# DCML-CPS - Module 5

# Data Analysis

Tommaso Zoppi

University of Trento – Povo (IT)

tommaso.zoppi@unifi.it

tommaso.zoppi@unitn.it

# Course Map

**Monitoring**

1. Basics and Metrology

2. Monitoring

**Testing**

3. Fault Injection

4. Robustness Testing

**Anomaly Detection**

5. Data Analysis

6. Supervised ML

7. Unsupervised ML

8. Meta-Learning

9. Error/Intrusion Detection

**Tools & Libs**

**Deep Learning**

RCL RESILIENT COMPUTING LAB

# Data Analysis

RCL
RESILIENT COMPUTING LAB
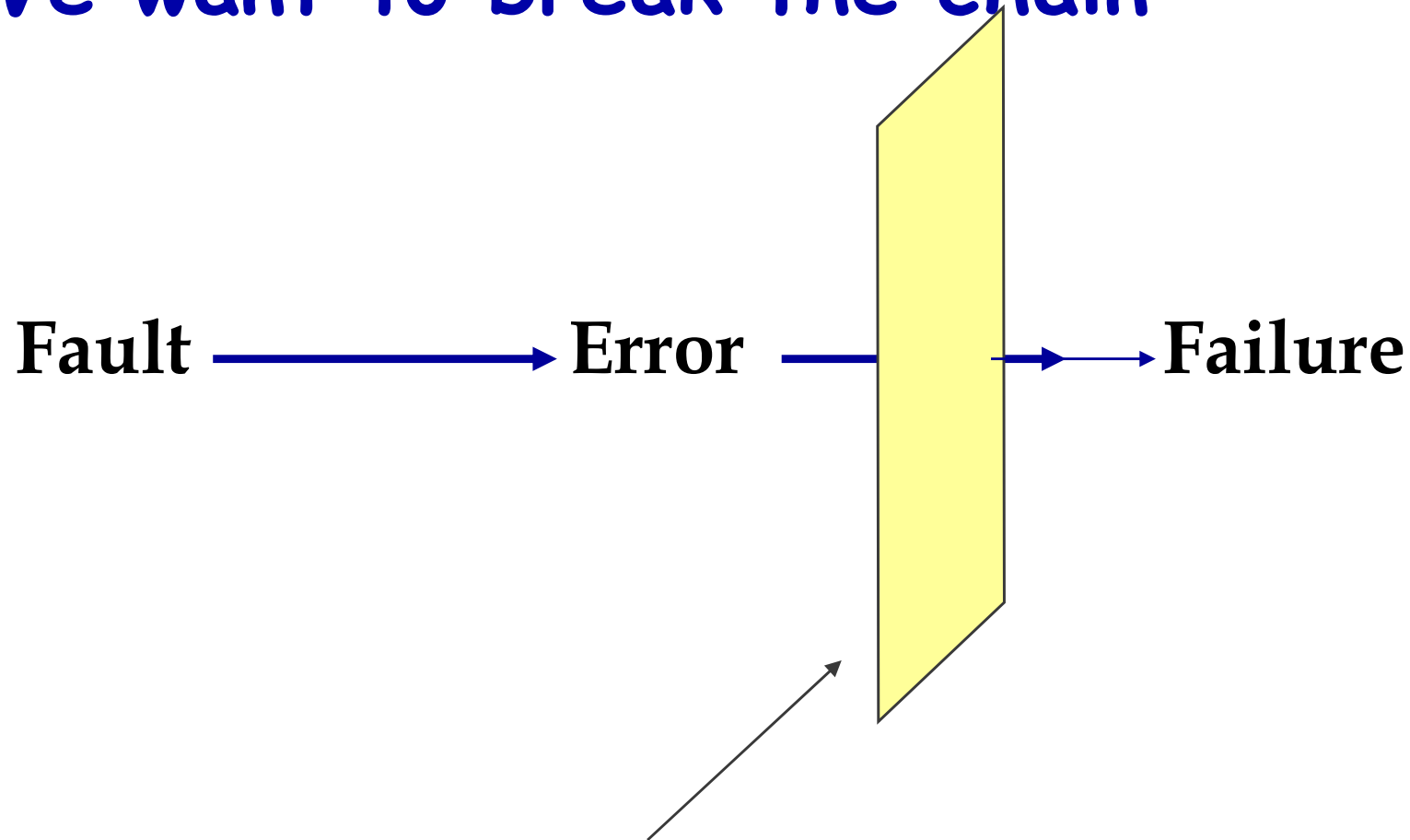
► We understood how to gather a nice amount of data

- As a result of monitoring activities

- application logs

- Also performing fault injection or robustness testing

► But, then... What are we going to do with these?

