

Aprendizagem Aplicada à Segurança

Mário Antunes

September 22, 2023

University of Aveiro

Table of Contents

Class Introduction

Grading

Class Schedule

Environment

Bibliography

Name: Mário Antunes

E-Mail: *mario.antunes@ua.pt*

Office: 19.2.15 (IT1)

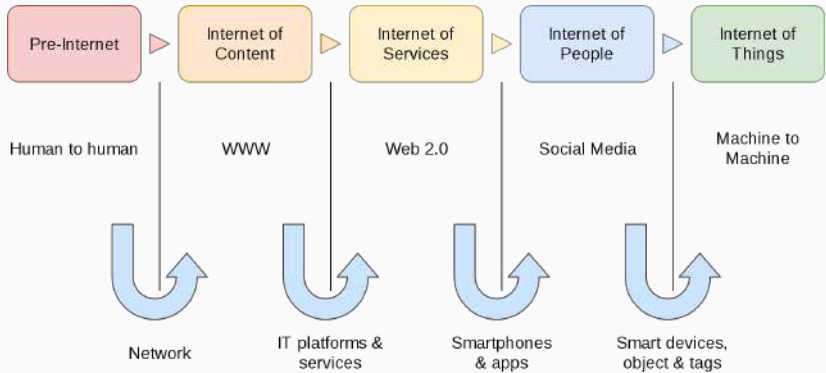


Class Introduction

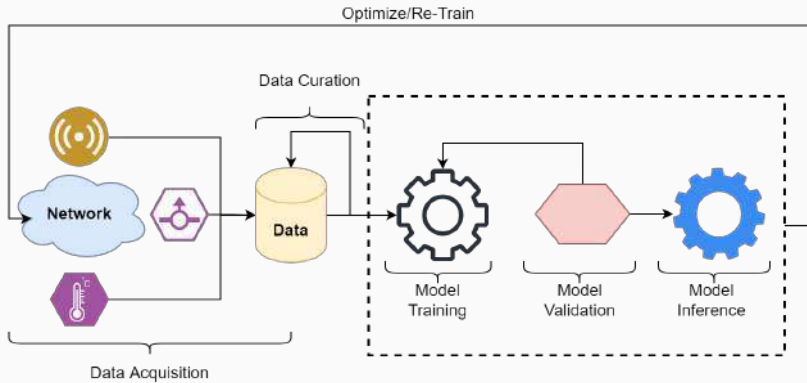
- Given the evolution of the threats
- And the complexity of the systems
- AI/ML are gaining traction as a useful tool



Class Introduction



Class Introduction



- 50% Theory + 50% Practice
- Discrete: 25% Mid-term Exam + 25% Final Exam + 20% Project Idea + 30% Project
- Final: 50% Final Exame + 50% Project

Class Schedule i

Date	Class	Topic
15/09/2023	1	Introduction
22/09/2023	2	
29/09/2023	3	SPAM Detector
06/10/2023	4	
13/10/2023	5	
20/10/2023	6	Anomaly Detection
27/10/2023	7	
03/11/2023	8	Mid-term Exam
10/11/2023	9	
17/11/2023	10	Malware Analysis
24/11/2023	11	
01/12/2023	12	
08/12/2023	13	
15/12/2023	14	Project
22/12/2023	15	



- All of the books are available here:
[*https://learning.oreilly.com/*](https://learning.oreilly.com/)

- [1] S. Halder and S. Ozdemir, *Hands-On Machine Learning for Cybersecurity: Safeguard your system by making your machines intelligent using the Python ecosystem*. Packt Publishing Ltd, 2018.
- [2] C. Chio and D. Freeman, *Machine Learning and Security*. O'Reilly, 2018.

- [3] A. Parisi, *Hands-On Artificial Intelligence for Cybersecurity: Implement smart AI systems for preventing cyber attacks and detecting threats and network anomalies*. Packt Publishing Ltd, 2019.
- [4] E. Tsukerman, *Machine Learning for Cybersecurity Cookbook*. Packt Publishing Ltd, 2019.
- [5] J. P. Mueller and R. Stephens, *Machine Learning Security Principles*. Packt Publishing Ltd, 2019.

