

# ARUNDO

## **Applied Machine Learning for Anomaly Detection on Equipment**

Kristin Hornæs  
Lukasz Mentel  
Trung Doan



Kristin



Lukasz



Trung

09:00 - 09:30	Welcome & Introduction to Anomaly Detection
09:30 - 10:00	Set-up environment
10:00 - 10:30	Walk through examples of AD methods in Jupyter notebooks
10:30 - 10:45	Coffee break
10:45 - 11:45	Improve models
11:45 - 12:00	Presentation of selected solutions

----- LUNCH BREAK -----

# BUZZWORD BINGO

Digital  
Transformation

Operational  
Intelligence

Streaming  
Analytics

IT Operations  
Analytics

Industry 4.0

Blended  
Analytics

# Increase efficiency and productivity



Data is the  
jetfuel



### **Decrease downtime**

Unexpected downtime on a single asset can cost upwards of a million dollars per day

## **INTERCONNECTED AMBITIONS**



### **Increase efficiency**

Increase profits despite a decreasing price per barrel



### **Scalable, actionable insight**

Make models which can be easily applied to the company's entire portfolio

# ARUNDO

provides software products to **enable**  
enterprise-scale machine learning and advanced  
analytics applications for **industrial companies**

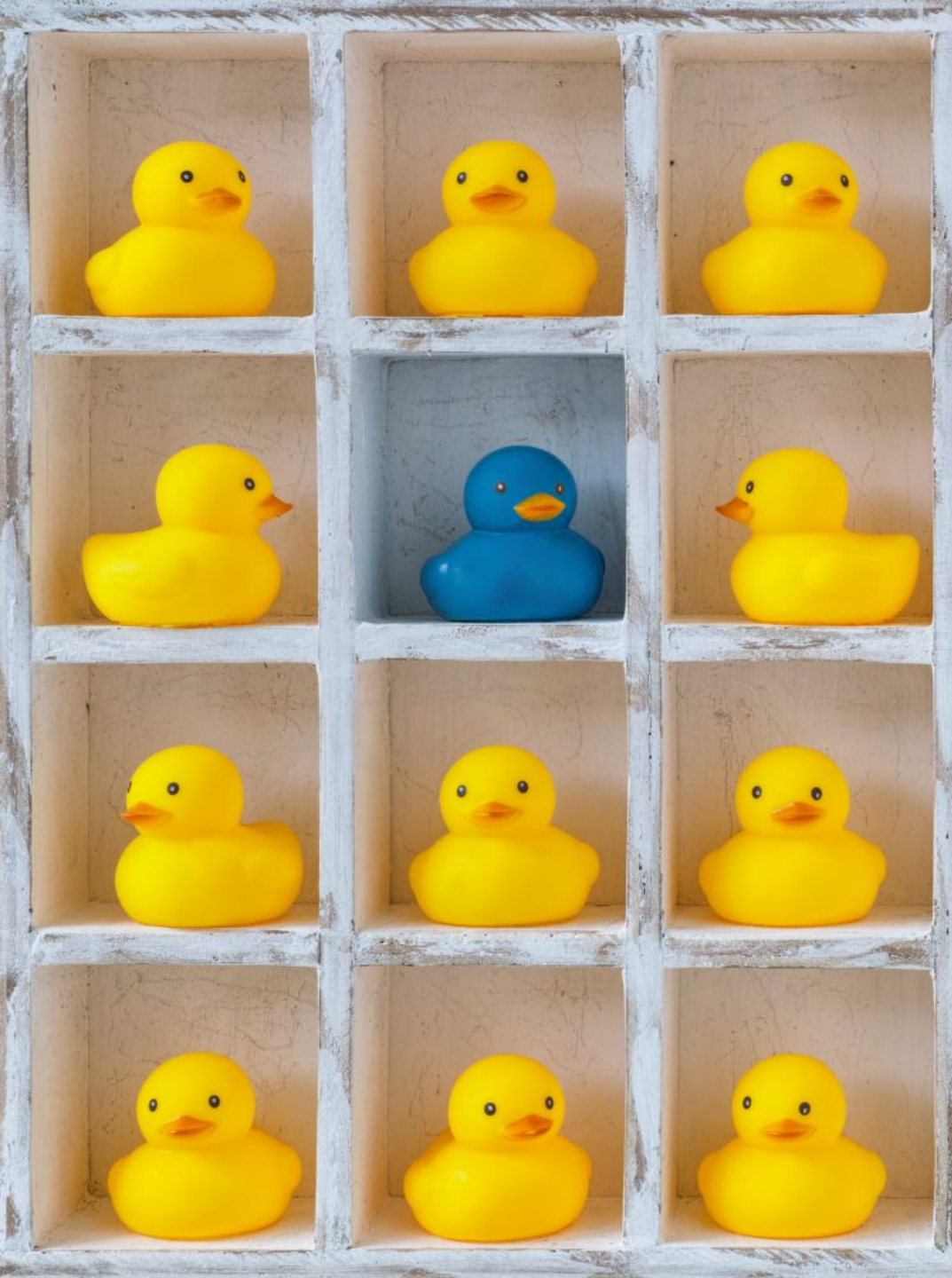




A large colony of brown fluffy penguin chicks is shown. In the center-right, an adult King penguin stands out, facing away from the camera and slightly to the right. It has a dark blue-grey body and a bright yellow patch on its neck. The text "What is an anomaly?" is overlaid in the center of the image.

What is an anomaly?