► **We want to break the chain**

Fault ⟶ Error ⟶ Failure

How can I detect the activation of a fault to activate fault tolerance mechanisms?

## ► Video from

– Machine Learning Basics | What Is Machine Learning? | Introduction To Machine Learning | Simplilearn

- https://www.youtube.com/watch?v=ukzFI9rgwfU

Machine Learning Basics _ What Is Machine Learning_ _ Introduction To Machine Learning _ Simplilearn.mp4

► Facebook recognizes your friend in a picture from an album of tagged photographs

- **Supervised** learning. Here Facebook is using tagged photos to recognize the person. Therefore, the tagged photos become the labels of the pictures and we know that when the machine is learning from labeled data, it is supervised learning.

RCL
RESILIENT COMPUTING LAB

► **Recommending new songs based on someone's past music choices**

- **Supervised** learning. The model is training a classifier on pre-existing labels (genres of songs). This is what Netflix, Pandora, and Spotify do all the time, they collect the songs/movies that you like already, evaluate the features based on your likes/dislikes and then recommend new movies/songs based on similar features.

► Scenario 3: Analyze bank data for suspicious-looking transactions and flag the fraud transactions

- **Unsupervised** learning. In this case, the suspicious transactions are not defined, hence there are no labels of "fraud" and "not fraud". The model tries to identify outliers by looking at anomalous transactions and flags them as 'fraud'.

RCL
RESILIENT COMPUTING LAB

# General Structure of a Dataset



**Feature (F)**

**Feature Set (FS)**

**Feature Value (FV)**

**Dataset (D)**

**Data Point (DP)**

► Overall, steps of a data analysis process are

– Generating/Getting the dataset

– Exploratory study:

- which are the features, label …

– Feature Selection

- Understanding which features are more informative
- Often embedded in the algorithm: representation learning

– Algorithm(s) Selection

- Based on experimental analyses, availability of implementations, computational complexity, …

– Choice of the Metric + Gather and Discuss Results

RCL
RESILIENT COMPUTING LAB

► **Dataset** (D): data structured as a set of **data points** (DP), described through a set of **Features** (F)

- Each data point is composed by features values (FV)
- A group of features defines a Feature Set (FS).
- Algorithms can be defined for the data analysis process

► They are being trained by feeding them with a Training Set T

- Each algorithm extracts the feature values from T

► The result of this process is an Algorithm **Model**

► **The model is used for testing.**

– Data points in the **validation set** V are given as input to the model separately

– The model outputs a **numeric score** AS that allows to decide on the «class» of the data point

– AS can then be converted into a boolean score (for binary classification), or into categories (for multi-class) through a **Decision Function** DF

► If **Ground Thruth** (label) is available, it is possible to calculate Metric Scores