Aprendizagem Aplicada à Segurança

Mário Antunes

September 22, 2023

University of Aveiro

Table of Contents

SPAM

SPAM Detection

Binary Classification

Text Mining

Natural Language Processing (NLP)

Classification Model

Model Evaluation

- The term “spam” is internet slang that refers to unsolicited commercial email (UCE).
- The first reported case of spam occurred in 1898, when the New York Times reported unsolicited messages circulating in association with an old swindle.
- The term “spam” was coined in 1994, based on a now-legendary Monty Python’s Flying Circus sketch, where a crowd of Vikings sings progressively louder choruses of “SPAM! SPAM! SPAM!”

SPAM



Dear Sir,

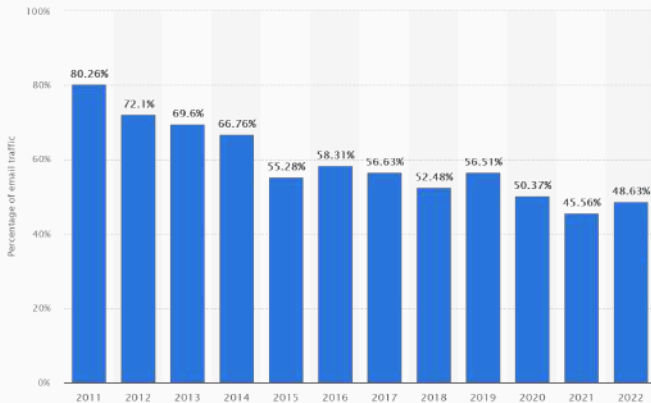
I am prince [REDACTED] from Nigeria. Your help would be very appreciated. I want to transfer all of my fortune outside if Nigeria due to a frozen account, If you could be so kind and transfer small sum of 3 500 USD to my account, I would be able to unfreeze my account and transfer my money outside of Nigeria. To repay your kindness, I will send 1 000 000 USD to your account.

Please contact me to proceed

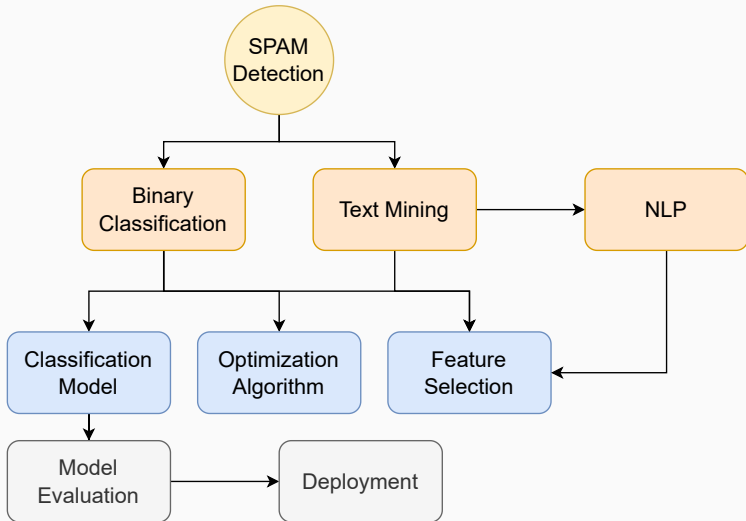
Prince [REDACTED]

- Huge list of *https://en.wikipedia.org/wiki/Anti-spam_techniques*
- From common sense to *Bayesian spam filtering*
- Unfortunately it is a costly battle

Fight against SPAM

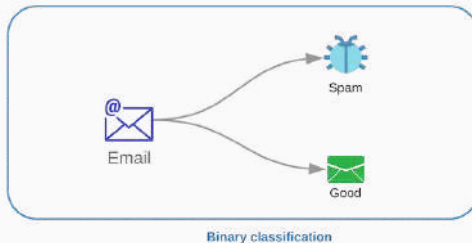


SPAM Detection



Binary Classification

- Binary classification is the task of classifying the elements of a set into two groups (each called class) on the basis of a classification rule.
- For this application one message can either be spam or ham.