**“DATAPOINTS, ITEMS,  
OBSERVATIONS OR EVENTS  
THAT DO NOT CONFORM TO  
THE EXPECTED PATTERN”**

# Examples of anomaly detection



Health  
monitoring



Video  
surveillance



Equipment  
monitoring



Fraud  
detection

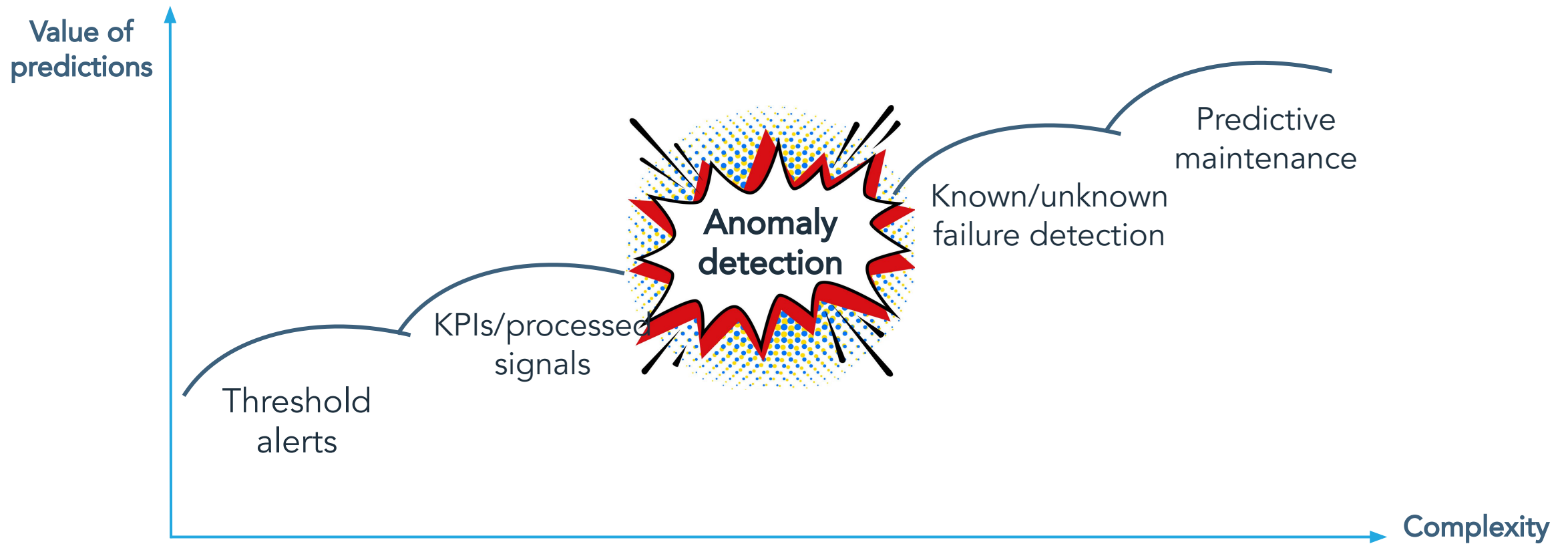


Intrusion  
detection



Spam  
filtering

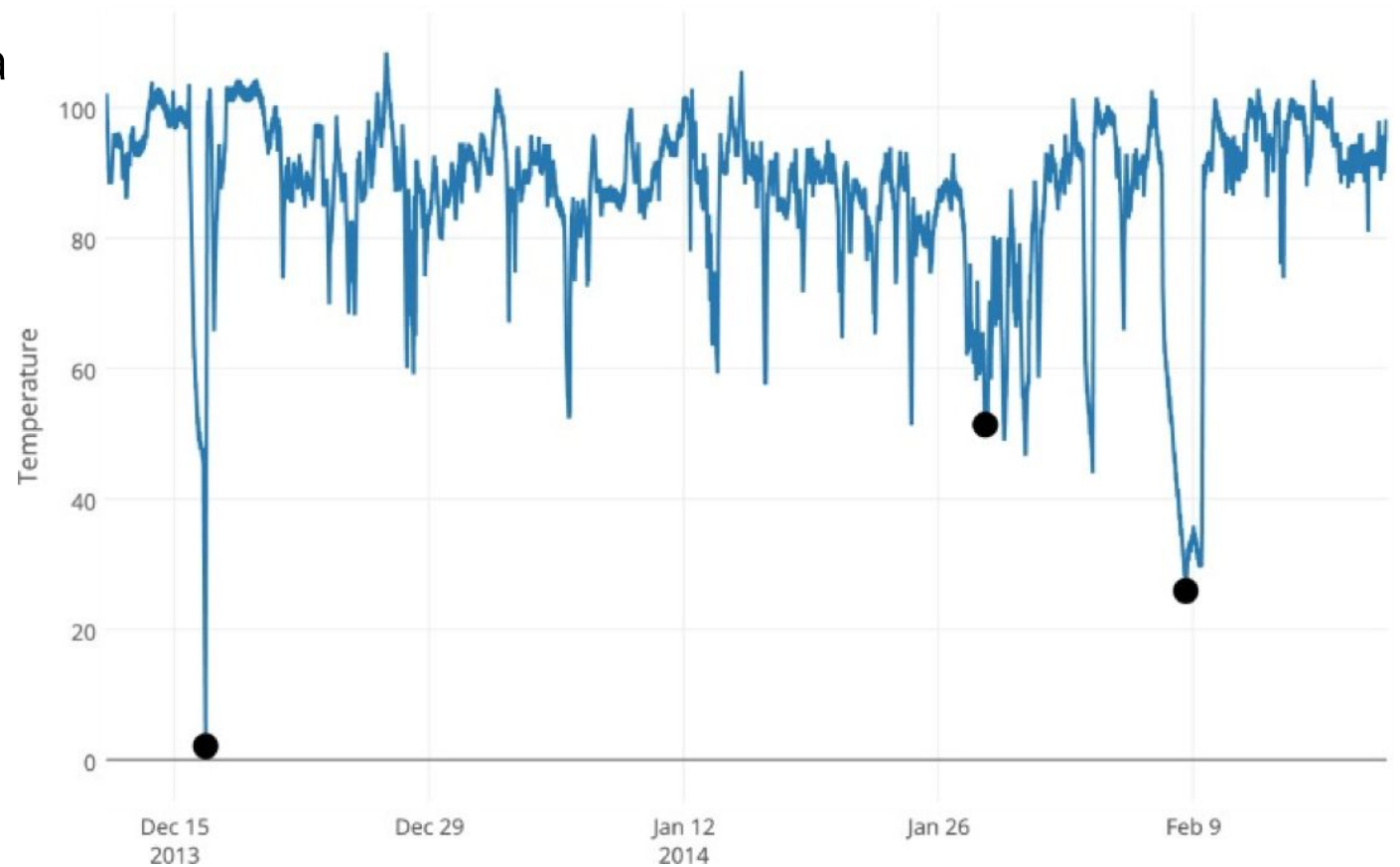
# The stages of data-driven equipment monitoring



