |
| Particle | 50 | 130K | 5% |
| Shuttle | 9 | 58K | 5% |

Oregon State UNIVERSITY OSU

# Pattern Space Generation

- Construct a forest of 250 random decision trees

- Each internal node is a threshold test on a feature

- Each tree node is a pattern i.e. an axis aligned hyper rectangle

- depth $(k)$ of the node determines the complexity of the pattern

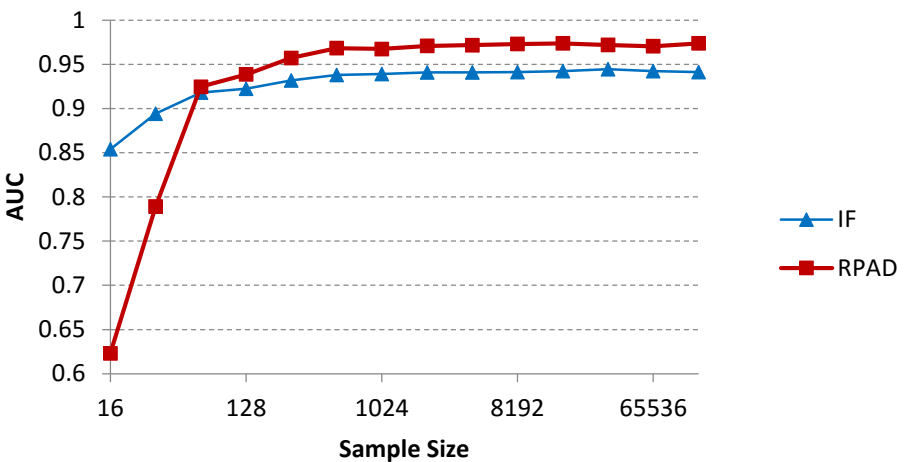- $\mathcal{H}_k$ : Set of patterns up to $k$ threshold tests, for example, $\mathcal{H}_2 = \{h_1, h_2\}$

$f_1 < 0.2$

$f_2 < 0.6$

$h_1$

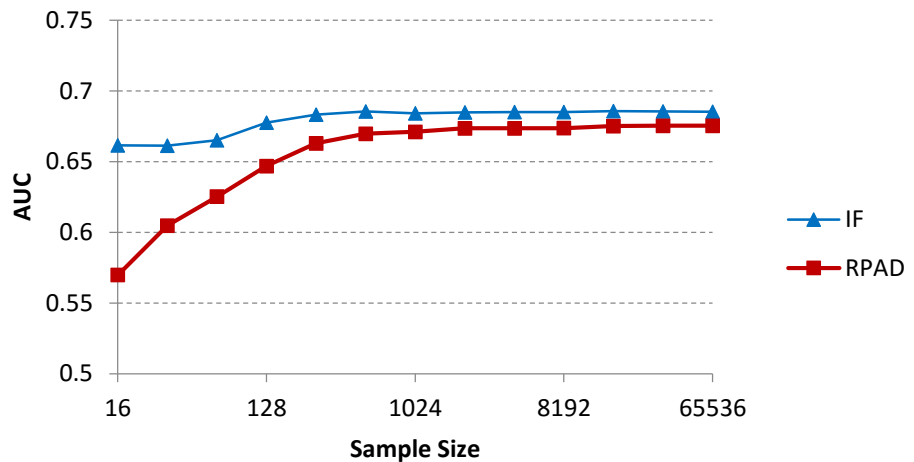$f_1 < 0.5$

$h_2$

$h_3$

$h_4$

# RAREPATTERNDETECT Learning Curve



Covertype

# Comparison

# Summary

- We developed a PAC framework to better understand the sample complexity of modern anomaly detection

- To the best of our knowledge, this is the first study of empirical learning curves for anomaly detection

