# Predictive models using causal networks

Ana Rita Nogueira, Carlos Ferreira, João Gama

INESCTEC · FCT Fundação para a Ciência e a Tecnologia



Simple example on the difference between correlation and causation (https://amplitude.com/blog/2017/01/19/causation-correlation): Although it can be thought that foul odours can be a focus of diseases, this is not always true: there are cases where there are diseases even without foul odours (for example, the transmission of diseases by a handshake). This is a classic example of mistaken identity.

## Causality is not the same as Correlation!

*"Correlation helps you predict the future, because it gives you an indication of what's going to happen. Causality lets you change the future."* https://neilpatel.com/blog/lean-analytics/

## Causal Inference

Process in which we compare the potential outcomes of an event, in which we have different conditions, but the same variables. This idea can be seen from two different perspectives



This model can be seen from an **Observational Perspective** (p(Wet Floor = yes | Rain = yes)) or a **Interventional Perspective** (p(Wet Floor = yes | do(Rain = yes)))

$$P(x_1, ..., x_n) = \prod_i P(x_i | pa_i)$$

Causal Bayesian Networks traditional representations

## Causal Bayesian Networks

The connection between two nodes represents a causal relationship.
Example: A⊥B|C*, can be represented as A←B→C (B is a common cause of A and C)

## Causal Bayesian Networks
### Local Discovery

Applies Independence tests to a target the variable to define the dependences
**Used to create local graphs (usually in big datasets)**
Example: PC-Simple

## Causal Bayesian Networks
### Global Discovery

Applies both independence tests to all the variables to define the dependences and orientation rules to direct those dependences
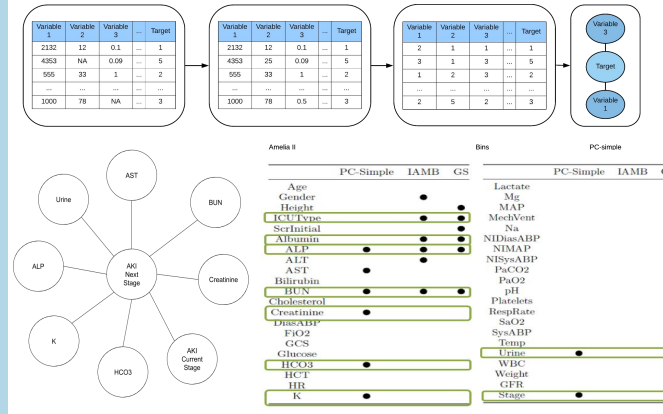**Used to create full graphs**
Example: PC



## Causal Discovery has other applications...

Causal discovery can also be used on for other machine learning tasks, suchs as feature engineering.
In such cases, supposed causal features, that represent the causal relationship between the target variable and the other variables (whether they are its parents and children or members of its markov blanket), can be created to feed more information about the variables' behaviour to a classification algorithm.

## Practical Applications

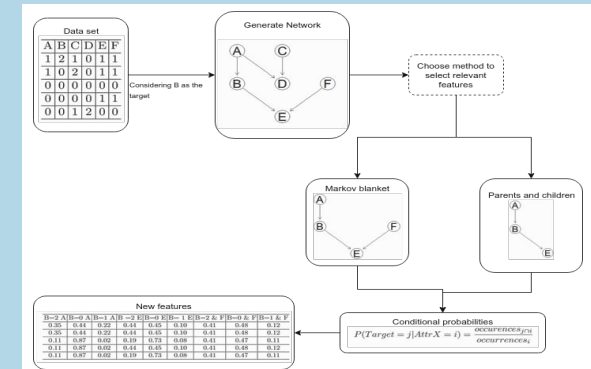Causal discovery can be applied in several areas, such as:

- Medicine
- Economics
- Climatology
- ...

For example, causal discovery can be applied to diagnose Acute kidney Injury disease. In this example, a preprocessing methodology is applied to prepare the data for PC-Simple (used for high-dimensional data). When compared to other local-discovery methods (IAMB and GS, for example) it is possible to see that PC-Simple uncovers more causal relationships. The variables highlighted are causally related to AKI (this was proven through existing literature)
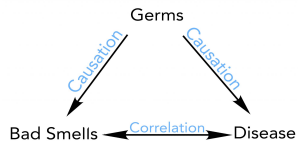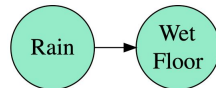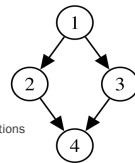
## SELECTED PUBLICATIONS

1. Nogueira A.R., Ferreira C.A., Gama J. (2017) Acute Kidney Injury Detection: An Alarm System to Improve Early Treatment. In: Kryszkiewicz M., Appice A., Ślęzak D., Rybinski H., Skowron A., Raś Z. (eds) Foundations of Intelligent Systems. ISMIS 2017. Lecture Notes in Computer Science, vol 10352. Springer, Cham
2. Nogueira A.R., Gama J., Ferreira C.A. (2018). A Full Causal Parallel Approach to Markov Blanket Variable Selection. Unpublished manuscript.
3. Nogueira A.R., Ferreira C.A., Gama J. (2018). Improving acute kidney injury detection with conditional probabilities. Intelligent Data Analysis. 22. 1355-1374. 10.3233/IDA-173626.
4. Nogueira, A. R., Gama, J., & Ferreira, C. A. (2020, April). Improving Prediction with Causal Probabilistic Variables. In International Symposium on Intelligent Data Analysis (pp. 379-390). Springer, Cham.
5. Nogueira, A. R., Gama, J., & Ferreira, C. A. (2021). Causal discovery in machine learning: Theories and applications. Journal of Dynamics & Games.