# Anomaly detection in **equipment monitoring**

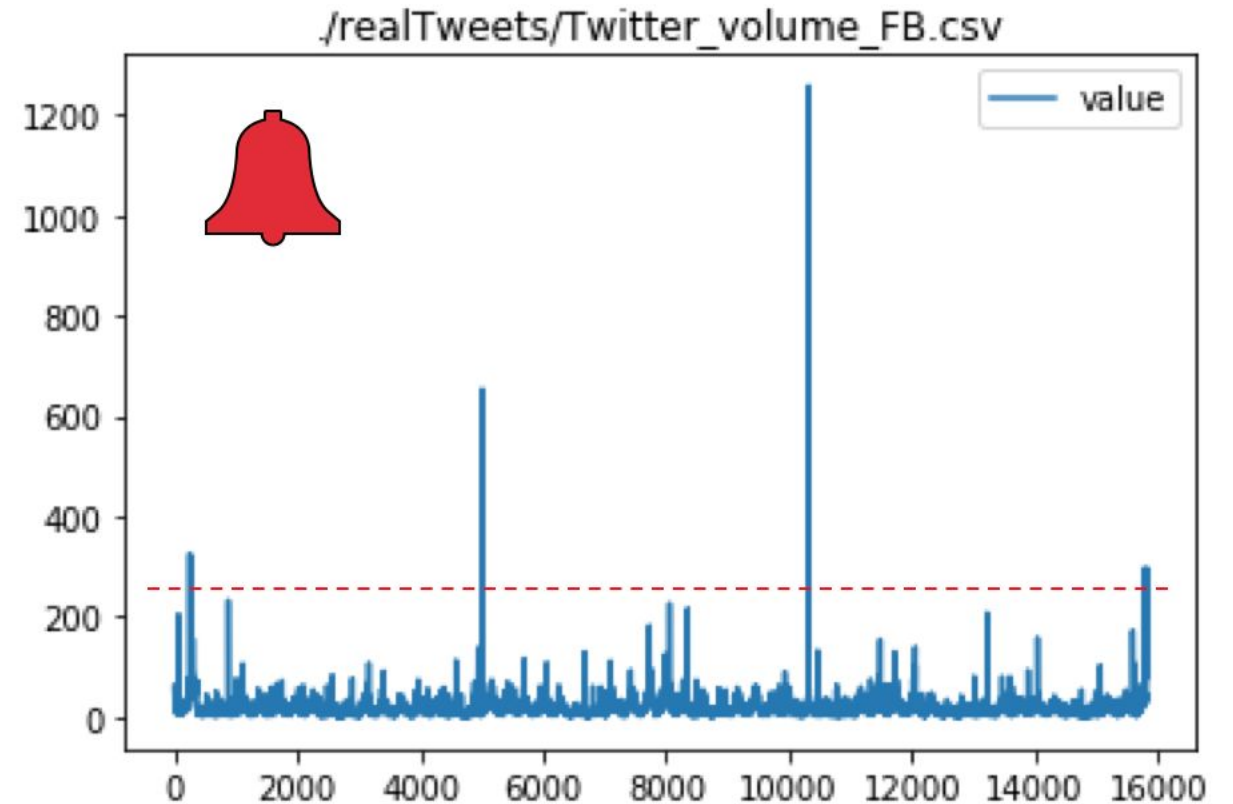
Previously unseen patterns can be a sign of:

- **misconfiguration**
- increasing mechanical **wear-out**
- **unforeseen** situations



## How can you detect anomalies?

- Define a threshold for each sensor channel
- Raise a notification once a specified threshold is violated