- Text mining is the process of deriving high-quality information from text.
- Combines concepts from Machine Learning, Linguistic and statistical analysis.
- In this area we will explore the methods used to rank words/tokens and the BoW model.

Bag of Words (Bow) model

	the	red	dog	cat	eats	food
1. the red dog →	1	1	1	0	0	0
2. cat eats dog →	0	0	1	1	1	0
3. dog eats food →	0	0	1	0	1	1
4. red cat eats →	0	1	0	1	1	0

Natural Language Processing (NLP)

- NLP gives the computers the ability to understand text.
- Combines *Syntax* and *Semantic* into the analysis.
- One famous examples are the Large Language Models (LLMs) that power OpenAI Chat GPT.

Classification Model

- SPAM detection is “considered” a toy example.
- As such, we will explore two of the simplest learning models: Naive Bayes and Logistic Regression.

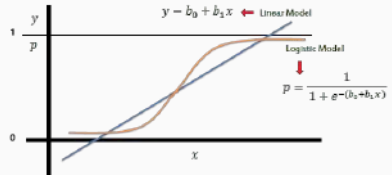
Likelihood of the Evidence given that the Hypothesis is True

Prior Probability of the Hypothesis

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

Posterior Probability of the Hypothesis given that the Evidence is True

Prior Probability that the evidence is True

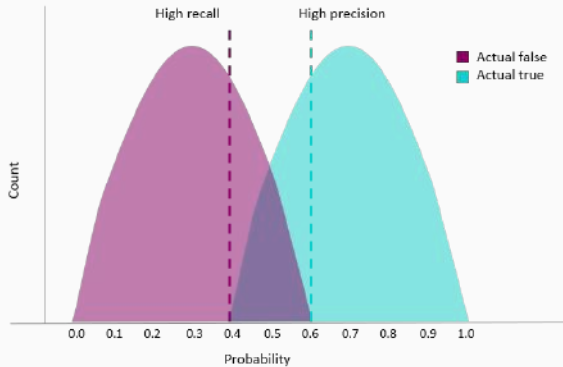


Model Evaluation

- Classification model can be evaluated using a confusing matrix
- The simplest methods to evaluate a model is through accuracy: $acc = \frac{TP+TN}{TP+TN+FP+FN}$

	Predicted Positive	Predicted Negative	
Actual Positive	TP <i>True Positive</i>	FN <i>False Negative</i>	Sensitivity $\frac{TP}{(TP + FN)}$
Actual Negative	FP <i>False Positive</i>	TN <i>True Negative</i>	Specificity $\frac{TN}{(TN + FP)}$
	Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$

Model Evaluation



Aprendizagem Aplicada à Segurança

Mário Antunes

October 14, 2023

Universidade de Aveiro

Context

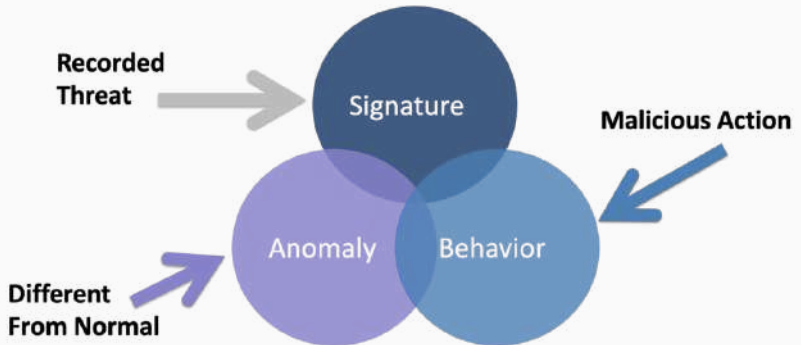
It is becoming difficult to identify Cybersecurity attacks. These attacks can originate internally due to malicious intent or negligent actions or externally by malware, target attacks, and APT (Advanced Persistent Threats).