# Reminiscence Therapy Improvement Using Emotional Information

Soraia M. Alarcão[1], Carolina Maruta[2], Manuel J. Fonseca[1]

[1] LASIGE, Faculdade de Ciências, Universidade de Lisboa, Portugal
[2] LEL, Centro de Estudos Egas Moniz, Faculdade de Medicina, Universidade de Lisboa, Portugal

## Reminiscence Therapy

- Revisits and stimulates memories from the past using multimedia (images, music, etc.)
- Promotes communication between people with dementia (PwD) and the rest of the world
- Uses preserved abilities to alleviate the experience of failure, and social isolation

## Existing Technological Solutions

- Images used through therapy sessions remain unchanged
- Burden of customizing the therapy is placed on caregivers
- Emotional information on images and the patient's emotions are not considered

## Requirement's Elicitation

- Online Survey / Interviews with Formal caregivers (EN, PT - *ongoing*)
- Online Survey with Informal caregivers (EN, PT, ES - 07/2017 to 12/2018)
  - 603 participants from 39 countries worldwide
  - Therapy performed at Home (92%) during less than 30 min (66%)
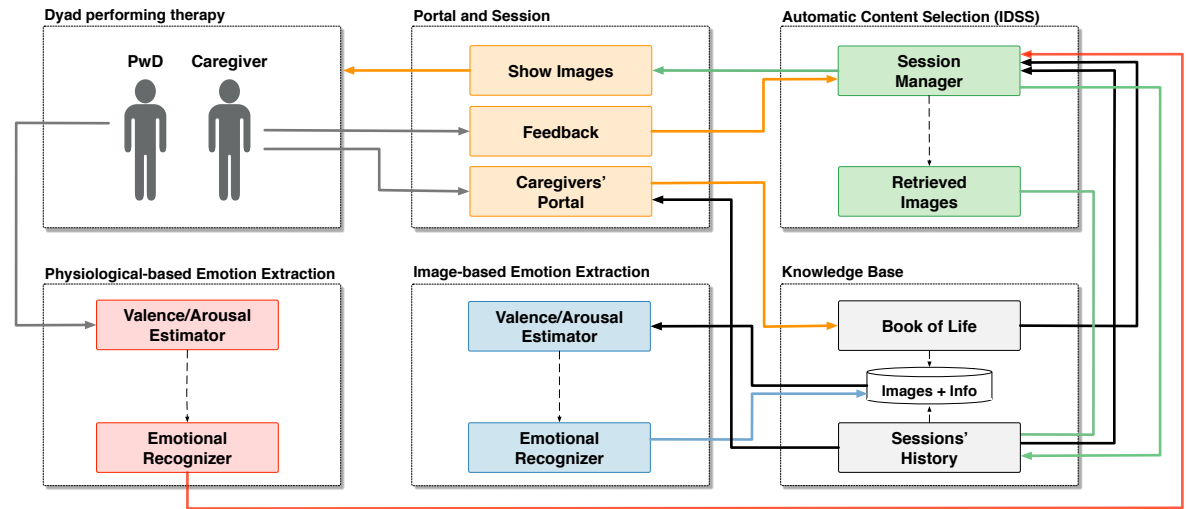  - Negative emotional reactions frequently occur (64%)

### Functional Requirements / Primary Outcomes



Functional Requirements:
- Deliver reminiscence therapy at home
- Automatically retrieve new images
- Gather personalized images
- Recognise PwD emotional reactions
- Adapt current (and future sessions)
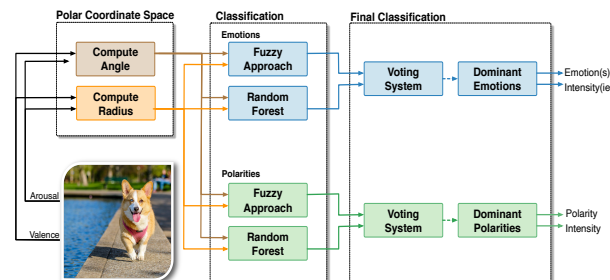
Primary Outcomes:
- Improve PwD cognitive function
- Reduce caregivers' stress
- Reduce caregivers' burden
- Reduce PwD aggressiveness

## User-Centered Technological Solution



## Image-based Emotion Extraction

- **Valence / Arousal Estimator** – ensemble of estimators using visual (color, shape, and texture) and semantic features (tags) with an active learning approach (*ongoing*)
- **Emotional Recognizer** – multi-label supervised classifier that uses valence and arousal to identify emotional polarities and discrete emotions conveyed by images



## Automatic Content Selection (IDSS)

- **Session Manager** – recommendation system based on autobiographical information, emotional reactions to images, and images' emotional information (*to do*)
- **Retrieved Images** – model-agnostic solution based on online learning with expert advice, which dynamically combines several Content-based Image Retrieval systems

# Recommender Systems for Scientific Fields

Marcia Barros[1,2], André Moitinho[2], and Francisco M. Couto[1]

1. LASIGE, Faculdade de Ciências, Universidade de Lisboa
2. CENTRA, Faculdade de Ciências, Universidade de Lisboa
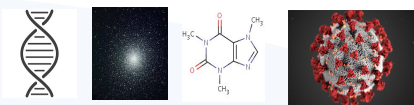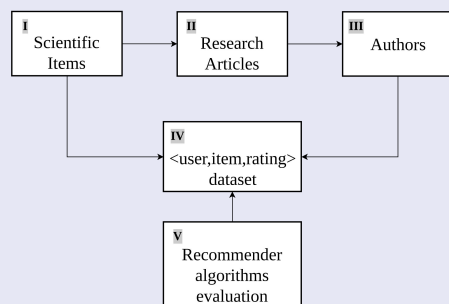
**Scientific Items**

Stored in Large and Complex Databases

- Widely used in movies, music and e-commerce
- Help in the discovery of new scientific items of interest
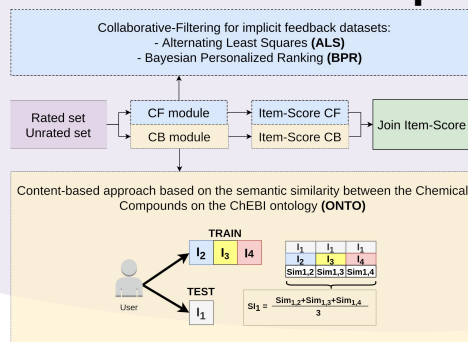- **Challenge:** lack of open access datasets with the users' preferences

**Recommender Systems**

**LIBRETTI**

I Scientific Items → II Research Articles → III Authors

IV <user,item,rating> dataset

V Recommender algorithms evaluation

- Users: authors of scientific publications
- Items: scientific entities
- Ratings: number of publications where an author mentioned an entity
- Case studies in Astronomy using clusters of stars, and Life and Health Sciences, using phenotypes, chemical compounds, diseases, and gene terms, particularly related to the COVID-19 disease

## Hybrid semantic recommender system for chemical compounds

Collaborative-Filtering for implicit feedback datasets:
- Alternating Least Squares (**ALS**)
- Bayesian Personalized Ranking (**BPR**)

Rated set / Unrated set

CF module → Item-Score CF
CB module → Item-Score CB
→ Join Item-Score

Content-based approach based on the semantic similarity between the Chemical Compounds on the ChEBI ontology (**ONTO**)