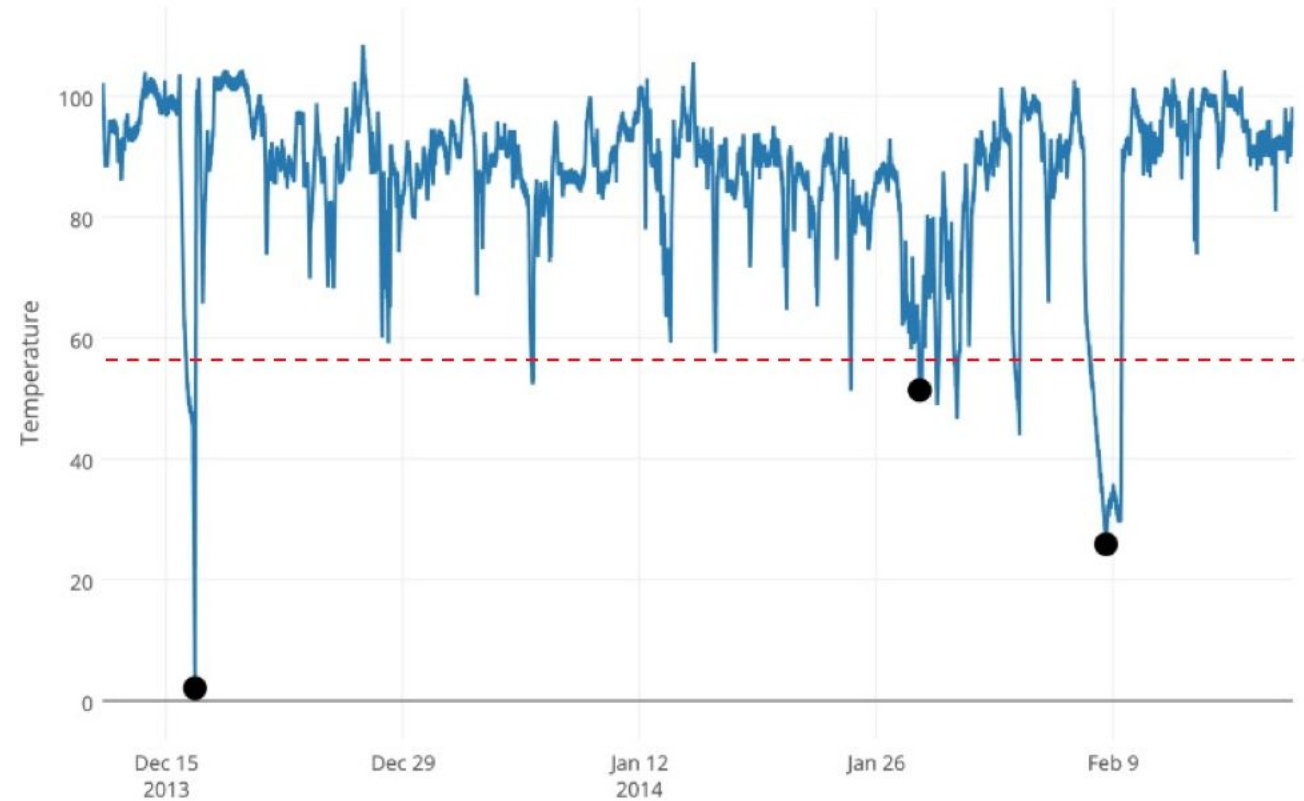


An oil rig can have upwards of 15.000 sensors  
(the newest have more than 50.000 sensors!!!)



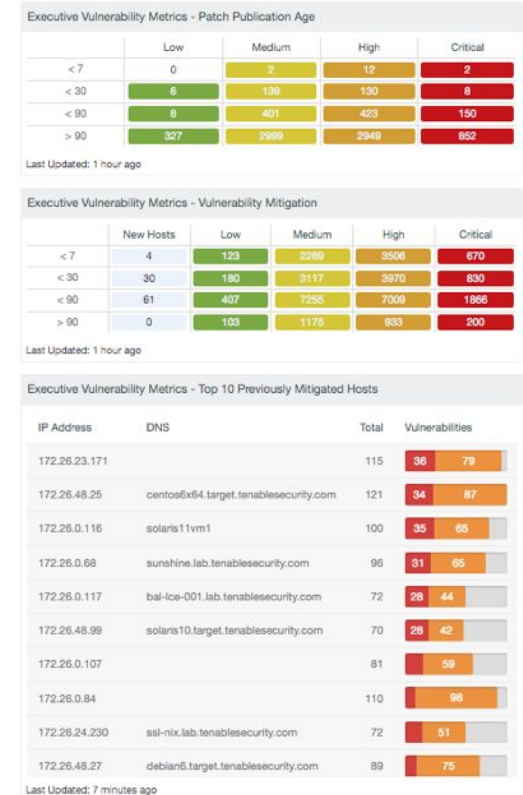
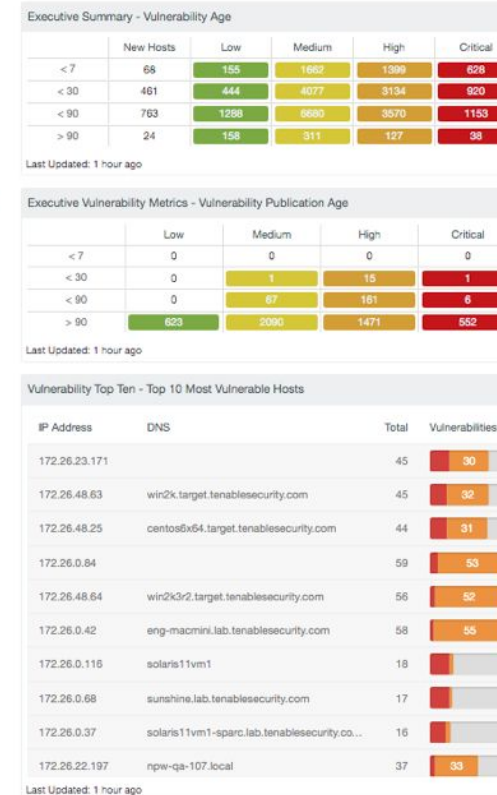
## How can you detect anomalies?