But insider threats are more challenging and can cause more damage than external threats because they have already entered the network.

These activities present unknown threats and can steal, destroy or alter the assets.

Earlier firewalls, web gateways, and some other intrusion prevention tools are enough to be secure, but now hackers and cyber attackers can bypass approximately all these defense systems.

Therefore with making these prevention systems strong, it is also equally essential to use detection. So that if hackers get into the network, the system should be able to detect their presence.



Signature detection requires knowing what to look for and comparing hashes or other strings to identify a match. Signature detection is a common feature found within antivirus and IPS/IDS products.

Behavior detection looks for malicious or other known behavior characteristics and alarms the SOC when a match is made. An example is identifying port scanning or a file attempting to encrypt your hard drive, which is an indication of ransomware behavior. Antimalware and sandboxes are examples of tools that heavily leverage behavior detection capabilities.

Anomaly detection it takes into consideration hot topics including big data, threat intelligence, and “zero-day” detection.

Anomalies

Anomaly detection, also called outlier detection, is the identification of unexpected events, observations, or items that differ significantly from the norm:

- Anomalies in data occur only very rarely
- The features of data anomalies are significantly different from those of normal instances

What is an anomaly?

Generally speaking, an **anomaly** is something that differs from a norm: a deviation, an exception. In software engineering, by anomaly we understand a rare occurrence or event that doesn't fit into the pattern, and, therefore, seems suspicious. Some examples are:

- sudden burst or decrease in activity;
- error in the text logs;
- sudden rapid drop or increase in temperature.

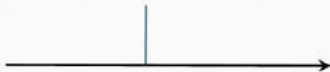
Common reasons for outliers are:

- data preprocessing errors;
- noise;
- fraud;
- attacks.

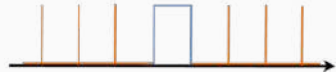
Anomalies can be broadly categorized as:

- Point anomalies: A single instance of data is anomalous if it's too far off from the rest.
- Contextual anomalies: The abnormality is context specific. This type of anomaly is common in time-series data.
- Collective anomalies: A set of data instances collectively helps in detecting anomalies.

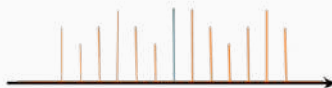
Types of Anomalies #2



(a) Point Anomaly



(b) Collective Anomaly



(c) Contextual Anomaly

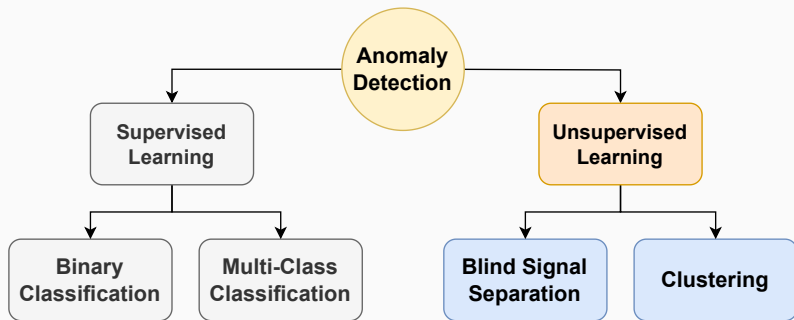
Examples

Network anomalies: Anomalies in network behavior deviate from what is normal, standard, or expected. To detect network anomalies, network owners must have a concept of expected or normal behavior. Detection of anomalies in network behavior demands the continuous monitoring of a network for unexpected trends or events.

Application performance anomalies: These are simply anomalies detected by end-to-end application performance monitoring. These systems observe application function, collecting data on all problems, including supporting infrastructure and app dependencies. When anomalies are detected, rate limiting is triggered and admins are notified about the source of the issue with the problematic data.