TRAIN
$I_2$ $I_3$ $I4$
$I_2$ $I_4$
$Sim_{1,2}$ $Sim_{1,3}$ $Sim_{1,4}$

User

TEST
$I_1$

$SI_1 = \frac{Sim_{1,2}+Sim_{1,3}+Sim_{1,4}}{3}$

## Content-based for Astronomical objects

**Open Clusters of Stars** → **Mapping to Gaia sources**

**Similarity between the clusters based on the Gaia features**

**Recommendation of the most similar clusters** ←

## Sequential Enrichment (SeEn) for sequence aware recommendations

Time

**Item1** → **Item2** → **Item_n** → **?**

U LISBOA · C Ciências ULisboa · FCT Fundação para a Ciência e a Tecnologia · LASIGE driven by excellence

# Using Machine Learning in Simulation-Based Data Analytics to Identify Quality-of-Life-increasing Interventions for Prostate- and Breast Cancer Patients

PhD student: Johannes Rust, DFKI Bremen, johannes.rust@dfki.de
Supervisor: Dr.-Ing. Serge Autexier, DFKI Bremen, serge.autexier@dfki.de

## Motivation

Many data analytics in the past years were driven by advancements in machine learning, especially Deep Neural Networks. While being able to automatically identify patterns and complex dependencies, neural networks are usually used as a black box replacing the whole data processing pipeline, thus providing little insight about the inference of a result. On the other side, classical data science disciplines use statistically founded methods. They include no machine learning or only interpretable methods such as linear regression. They are directly interpretable but strongly rely on human expertise. Research to use Deep Learning in combination with simulations was done in some fields [1]. Some of them, like SHAP[2][3] or LIME [4] were used to sample predictions from a model to approximate Feature Attributions. However, these methods do not consider for the domain and structure of the data. When working with data collected in observational studies, for example in the medical field, the influence of an event, exposure or treatment is measured with regards to the outcome variable. We want to make use of this property to make interpretable predictions about the influence of a medical intervention on a patients QoL in the context of the ASCAPE project.

We present an approach to use machine learning models on tabular data in the context of observational studies using methods commonly used in classical data analysis for machine-learning-based simulations. By first analyzing data for correlations and dependencies between interventions and other variables, we provide knowledge interpretable by the user which is then used to run simulations on machine learning models.

## Concept

The approach is applied on a tabular dataset containing input features, variables indicating if an intervention has taken place for the respective sample, and an outcome variable. Our method aims to identify interventions that have a desirable influence on the outcome variable. We split this process into two steps. First, we measure analytically how input features are influenced by each intervention – in the medical context this is known as the average treatment effect (ATE). In the second step, we use the ATE to run Monte-Carlo-Simulations with Machine Learning Models, providing more insight to the result than simply inferring a single prediction.

For each intervention, the database is split into a cohort containing instances that have been exposed to or treated with the respective intervention (test cohort) and one that has not (control cohort). To reduce bias and imbalance of the control and test cohorts, cohort matching is used. After the matching process, the ATE produced by the intervention t is determined for each variable v:

$$ATE_t^v = \frac{1}{N} * \sum_i v_1(i) - v_0(i)$$

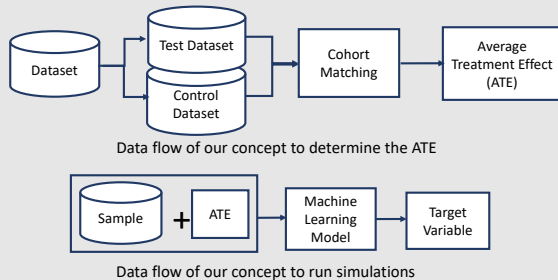The ATE is used to create simulated versions of new sample data as if an intervention has been performed:

$$S_{simulated} = S_{original} + \sigma * ATE_t$$

The factor σ allows us to create simulated samples, where the ATE is higher (σ>1) or lower (σ<1) than expected.

One or more machine learning models are trained to predict the outcome variables based on the input variables. For each possible intervention, the predictions of the ML model M based on the real and the simulated sample are evaluated by comparing the prediction results of both samples:
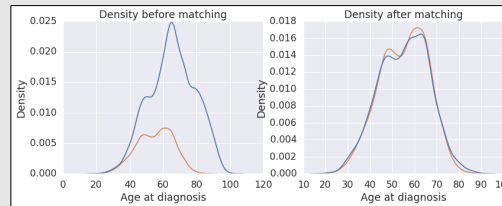
$$Intervention\ Influence = Model(S_{simulated}) - Model(S_{real})$$

For more insight to the model, we can also scale the ATE with a factor σ and visualize how the models' predictions change when the treatment effect is stronger or weaker than expected.



Data flow of our concept to determine the ATE



Data flow of our concept to run simulations
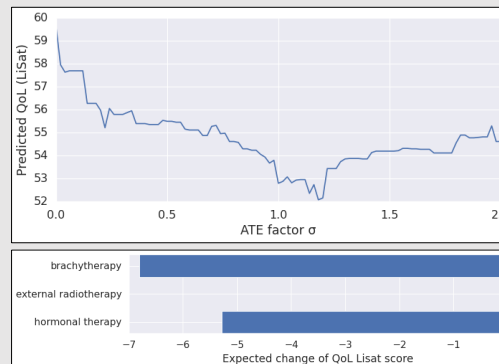
## Analysis of Cohort Matching

We evaluated the performance and suitability of propensity score matching and Mahalanobis-distance matching for our dataset. Both approaches are commonly used when performing cohort matching. Mahalanobis matching pairs samples based on the Mahalanobis-distance of their features [6]. Propensity score matching matches samples on a single scalar value, the propensity score [5]. The propensity score is an estimation on how likely a sample is to belong to the control cohort, which in our scenario is the estimated likeliness of a patient to have received a particular intervention or not. Because propensity score matching was criticized being  as unreliable on some datasets [7], we evaluated the consistency of both methods based on the mean average error (MAE) of the propensity scores of matches. Propensity score matching itself created matches which MAE was only 5.93e-05. When evaluating the propensity scores of matches that were made based on Mahalanobis matching, the MAE was 0.0292. When matches were randomly reassigned, the MAE was 0.1418, which we used as a baseline as an unsuccessful matching. Since the propensity scores after Mahalanobis matching were a lot closer to the Propensity score matchings results than it was to the baseline value, we concluded that the two matching techniques provide somewhat consistent results in our setup.