- Causes many false alerts
- Does not take into account the joint characteristics of multiple channels



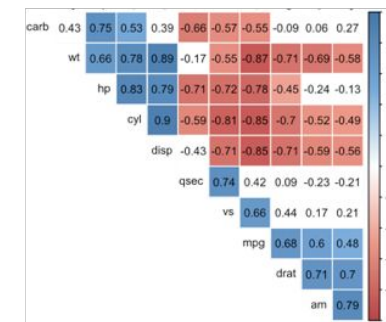
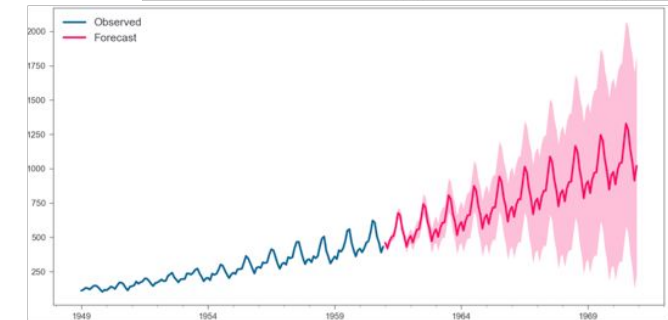
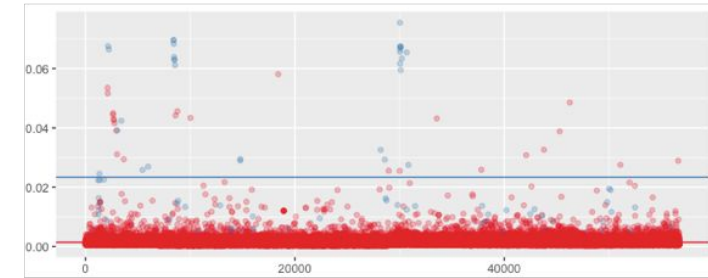
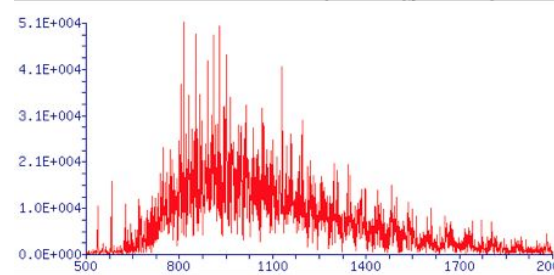
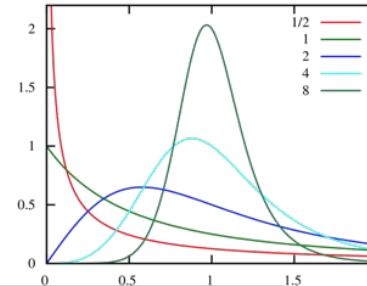
# Manual analysis is immensely time-consuming and unreliable

- **Massive** amount of multi-sensor data
- **Complex** systems
- **Rare** faults



# Multivariate anomaly detection

- No prior knowledge about anomalies
- No precise boundary
- Data often contain noise
- Normal behaviour keeps evolving
- Temporal dependencies
- Highly unbalanced classes
- High dimensionality and multimodal dependencies

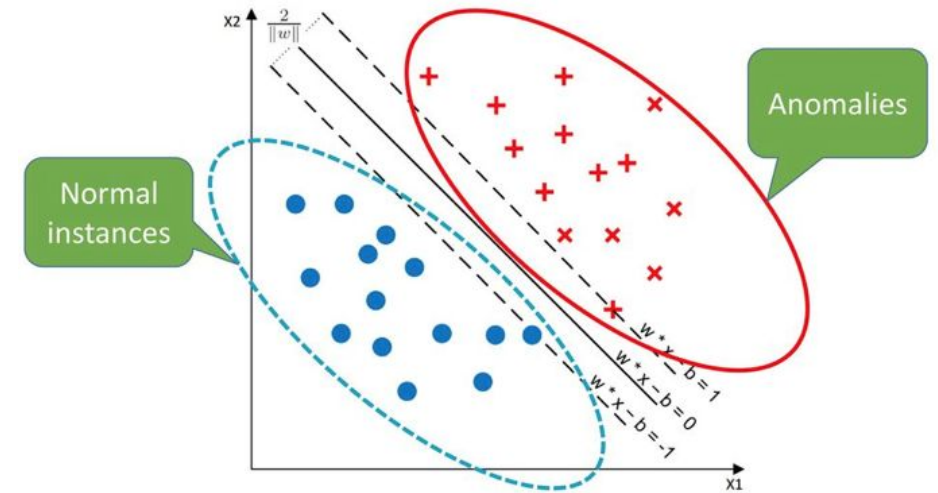


# Approaches

Data

Labeled data?  
Predict Y from X?

# Approaches



Data

Labeled data?  
Predict Y from X?