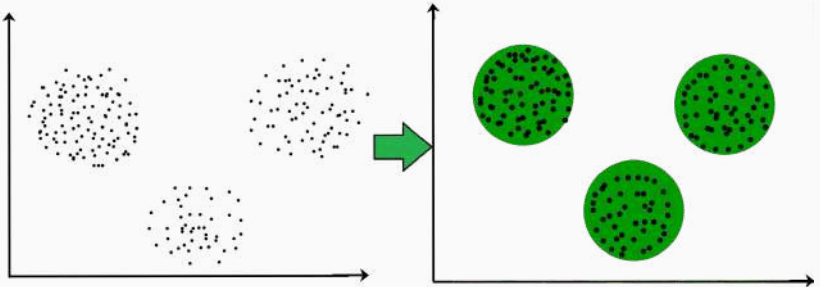
Web application security anomalies: These include any other anomalous or suspicious web application behavior that might impact security such as XSS attacks or DDOS attacks.

Anomaly Detection



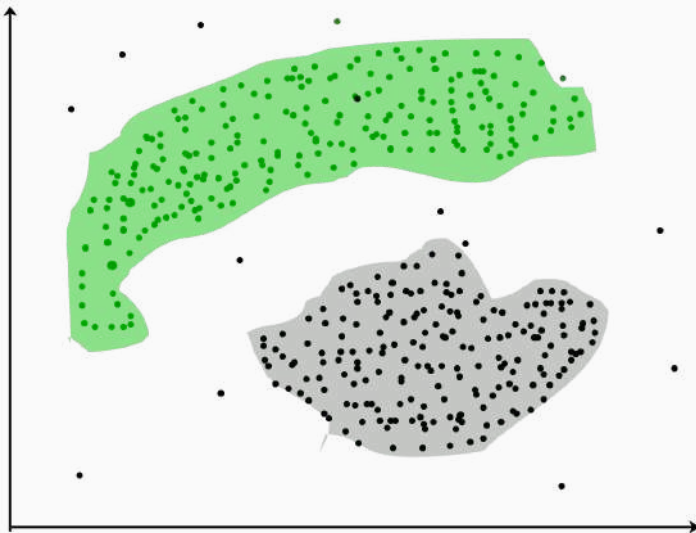
Clustering

Type of **unsupervised learning method**. Generally, it is used as a process to find meaningful structure, explanatory underlying processes, generative features, and groupings inherent in a set of examples.



- **Density-Based Methods:** These methods consider the clusters as the dense region having some similarities and differences from the lower dense region of the space. These methods have good accuracy and the ability to merge two clusters.
- **Hierarchical Based Methods:** The clusters formed in this method form a tree-type structure based on the hierarchy. New clusters are formed using the previously formed one.
- **Partitioning Methods:** These methods partition the objects into k clusters and each partition forms one cluster. This method is used to optimize an objective criterion similarity function such as when the distance is a major parameter.

Clustering: Anomaly Detection

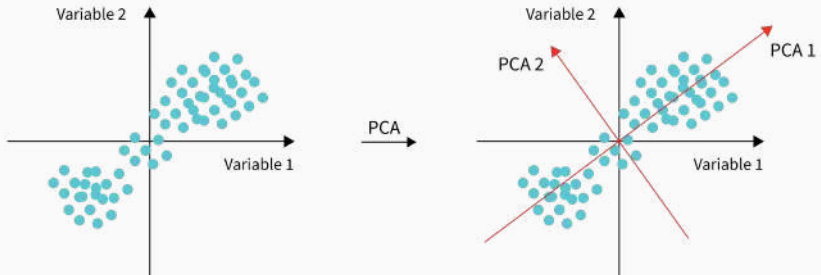


Blind Source Separation (BSS) refers to a problem where both the sources and the mixing methodology are unknown, only mixture signals are available for further separation process.

In several situations it is desirable to recover all individual sources from the mixed signal, or at least to segregate a particular source.

Blind Source Separation: PCA

**** Principal component analysis****, or PCA, is a statistical procedure that allows you to summarize the information content in large data tables by means of a smaller set of “summary indices” that can be more easily visualized and analyzed.