Example of the distribution of patient ages before and after propensity score matching. The blue graph shows the distribution of the control cohort, the orange graph shows  the distribution of the test cohort.

## Experiments

The approach is used in a real scenario to identify medical interventions that improve the future Quality-of-Life (QoL) of breast- and prostate cancer patients in the context of the "ASCAPE" project. The program aims to design a platform for medical providers and patients that trains machine learning models that predict the patients QoL or the risk of QoL-related issues. The predictive models are made available over a cloud, allowing medical providers to exchange knowledge, while keeping their patients medical data private. To identify medical interventions based on a patients' data and the available machine learning models, we use our approach to identify interventions that are expected to increase the patients' QoL. We use two datasets with retrospectively collected data from breast cancer and prostate cancer patients. The datasets contain medical as well as socio-economic data like household income, which might have an influence on the patients' QoL as well.



Predicted QoL Lisat score of a model simulated depending on the factor σ of the treatment effect in the interval of [0, 2.0]. σ=1 means the ATE is as strong as expected, σ=0.5 means it is only half a strong, etc.



Predicted change of a patients QoL score after a treatment. In this case, none of the three treatments shown increases the patients QoL.

## Federated Learning

Our approach can be implemented to work with little human input in a federated scenario, providing suggestions for medical interventions to medical experts and patients and being continually updated whenever new anonymized patient data is added to the study dataset. While there are various approaches and implementations to train Machine Learning Models like Neural Networks in a federated manner, cohort matching and the calculation of the average treatment effect must be performed individually by each federation partner. Each federation partner performs cohort matching on his own dataset for each intervention. The ATE is calculated and sent to a central federated learning coordinator. Since the ATE is calculated additively, multiple ATEs can easily be combined by simply averaging them:

$$ATE_{all} = \frac{1}{\#ATEs} * \sum_i ATE_i$$

The ATEs can be collected for each variable independently, which means that not every federated learning partner must have data about a certain intervention but can still profit from the measured ATEs provided by other partners.

## Key Findings and Conclusions

We presented a way to use machine learning models in a simulation-based setup. By using classical data analysis usually used in medical sciences, we could identify the average treatment effect. For observational cross-over-studies, we showed that cohort matching with Mahalanobis Distance Matching and Propensity Score Matching is suitable and also needed to reduce biases and imbalances in our dataset. We furthermore showed that the average treatment effect can be used as a source to effectively explore the feature space of a machine learning model. Given that the machine learning model can be trained in a federated manner, our approach can as well be used in this domain.

However, evaluations on two datasets from the medical fields showed that the approach is limited by the ML model performance. Further open research questions lie in how the ML model's inner decision process itself can be made more transparent and how our approach can be evaluated regarding its accuracy.
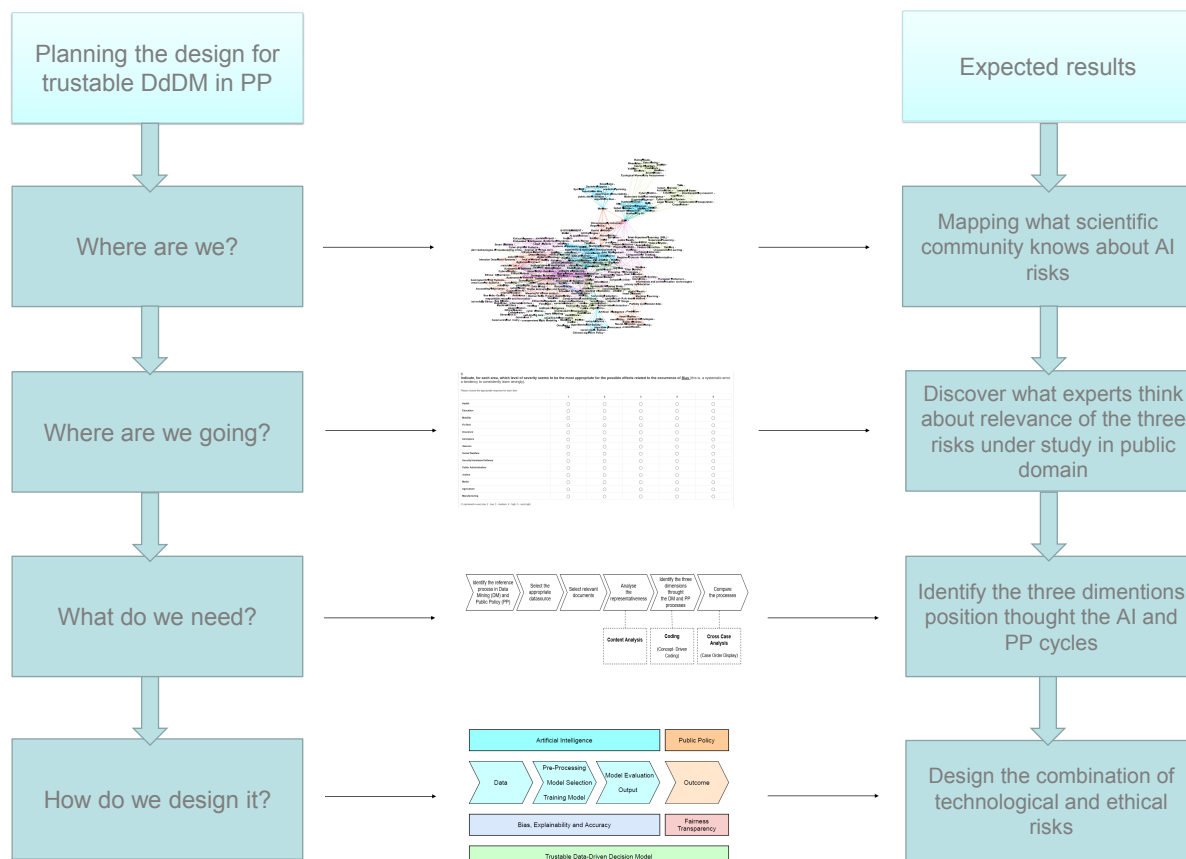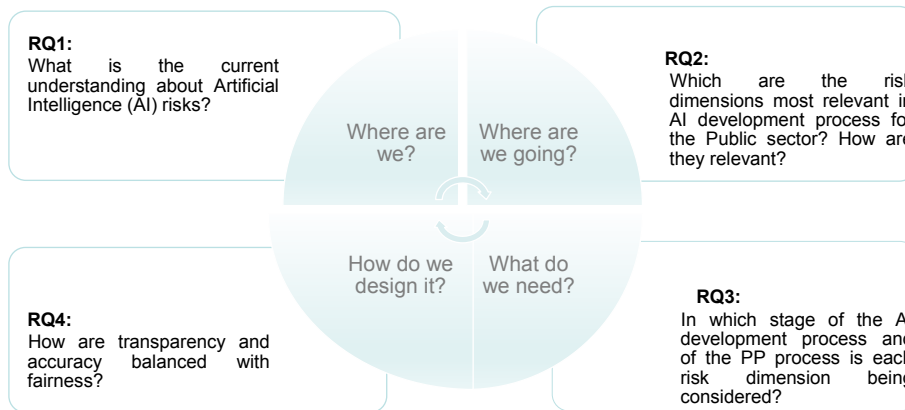
## References