Supervised  
Learning

Develop predictive model  
based on both input and  
output data

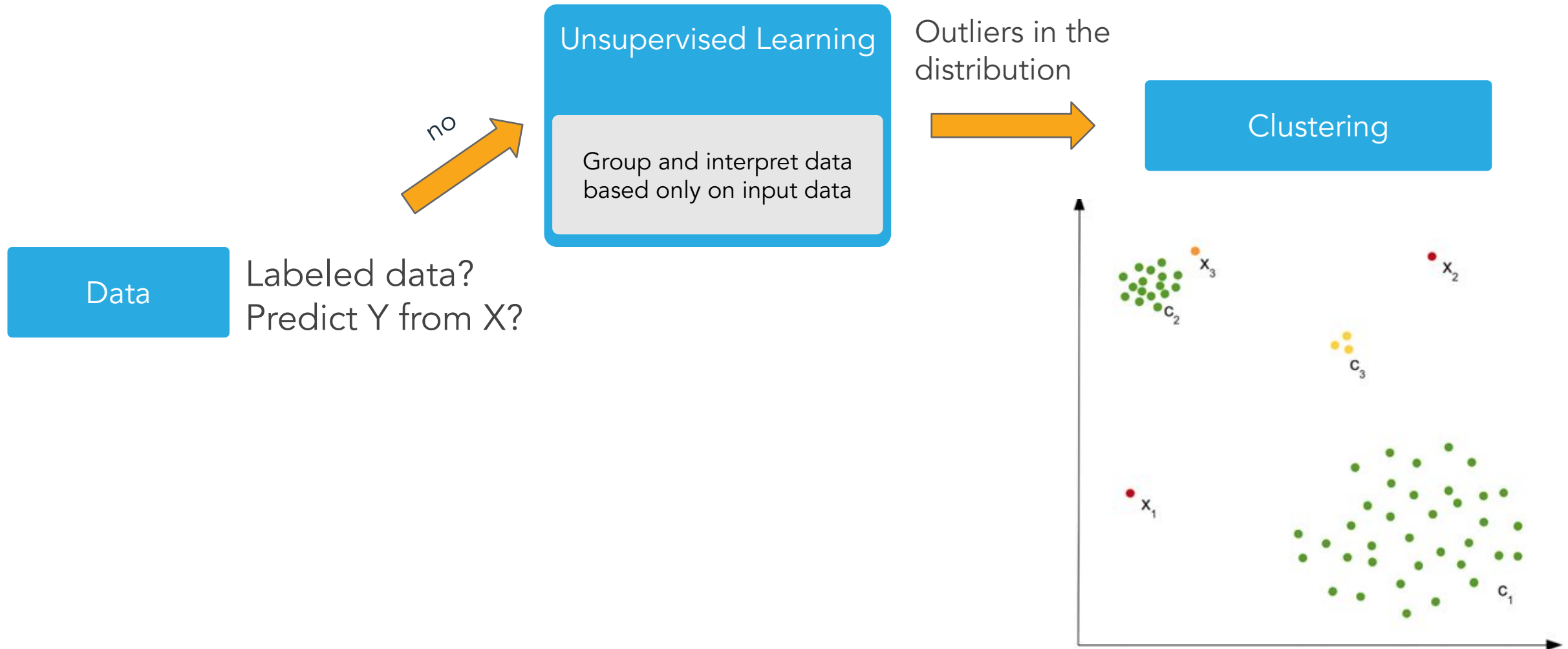
Partition data  
based on labels



Classification

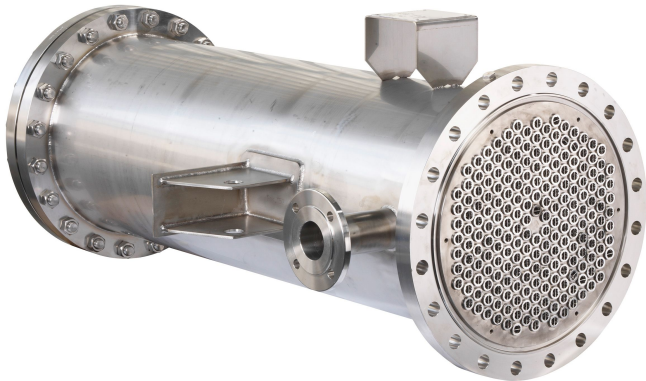
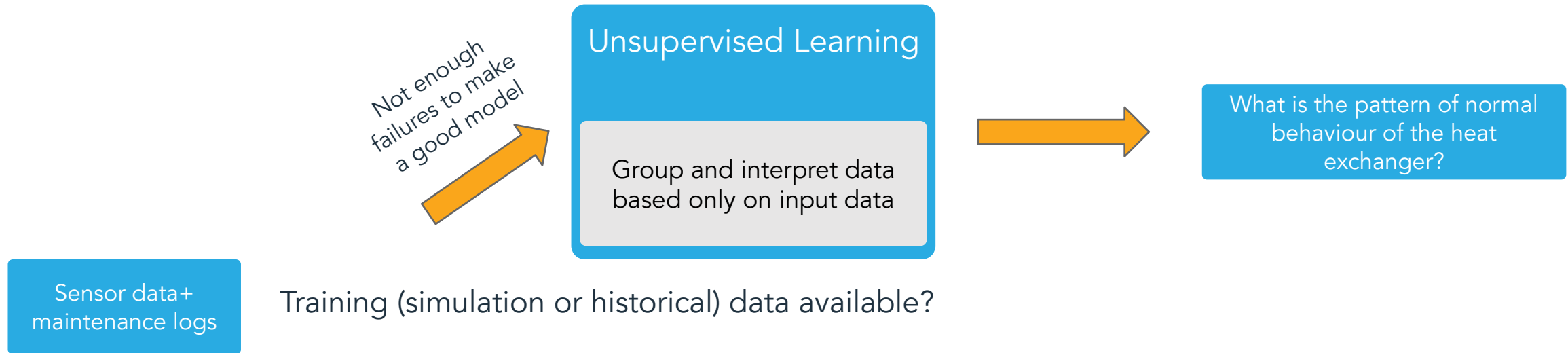
Regression

# Approaches





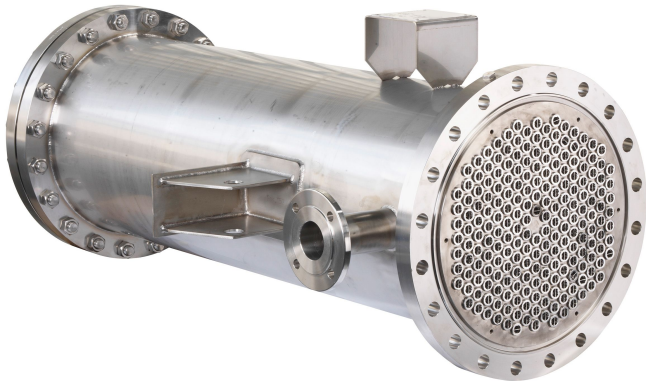
# Real life example - Leakage on heat exchangers



# Real life example - Leakage on heat exchangers

Sensor data+  
maintenance  
logs

Training (simulation or historical) data available?



Yes but not a  
lot of failures

Supervised  
Learning

Develop predictive model  
based on both input and  
output data

discrete

Has my heat exchanger sprung  
a leak?

continuous

What is the predicted  
performance of my heat  
exchanger?