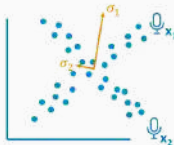
Independent Component Analysis (ICA) is a powerful technique in the field of data analysis that allows you to separate and identify the underlying independent sources in a multivariate data set.



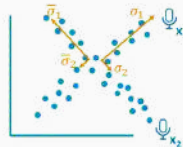
PCA finds main directions in data: the *principal components*.



PCA fails for data sets where we have more than one principal direction



ICA solves this problem for us by focusing on independent components rather than principal components



Blind Source Separation: NNMF

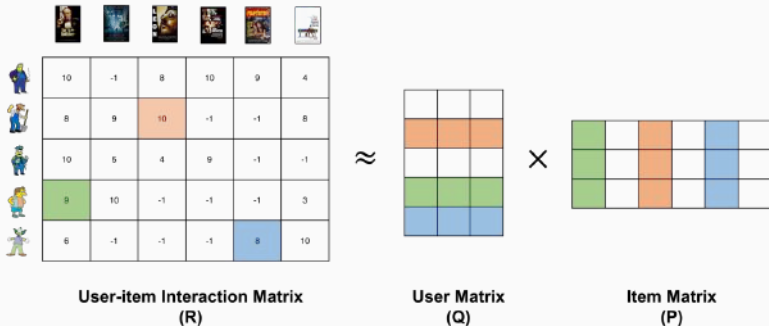
- **Non-negative matrix factorization (NNMF)** is a group of algorithms in multivariate analysis and linear algebra where a matrix V is factorized into two matrices W and H , with the property that all three matrices have no negative elements.
- This non-negativity makes the resulting matrices easier to inspect. Also, in applications such as processing of audio spectrograms or muscular activity, non-negativity is inherent to the data being considered.

The diagram illustrates the Non-negative Matrix Factorization (NNMF) process. It shows three matrices represented by grids of cells:

- Matrix W:** A 4x2 grid (4 rows, 2 columns).
- Matrix H:** A 2x6 grid (2 rows, 6 columns).
- Matrix V:** A 4x6 grid (4 rows, 6 columns).

The matrices are arranged in the equation: $W \times H \approx V$. The multiplication symbol \times is placed between W and H, and the approximation symbol \approx is placed between H and V.









Blind Source Separation: Anomaly Detection



Auto Encoders

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

-  Input Cell
-  Backfed Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Capsule Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Gated Memory Cell
-  Kernel
-  Convolution or Pool

Perceptron (P)



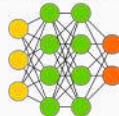
Feed Forward (FF)



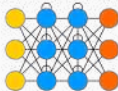
Radial Basis Network (RBF)



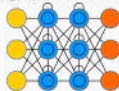
Deep Feed Forward (DFF)



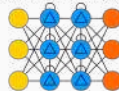
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



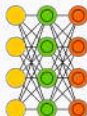
Gated Recurrent Unit (GRU)



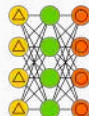
Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)



Hopfield Network (HN)



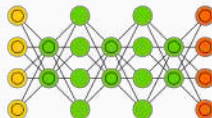
Boltzmann Machine (BM)



Restricted BM (RBM)

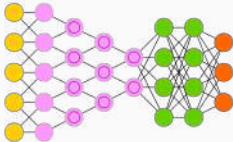


Deep Belief Network (DBN)

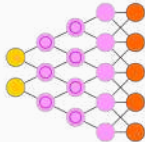


Neural Networks #2

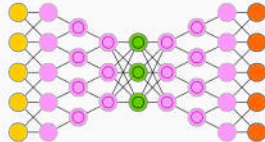
Deep Convolutional Network (DCN)



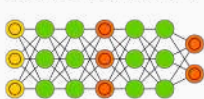
Deconvolutional Network (DN)



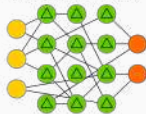
Deep Convolutional Inverse Graphics Network (DCIGN)



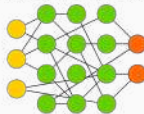
Generative Adversarial Network (GAN)