[1] von Rueden, Laura, et al. "Combining machine learning and simulation to a hybrid modelling approach: Current and future directions." *International Symposium on Intelligent Data Analysis*. Springer, Cham, 2020.
[2] Lundberg, Scott, and Su-In Lee. "A unified approach to interpreting model predictions." *arXiv preprint arXiv:1705.07874* (2017).
[3] Lundberg, Scott M., et al. "From local explanations to global understanding with explainable AI for trees." *Nature machine intelligence* 2.1 (2020): 56-67.
[4] Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. ""Why should i trust you?" Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016.
[5] Caliendo, Marco, and Sabine Kopeinig. "Some practical guidance for the implementation of propensity score matching." *Journal of economic surveys* 22.1 (2008): 31-72.
[6] Rubin, Donald B. "Bias reduction using Mahalanobis-metric matching." *Biometrics* (1980): 293-298.
[7] King, Gary, and Richard Alexander Nielsen. "Why propensity scores should not be used for matching." (2019).

# Trustability in Data-Driven Decision Models (DdDM) for Public Policy (PP)

Sónia Teixeira, José Coelho Rodrigues, João Gama

**Motivation:** The adoption of algorithmic systems, in particular the data-driven decision models, in Public domain has raised questions about the risks associated with this technology, from which ethical problems may emerge.

**Objective:** Principles to design trustable data-driven decision models for Public Policy.

**RQ1:**
What is the current understanding about Artificial Intelligence (AI) risks?

**RQ2:**
Which are the risk dimensions most relevant in AI development process for the Public sector? How are they relevant?

Where are we?

Where are we going?

How do we design it?

What do we need?

**RQ4:**
How are transparency and accuracy balanced with fairness?

**RQ3:**
In which stage of the AI development process and of the PP process is each risk dimension being considered?

| Planning the design for trustable DdDM in PP | | Expected results |
|---|---|---|
| Where are we? |  | Mapping what scientific community knows about AI risks |
| Where are we going? |  | Discover what experts think about relevance of the three risks under study in public domain |
| What do we need? |  | Identify the three dimentions position thought the AI and PP cycles |
| How do we design it? |  | Design the combination of technological and ethical risks |

# SELECTED PUBLICATIONS

1. Teixeira, S.; Rodrigues, J.; Gama, J.; Veloso, B.. "Challenges of Data-Driven Decision Models: Implications for Developers and Public Policy Decision Makers". In Advances in Urban Design and Engineering - Perspectives from India. Springer Nature, 2020. (Accepted for publication)
2. Teixeira, S.; Gama, J.; Amorim, P.; Figueira, G.. "Trustability in Algorithmic Systems Based on Artificial Intelligence in the Public and Private Sec-tors", ERCIM News, 2020, https://ercim-news.ercim.eu/images/stories/EN122/EN122-web.pdf.
3. Teixeira, S.; Rodrigues, J.; Gama, J.. "The Risks of Data-Driven Models as Challenges for Society". IEMS '20 — 11th Industrial Engineering and Management Symposium: The Impact of DEGI Research on Society, Porto, 2020

# Automated Privacy Preserving Strategies

**Tânia Carvalho | Nuno Moniz**
Faculdade de Ciências da Universidade do Porto, TekPrivacy, INESC TEC

## Motivation

- Suppose a hospital needs a research team to analyse data to help fight Covid-19.
- The hospital has to provide the data for further processing and analysis.
- However, the hospital cannot share the data because it contains personal information.
- The hospital needs to guarantee its patients' privacy by de-identifying such data.
- Finally, data can be shared with the team that will analyse it.

## Objective

Create a framework that automatically determines and applies a de-identification approach with high privacy assurances without significant impacts in predictive utility.

## Approach

- Data de-identification:
  - Ensure that the private information of na individual is not compromised;
  - Ensure enough predictive utility for data to be used without too much information and performance loss;
  - Apply transformation techniques to reduce detail, suppress or distort information.

- However, there are some disadvantages in the de-identification process:
  - Multiple steps and tests;
  - The learning process requires finding the best statistical model and its configurations;
  - Significant time budget.

- To overcome such obstacles, we propose to automate the de-identification process using AutoML.

## AutoML

- Automate the machine learning workflow, p.e., hyperparameterisation, model and feature selection.

- Meta-learning and Optimisation will be explored and evaluated concerning their ability to ensure that data is de-identified and maintaining as much as possible the predictive utility of the data.
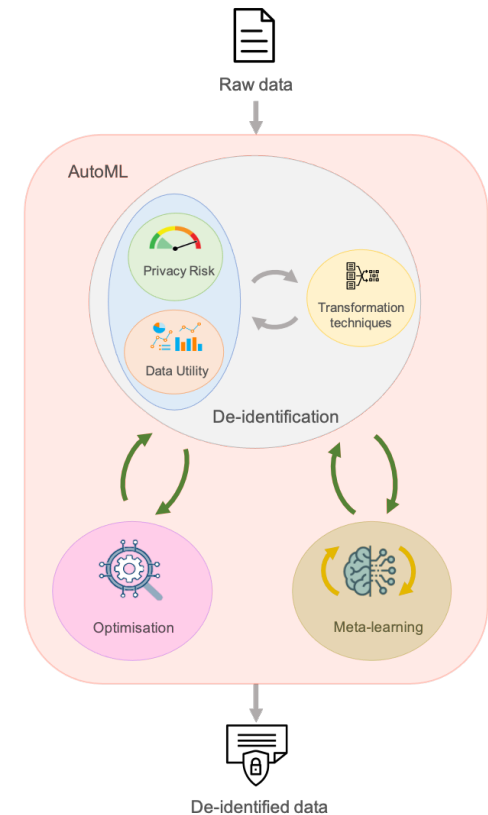


Fig.1 Workflow of our proposed methodology.