# Approaches

## Supervised Methods

Random Forest

Support Vector Classification

KNN Classifier

Logistic Regression

## Unsupervised Methods

One Class SVM

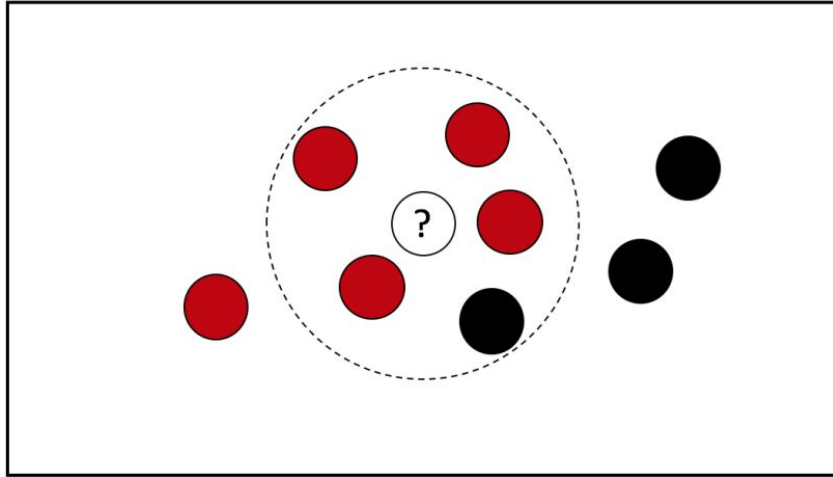
Isolation Forest

Elliptic Envelope

Local Outlier Factor

Autoencoder

# $k$ -NN Classifier



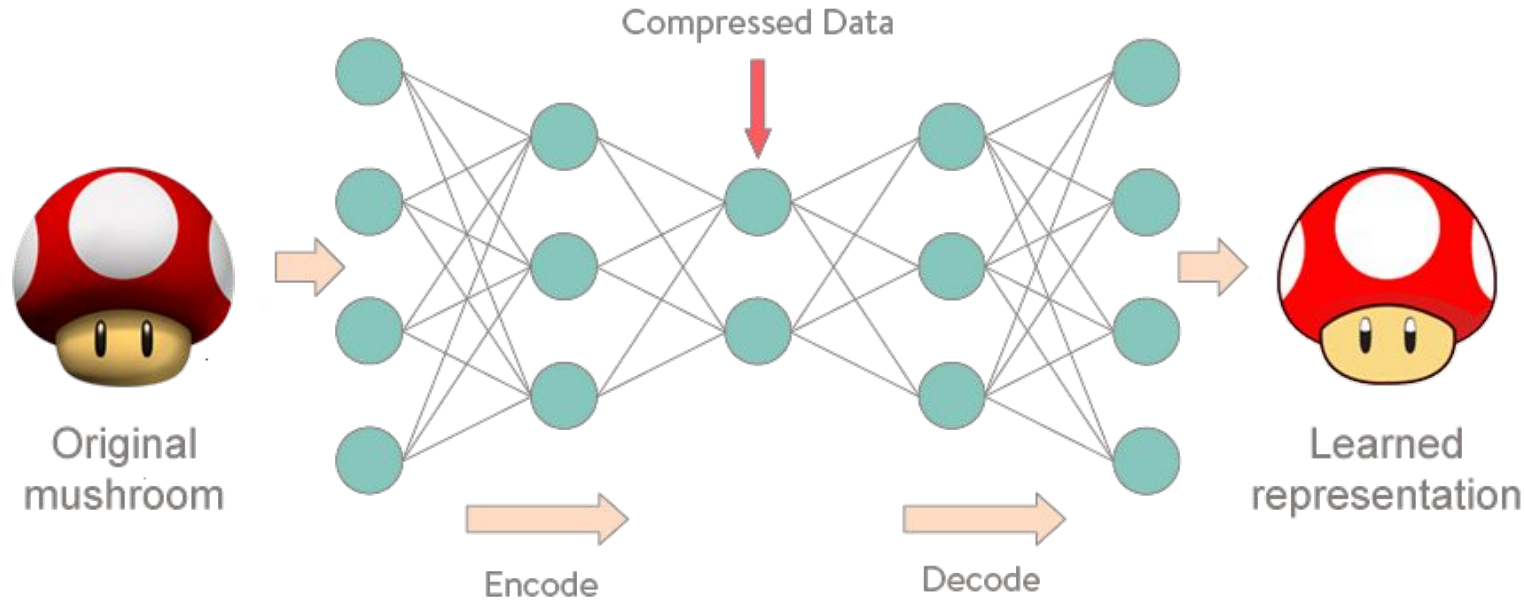
**$k$ -NN classifies a data point based on how its neighbors are classified.**

If a data point is surrounded by 4 red points and 1 black point, that data point is likely a red point by majority vote.