- A simple PAC-RPAD algorithm is competitive with a state-of-the-art algorithm

Oregon State UNIVERSITY OSU

# Questions?

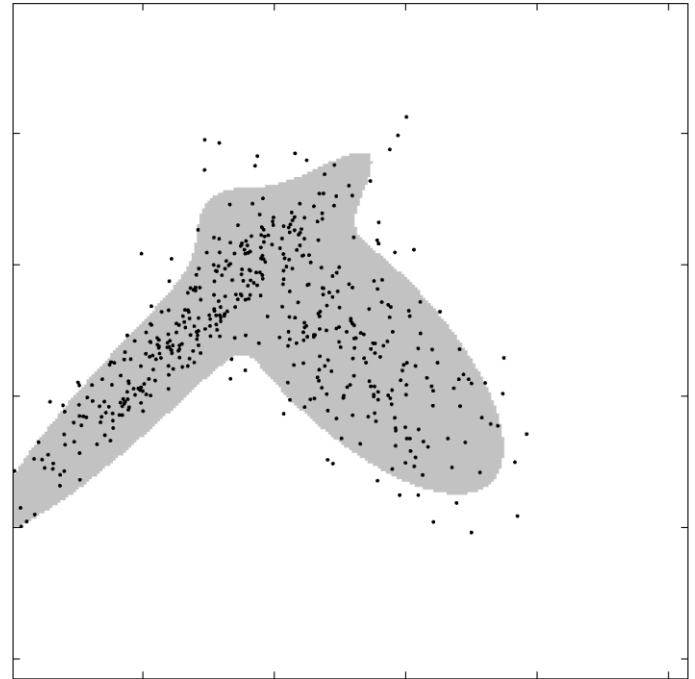Oregon State **OSU**

# Extra Slides

Oregon State OSU

# Prior Work

- Sample Complexity for Anomaly Detection:

  ✓ One Class SVM (Scholkopf et al. 2001)
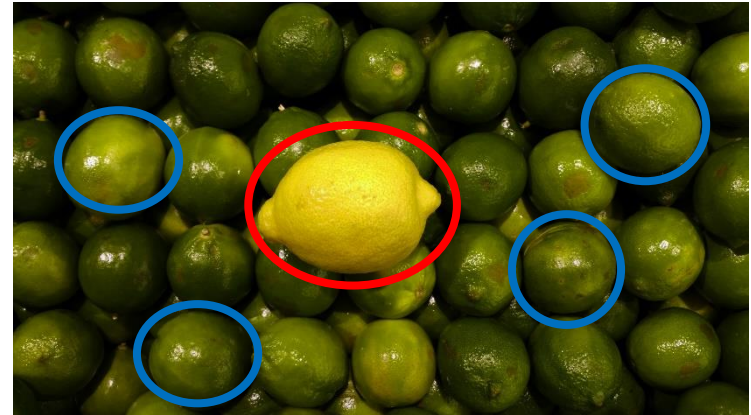
  ✓ Learning Minimum Volume Sets
    (Scott & Nowak 2006)



- Find a region in the input space that capture the normal points
- NOT competitive with pattern based approaches (Emmott et al. 2013)

Oregon State OSU

# Rare Pattern Anomaly Detection (RPAD)

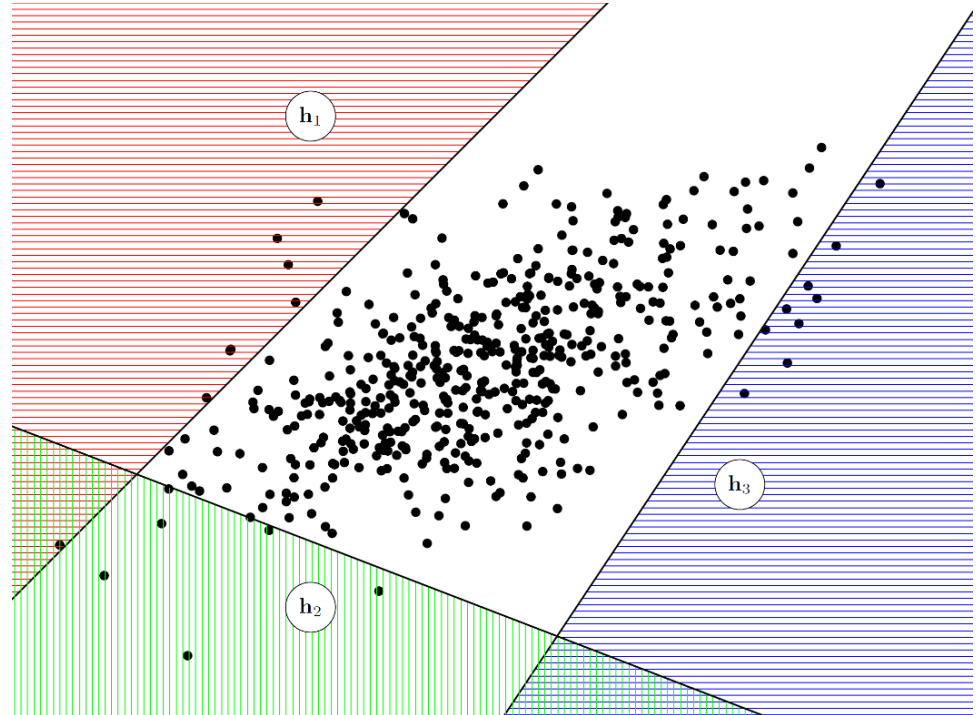- A pattern simply can be a specific color or size