# PERSONALIZE COUNTERFACTUAL EXPLANATION INTERACTIVELY IN REAL-TIME USING GENERATIVE ADVERSARIAL NETWORKS (GANS)

**Ahmed A. Fares,**[1,2] **João Mendes Moreira ,**[1,2] **Filipe F. Correia,**[1,2]
[1]**Faculty of engineering, University of Porto, Portugal,** [2]**INESC TEC, Porto, Portugal.**

## MOTIVATION

The explainability of machine learning in decision-making is considered one of the essential factors to trust machine learning models. The predictions made by black-box models should always be verified whenever it affects human affairs. The decisions should be well explained to the relevant stakeholders, who deal directly with the system from technical and operational users until end-users, based on each one's knowledge and background. Arya et al., 2019.

## WHAT IS EXPLAINABLE ARTIFICIAL INTELLIGENCE (*XAI*)?

Users of AI applications are interested in getting accurate predictions as long as being convinced by the reason(s) behind receiving that particular value.

*XAI* offers several types of explanations: Verma, Dickerson, and Hines, 2020

- ► Feature importance (Local surrogate (LIME), Shapley Values).
- ► Leading examples (KNN predictions of the training data).
- ► Counterfactual explanations (How an instance has to be modified to change its prediction significantly?).

## COUNTERFACTUAL EXPLANATIONS (EVALUATION CRITERIA)

- ► Realism (How the counterfactual is far from the distribution of real data)
- ► Actionability (How counterfactual relative to the input data point).
- ► Latency (ms) (Computational time needed to generate counterfactuals).

State-of-the-art methods comparison using Pima Indians Diabetes dataset (lower values correspond to better solutions). Nemirovsky et al., 2020

|  | RGD | CSGP | CounterRGAN |
|---|---|---|---|
| **Realism** | $2.20 \pm 0.24$ | $2.03 \pm 0.11$ | $\mathbf{1.79 \pm 0.11}$ |
| **Latency (ms)** | $1195 \pm 5.65$ | $3211 \pm 11.65$ | $\mathbf{42.74 \pm 4.28}$ |
| **Actionability** | $1.64 \pm 0.20$ | $\mathbf{1.14 \pm 0.19}$ | $6.91 \pm 0.43$ |

## IMPROVE ACTIONABILITY

- ► **WHY?**

  Realism is not enough. Actionability describes whether the suggested changes are interpretable and reasonable according to the input case. Allowing the user to ask "what if" questions, would direct the search process and reduce the steps to match the desired class.
- ► **HOW?**
  - ► User interaction with GANs. Nemirovsky et al., 2020
  - ► Interactive reinforcement leaning. Arzate Cruz and Igarashi, 2020
  - ► Variational autoencoder. Kingma and Welling, 2019
  - ► Multi-objective optimization. Dandl et al., 2020

## WHY INTERACTIVE EXPLAINABILITY?

- ► Personalized explanations (different explanations based on the individual user's profile).
- ► Inferring the explanations in real–time (no need to explicitly build them in the system).
- ► Mirror the user's mental model onto the system (online or offline)

## PROPOSED APPROACH

The majority of existing techniques make rough approximations, which leads to limited performance regarding actionability. To overcome these limitations, we are going to start the experiments by merging user feedback with the Generative Adversarial Networks (GANs) methodology to generate counterfactuals. This method has a high potential to improve the actionability to be used in real-time while considering realism.

## REFERENCES

Arya, Vijay et al. (2019). *One Explanation Does Not Fit All: A Toolkit and Taxonomy of AI Explainability Techniques*. arXiv: 1909.03012 [cs.AI].

Arzate Cruz, Christian and Takeo Igarashi (2020). "A Survey on Interactive Reinforcement Learning: Design Principles and Open Challenges". In: *Proceedings of the 2020 ACM Designing Interactive Systems Conference*. DIS '20. Eindhoven, Netherlands: Association for Computing Machinery, pp. 1195–1209. ISBN: 9781450369749. DOI: 10.1145/3357236.3395525. URL: https://doi.org/10.1145/3357236.3395525.

Dandl, Susanne et al. (2020). "Multi-Objective Counterfactual Explanations". In: *Lecture Notes in Computer Science*, pp. 448–469. ISSN: 1611-3349. DOI: 10.1007/978-3-030-58112-1_31. URL: http://dx.doi.org/10.1007/978-3-030-58112-1_31.

Kingma, Diederik P. and Max Welling (2019). "An Introduction to Variational Autoencoders". In: *Foundations and Trends® in Machine Learning* 12.4, pp. 307–392. ISSN: 1935-8245. DOI: 10.1561/2200000056. URL: http://dx.doi.org/10.1561/2200000056.

Nemirovsky, Daniel et al. (2020). *CounteRGAN: Generating Realistic Counterfactuals with Residual Generative Adversarial Nets*. arXiv: 2009.05199 [cs.LG].

Verma, Sahil, John Dickerson, and Keegan Hines (2020). *Counterfactual Explanations for Machine Learning: A Review*. arXiv: 2010.10596 [cs.LG].

# LASIGE

data and systems intelligence

## Learning Prognostic Models Using a Mixture of Biclustering and Triclustering

Diogo F. Soares and Sara C. Madeira

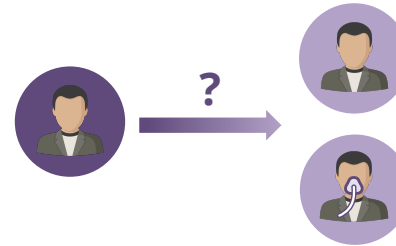LASIGE, Faculdade de Ciências, Universidade de Lisboa, Portugal

U LISBOA    Ciências ULisboa

FCT Fundação para a Ciência e a Tecnologia    LASIGE driven by excellence

## Motivation

- Subspace Clustering (biclustering and triclustering) allows to extract discriminative patterns from data
- Learning from high dimensional data can be hard task
  - Curse of Dimensionality
  - Feature Selection
- Some features can be relevant only when grouped with some other features
  - Feature selection prevent finding these groups
- Subspace clustering can:
  - overcome limitations of feature selection
  - improve model performance and interpretability
- In clinical domain, biclustering and triclustering:
  - are promising approaches to identify groups of patients with correlated features along time
  - retrieve disease progression patterns