Tune for  $k$  - the number of nearest neighbors to include in the majority voting process



# Autoencoders



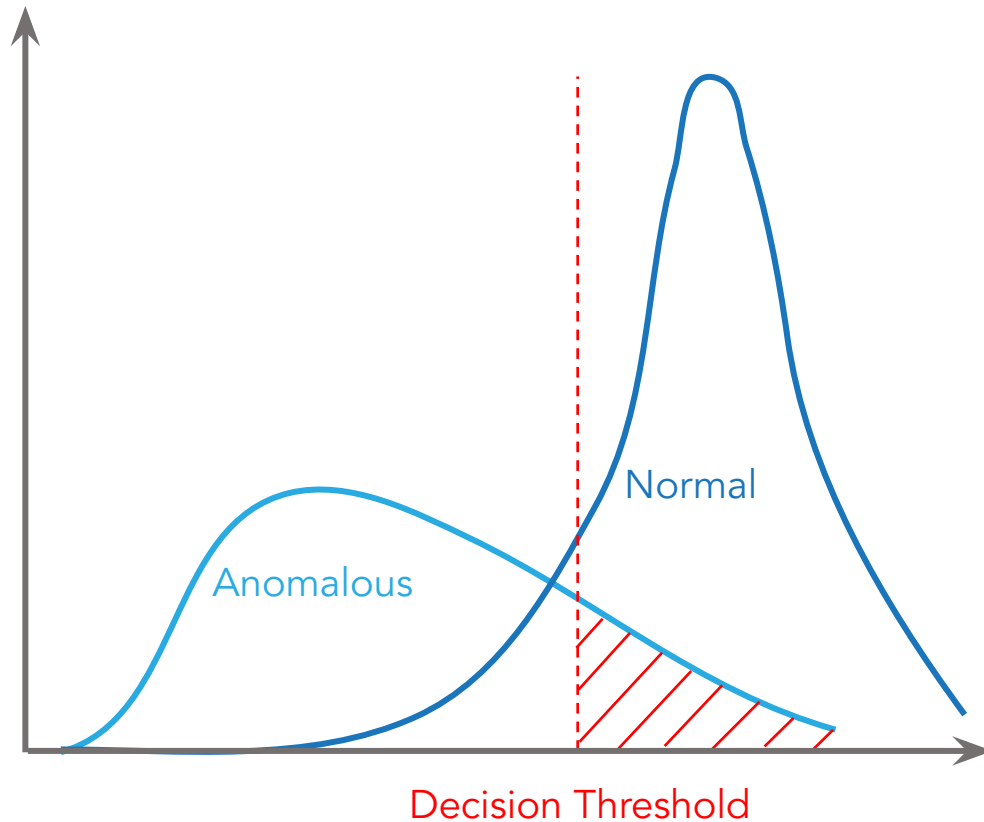
The job of those models is to predict the input, given that same input.

Learns a representation of the training data and recreates the input at the output layer.

Used for data compression and learning generative models from data.

We optimize the parameters of our Autoencoder model in such way that the reconstruction error is minimized.

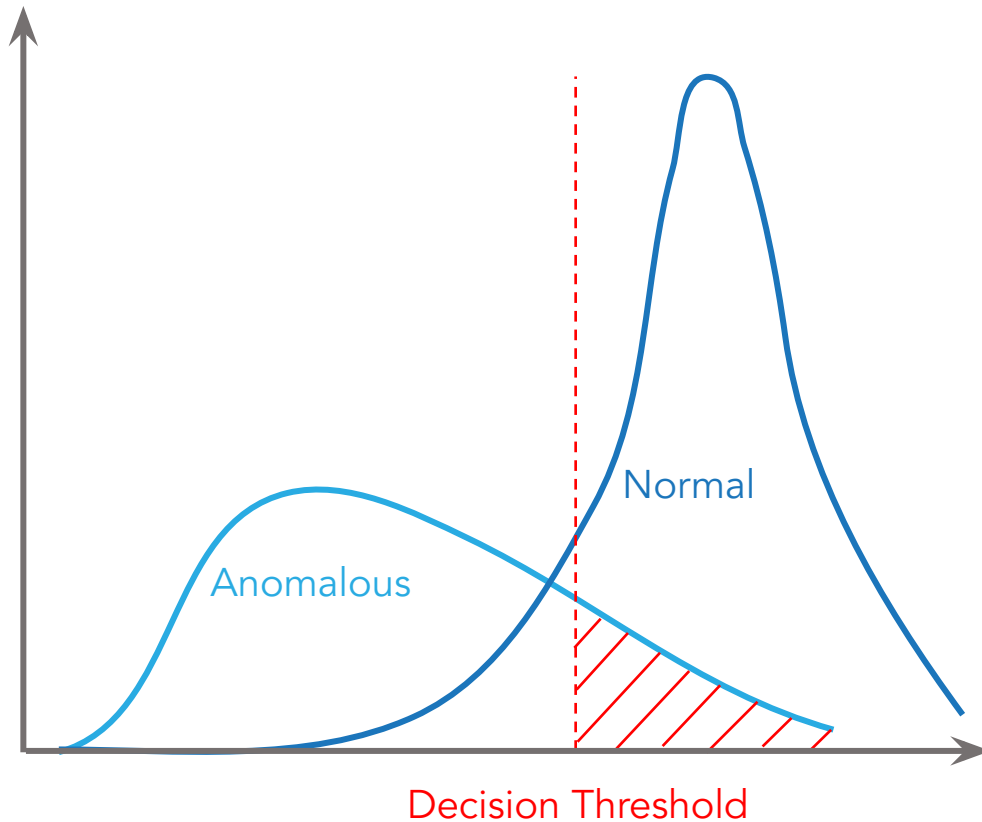
# Machine Learning Predictions are not 100% Accurate



- False positives (ML predicts anomalous but system is normal) can and will occur
- False negatives (ML predicts normal but system is anomalous) can and will occur



# Machine Learning Predictions are not 100% Accurate



- Choose trained algorithms which minimize these effects as much as possible
- Identify any additional sources of data which may further help minimize these effects
- Evaluate the expected probabilities of each scenario using independent historical data
- Retrain the system if there is evidence that probabilities change significantly with time

# Workshop details

## Task

- Make a model to identify anomalies based on the dataset provided
- The model will be tested on a separate test dataset
- The best model will be chosen based on f1 score

## Practicalities

- Log on to jupyter hub using your github account
  - <https://wids.arundo.com>

## Notebooks

- Opening and looking at the data
  - 01-Model Development
- Improve baseline model results
  - Please add comments to explain your code where possible