Color



Size

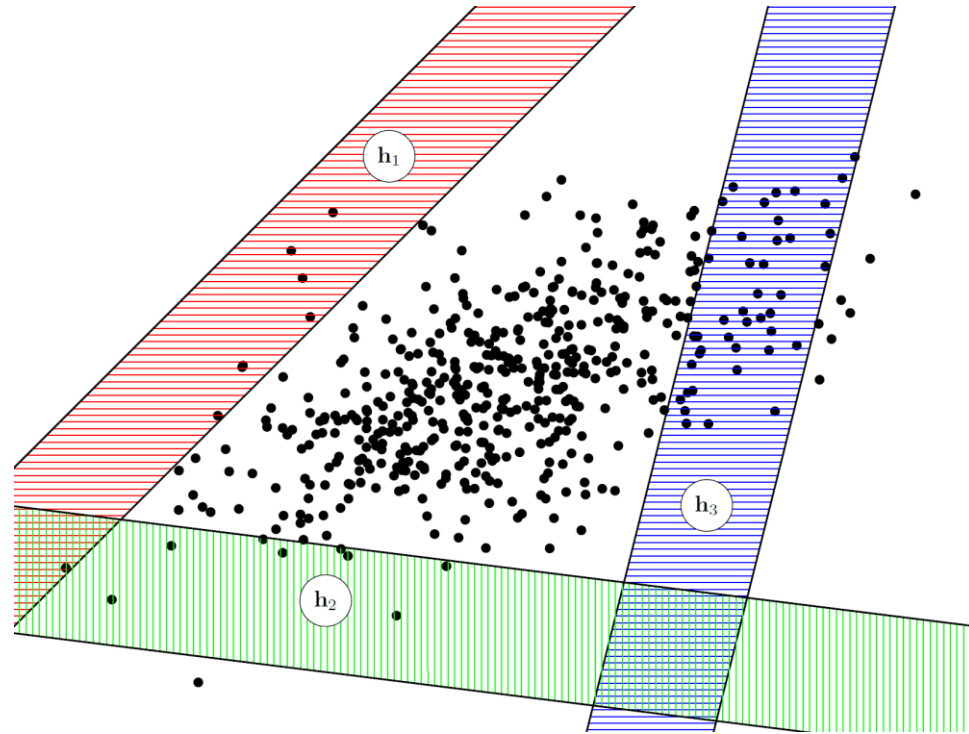- Identifies anomaly based on the characteristics of rare patterns

# Half-spaces

- A half-space pattern is an oriented $d$-dimensional hyperplane

- The half-space mass algorithm [Chen et al. 2015] operates in this pattern space



- **Anomaly score** : Mean frequency estimates of random half-spaces containing the query point $x$

# LODA

- Construct $T$ sparse random projections in of $\mathcal{R}^d$

- Each time, Estimate 1D histogram density from projected input data

- **Anomaly score:** geometric average of the $T$ densities corresponding a query point



- Each bin of the histograms corresponds to a stripe in $\mathcal{R}^d$

- The perpendicular direction of the projection defines the orientation of the stripe

- Bin width corresponds to the width of the stripe