## ALS Case Study

- Amyotrophic Lateral Sclerosis is a highly heterogeneous disease
- Patients develop respiratory insufficiency
- Rapid administration of non-invasive ventilation is effective in prolonging and improving quality of life

**Predict if a given ALS patient will need NIV within 90 days of the last clinical appointment, using data from patient's follow-up**

## Data

Lisbon ALS Clinic Dataset Hospital Santa Maria

940 patients

9 static features
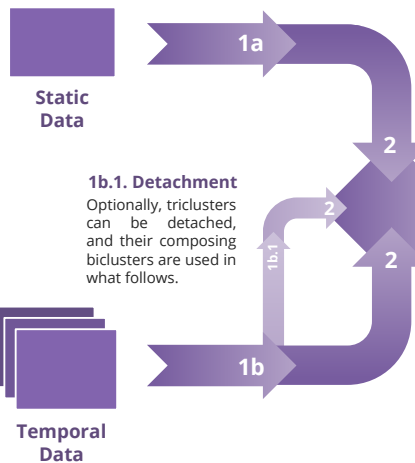12 temporal features
3, 4 and 5 consecutive snapshots (CS)

**Scan to know more about Data Preprocessing**

# BICLUSTERS AND TRICLUSTERS IMPROVE MODEL INTERPRETABILITY

### 1a. Biclustering
Biclustering two-way static data to obtain biclusters that will be used as features.

### 4. Learning Instances
The computed similarities matrices are merged into one with objects and the discovered patterns as features. This matrix, coupled with the correspondent labels, is fed to the classifier.

### 2. Virtual Patterns
Compute the bicluster/tricluster most representative pattern, which is the mean object within-cluster and represents its tendency.

Static Data

**1b.1. Detachment**
Optionally, triclusters can be detached, and their composing biclusters are used in what follows.

**Predictive Model**

### 3. Similarities Matrices
Similarities matrices are computed using the virtual patterns and the object patterns (using same features) according to 2 proposed methods:
1. Euclidean Distance
2. Pearson Correlation

Temporal Data

**Using the Model**

### 1b. Triclustering
Triclustering three-way static data to obtain triclusters that will be used as features.
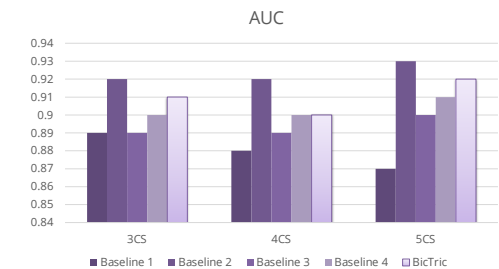
### 5. Classifier
Finally, we learn the predictive model using a Random Forests and evaluating using a repeated stratified 10-fold cross-validation.

## Results

**Baseline experiments were performed using:**
1. Only the last appointment of each set + static features
2. All appointments and static features
3. Triclusters only (discarding static features)
4. Triclusters together with static features

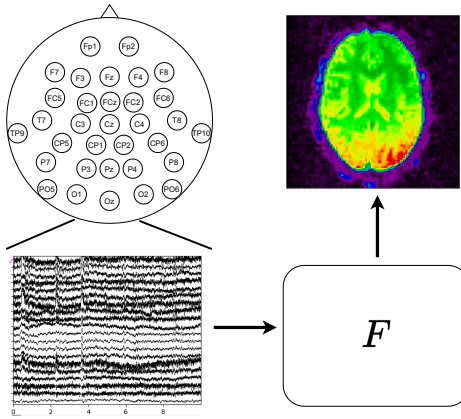**Our approach did not outperformed baseline 2 in predictability but achieved almost the same results.**

AUC

Baseline 1    Baseline 2    Baseline 3    Baseline 4    BicTric

**Most relevant patterns used by the model help clinicians to better understand the patient prognostic prediction**

# EEG to FMRI syntehsis

David Calhas[1,2]

[1]Instituto Superior Técnico, Universidade de Lisboa
[2]INESC-ID, Lisboa

## Introduction

Electroencephalography (EEG) measures the **neural activity at the scalp level**. It is known for its cheapness and ambulatory traits. Functional Magnetic Resonance Imaging (fMRI) measures the **blood supply of the brain**, which in contrast to EEG, is expensive [1]. Efforts have been made to relate these two modalities through connectivity properties [2] and we are even starting to see the application of machine learning (ML) methods to bridge this gap [3]. My **hypothesis** is that with the advances of ML in the last decade [4, 5], **the gap between EEG and fMRI is able to be bridged** with a consequent great impact on **cost reductions and availability**.

The objectives of this work are the following:

- Model a mathematical function that **captures the style features of the fMRI signal**, $\vec{y}$;
- Use well known neural architectures to perform the **encoding mapping function between EEG and fMRI**;
- Use automated machine learning techniques to avoid a learning bias, so that **the mapping between EEG and fMRI** does not rely on domain knowledge, and **purely on the algorithmic perspective**;

- Define $F : \vec{y} = F(\vec{x})$, being $\vec{x}$ a set of EEG features and $\vec{y}$ an fMRI volume.
- **Evaluate the synthesized fMRI signal** in classical disease diagnostic settings.

## Data

Nowadays, an increasing number of **neuroimaging datasets are becoming pubicly available**. This project will be validated on publicly available datasets, published by previous studies, that contain **simultaneous EEG and fMRI recordings**, either resting state or task based settings. **Openneuro** offers an extensive database with published datasets, with several containing simultaneous EEG and fMRI recordings.

## Methods

***Learning fMRI style signature***. During the last decade, there have been efforts to solve known differential equations with the help of data driven methods, such as neural networks (NN). Chen et al. [4] proposed a new **initial value problem solver**, that is able to **solve a differential equation** quickly, allowing its use with the optimization of more parameters, as it happens in the learning phase of a NN. We propose modelling latent representations of fMRI signal using a **neural ordinary differential equation** (neuralODE). Under the hypothesis that **an entire latent fMRI signal representation** (spatial line traversing the volume or a single voxel) **can be represented only using the initial value**.
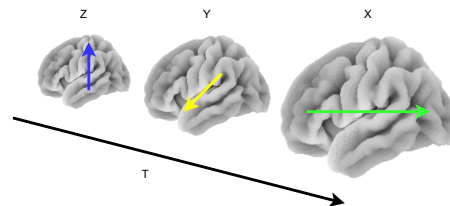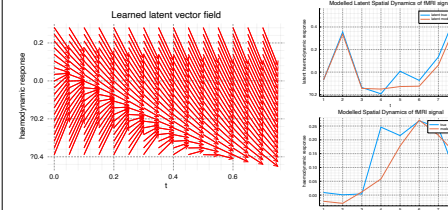


Figure 1: **fMRI dimensions.**

***Shared latent space***. EEG and fMRI have a range of studies that extract neural correlates. Since the two evaluate **very different brain processes**, at different time scales, the **transfer function between the two is not linear**. Mapping the two signals, $\vec{x} \in \mathbb{R}^{S_1}$ and $\vec{y} \in \mathbb{R}^{S_2}$, is not trivial because of the **high dissimilarity in structure**, i.e. $S_1 \neq S_2 \wedge \nexists \epsilon : \epsilon(S_1) = S_2$ with $\epsilon$ applying a small change in structure.