RCL
RESILIENT COMPUTING LAB
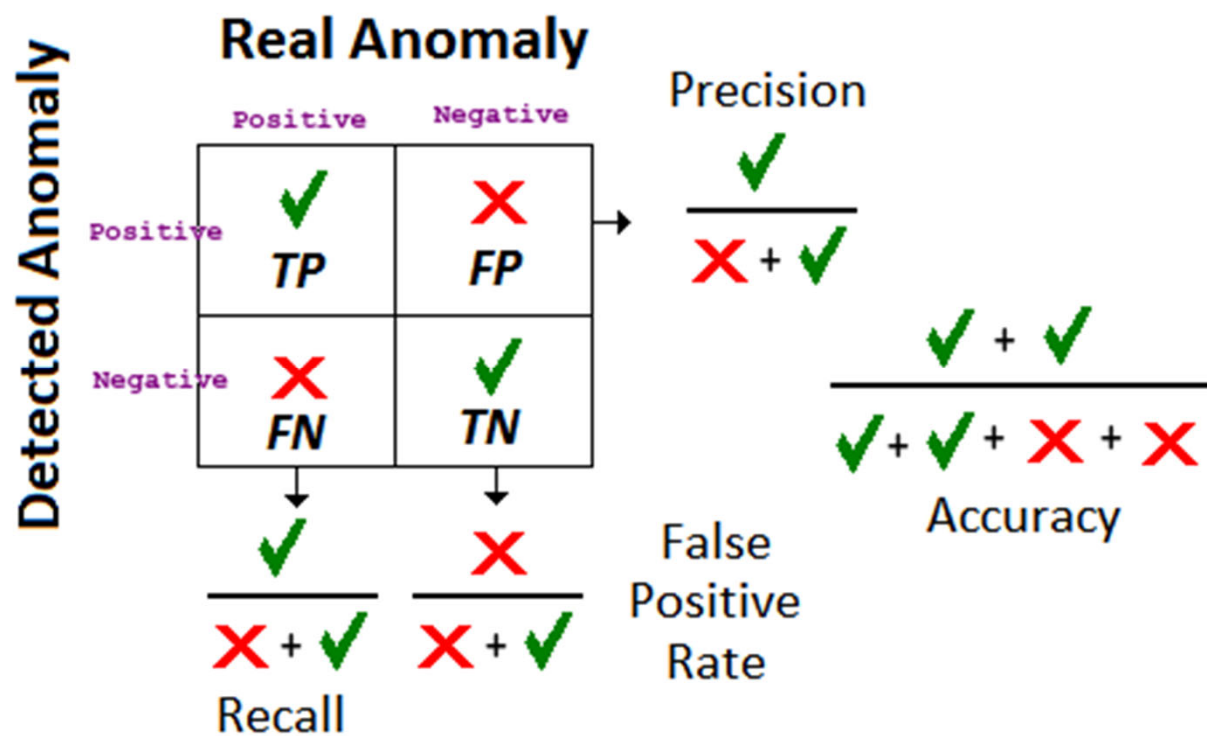
# Classification Metrics

► **The suitability and the effectiveness of such techniques are usually evaluated and compared depending on specific metrics.**

- True Positives (TP)
- True Negatives (TN)
- False Positives (FP)
- False Negatives (FN)

► Single metrics might be aggregated in order to build complex ones

► F-Score($\beta$) is obtained by combining Precision and Recall,

– weighting them by using a parameter $\beta$.

► F-Score(1) is often referred as F-Measure

– $\beta > 1$ favours recall (low FNs)
– $\beta < 1$ favours Precision (low F~)

$$F_\beta = \frac{(1 + \beta^2)PR}{(\beta^2 P) + R}$$

RCL
RESILIENT COMPUTING LAB

## **However**…. With unbalanced datasets the metrics above have severe weaknesses

► Example

– a dataset which contains 5% of normal data and 95% of anomalies

– A silly algorithm which always answers "anomaly"

► What do we get?

– TP: 95%, FP: 5%, FN: 0%, TN: 0%

– Accuracy: 95%, Precision: 95%, Recall: 100%, F1: 97.5%

► They seem very good scores, correct?

► They seem very good scores, correct?

  – But the algorithm does not "deserve" very good scores (it is a very silly one)

► Therefore, with unbalanced datasets where the ratio of normal/anomaly data points differs a lot, it is better to use a different metric

► **Matthews Correlation Coefficient (MCC)**

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

– Ranges from -1 to 1

– 1 means "perfectly correct classification"

– -1 means "perfectly wrong classification"

- Desirable as well, we just need to do the opposite with respect to what the algorithm says

– 0 means "random guessing", or very bad classification

► In the example, MCC: 0 -> appropriate for a silly algorithm!

► Further readings on metrics

– Hossin, M., & Sulaiman, M. N. (2015). A review on evaluation metrics for data classification evaluations. International Journal of Data Mining & Knowledge Management Process, 5(2), 1.

– Chicco, D., & Jurman, G. (2020). The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. BMC genomics, 21(1), 1-13.