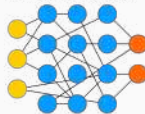
Liquid State Machine (LSM)



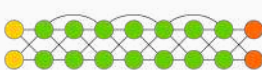
Extreme Learning Machine (ELM)



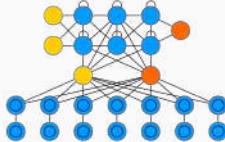
Echo State Network (ESN)



Deep Residual Network (DRN)



Differentiable Neural Computer (DNC)



Neural Turing Machine (NTM)



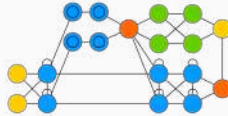
Capsule Network (CN)



Kohonen Network (KN)



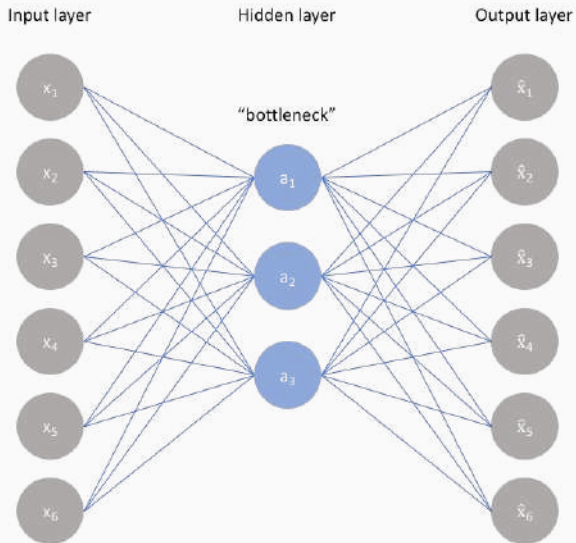
Attention Network (AN)



- Autoencoders are an unsupervised learning technique in which we leverage neural networks for the task of representation learning.
- Specifically, we'll design a neural network architecture such that we impose a bottleneck in the network which forces a compressed knowledge representation of the original input.

If the input features were each **independent** of one another, this compression and subsequent reconstruction would be a very **difficult task**. However, if some sort of structure exists in the data (ie. correlations between input features), this structure can be learned and consequently leveraged when forcing the input through the network's bottleneck.

Auto Encoders #2



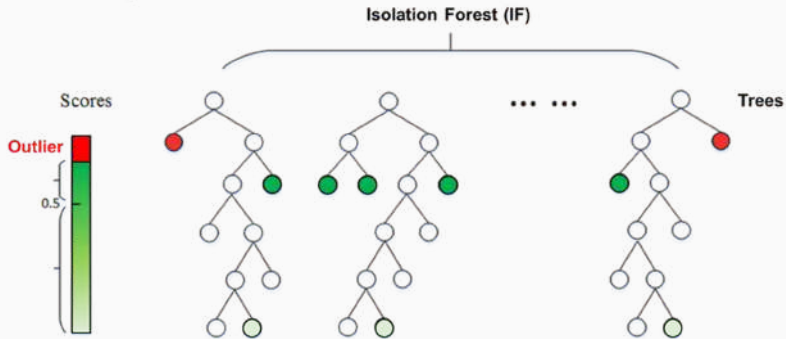
- As visualized, we can take an unlabeled dataset and frame it as a supervised learning;
- This network can be trained by minimizing the reconstruction error;
- The bottleneck is a key attribute of our network design; without the presence of an information bottleneck, our network could easily learn to simply memorize the input values by passing these values along through the network.

Other Methods

- IsolationForest **isolates** observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.
- Since recursive partitioning can be represented by a tree structure, the number of splittings required to isolate a sample is equivalent to the path length from the root node to the terminating node.
- This path length, averaged over a forest of such random trees, is a measure of normality and our decision function.

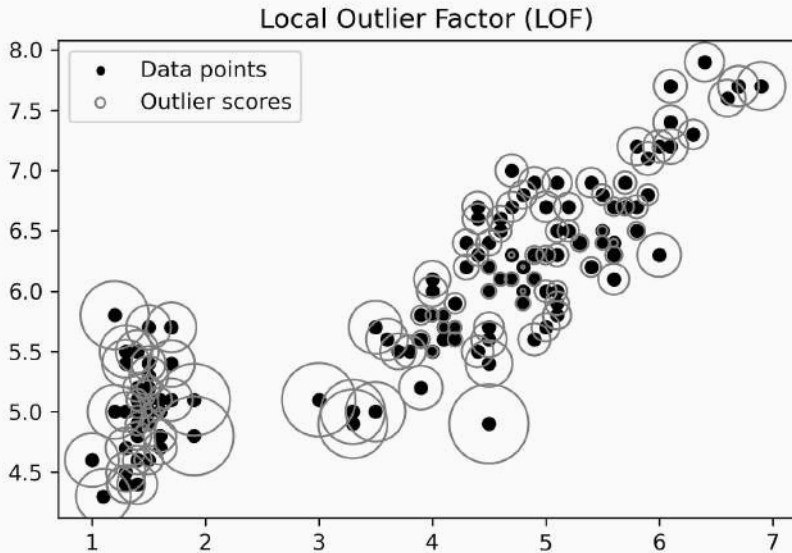
Random partitioning produces noticeably shorter paths for anomalies. Hence, when a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be **anomalies**.

Anomaly Detection - IsolationForest #2



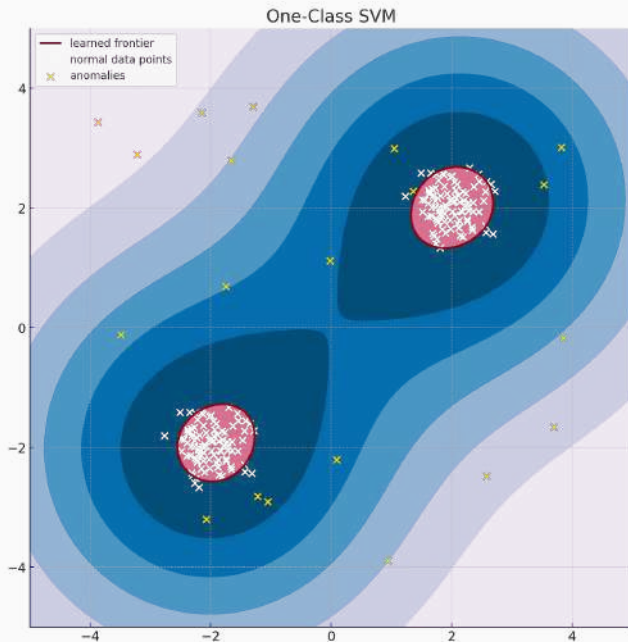
- Local Outlier Factor (LOF) measures the local deviation of the density of a given sample with respect to its neighbors.
- It is local in that the anomaly score depends on how isolated the object is with respect to the surrounding neighborhood.
- More precisely, locality is given by k -nearest neighbors, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors. These are considered outliers.

Anomaly Detection - Local Outlier Factor #2



- Many approaches are based on the estimation of the density of probability for the normal data. Anomalies corresponds to those samples where the density of probability is “very low”.
- Now, SVMs are max-margin methods, i.e. they do not model a probability distribution. Here the idea is to find a function that is positive for regions with high density of points, and negative for small densities.
- One-Class SVM is similar, but instead of using a hyperplane to separate two classes of instances, it uses a hypersphere to encompass all of the instances. Now think of the “margin” as referring to the outside of the hypersphere – so by “the largest possible margin”, we mean “the smallest possible hypersphere”.

Anomaly Detection - OneClassSVM #2



Tips

- Feature scaling is the process of normalizing the range of features in a dataset.
- Real-world datasets often contain features that are varying in degrees of magnitude, range, and units.
- Therefore, in order for machine learning models to interpret these features on the same scale, we need to perform feature scaling.

Resources

Neural Networks Zoo

AutoEncoders

Principal components analysis