***Neural architecture search (NAS)***. Automated machine learning is capturing a lot of attention due to it being **a way of removing the learning bias**. In this work, we explore this avenue and aim at developing a method that is **capable of generating neural network architectures**. With such a method, one will then perform a **search among the generated networks and choose the one that performs best**, $F_1$. This architecture is the candidate to map, $\vec{x}$, such that $S$ is the shared space and $F_1(\vec{x}) \in \mathbb{R}^S \wedge F_2(\vec{y}) \in \mathbb{R}^S$, with $F_2$ being a Resnet like architecture.

## Preliminary results



a: **Learned vector field.**

b: **Latent and original space representation.**

Currently, we are able to **learn a mathematical function that contains the style properties of the fMRI** signal. We define four different settings: $x$-axis, $y$-axis, $z$-axis and *temporal*-axis. And independently learn from each one of them in **an encoder decoder architecture**. The **feature representation from the encoder is modeled by a neuralODE and the modeled neuralODE signal is mapped back to the original space by the decoder**.

In Figure 2a is shown the **vector field learned from adjacent lines that traverse the fMRI volume in the $z$-axis direction**. A latent space regression, as well as its corresponding decoded representation, are shown in Figures 2b. This model can be used to either **model complete brain volumes using only the near the scalp haemodynamic response** and/or **correct a modelled brain volume**.

## Next steps

As of now, we are validating the significance of the learned mathematical function that describes the style of fMRI dynamics. The results are meaningful, but **are yet to be generalized to the whole fMRI volume**. The next step is to **combine this mathematical function with a simple autoencoder task or even a direct mapping from EEG to fMRI**.

## References

[1] Hans Op de Beeck and Chie Nakatani. *Introduction to Human Neuroimaging.* Cambridge University Press, 2019.

[2] Rodolfo Abreu, João Jorge, Alberto Leal, Thomas Koenig, and Patrícia Figueiredo. Eeg microstates predict concurrent fmri dynamic functional connectivity states. *Brain topography*, 2021.

[3] Xueqing Liu, Linbi Hong, and Paul Sajda. Latent neural source recovery via transcoding of simultaneous eeg-fmri. *arXiv*, 2020.

[4] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. *NeurIPS*, 2018.

[5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

## Acknowledgements

## Contact Information

- Email: david.calhas@tecnico.ulisboa.pt

# GLYCOPROTEOGENOMICS CHARACTERIZATION IN COLORECTAL CANCER USING DEEP LEARNING

Daniel Mateus Gonçalves

Supervisors: Rafael Costa, Rui Henriques, Alexandre Ferreira

Instituto Superior Técnico - Universidade de Lisboa

## Summary

Colorectal cancers (CRC) are amongst the top-ranked tumors in terms of mortality due to late diagnosis and significant molecular heterogeneity, which has frustrated hopes for accurate patient stratification and the introduction of new therapeutics. Several attempts have been made to provide machine learning models built upon genomics or transcriptomics data alone. However, a precise integration of multi-omics data is required to improve the stratification capacity of these models and provide cancer-specific molecular targets. The purpose of this work is to devise deep learning to analyze the various omic data layers in an integrative way. The focus will be channeled towards glycoproteogenomics, considered a promising field in pan-cancer research. This setting may contribute to fulfilling important systems oncology research gaps leading to the identification of novel biomarkers and new targets to fight CRC, helping to guide therapeutic interventions.

***Keywords:*** Systems Oncology; Multi-Omics Analysis; Glycoproteogenomics; Cancer neoantigens; Machine Learning; Deep Learning.

## State of the Art

**Colorectal Cancer** (CRC) is the third most common cancer type and the fourth leading cause of cancer-related death in the world [1]. Studies reveal extensive **molecular heterogeneity** upon the analysis of genomic, transcriptomic, and proteomic data [4, 10]. However, the growing amounts of data have yet to bring new biomarkers and drug targets to clinical practice [9]. The differences between tumor and normal tissue have not been systematically characterized in large cohorts, pointing to the importance of the **integrative analysis** as the next step [9].

The use of machine learning and the challenges in multi-omics analysis are well known, pointing to Unsupervised Deep Learning as a promising route [3, 2]. **Autoencoder Neural Networks** (Fig. 1) are described as powerful tools in the scene of single or multi-omics data analysis due to their superior feature extraction capabilities.
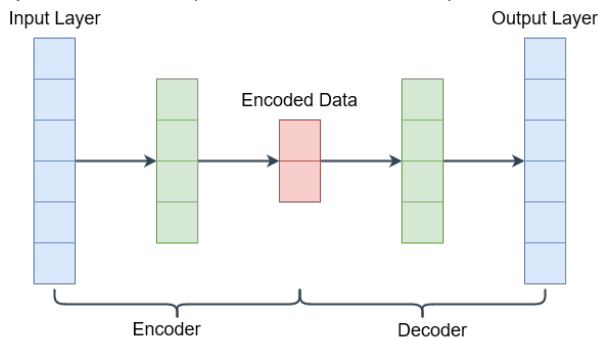


Fig. 1: Illustrative Autoencoder network architecture.

More recently, **Variational Autoencoders** (VAE) have been proposed, leveraging previous capabilities with a generative approach [7]. This method is capable not only to encode complex data sources but also to decode it in a way where the distance between similar generated instances is minimized. VAEs have been successfully used in recent studies with multiple **unsupervised objectives**, like cancer identification, molecular subtype identification, and survival analysis [6].

Another deep learning method described in the literature are **Self Organizing Maps** (SOM) [8]. SOM networks (Fig. 2) are trained in an unsupervised manner and are useful for their dimensionality reduction and **unique visualization capabilities**, contributing to more explainable and interpretable Artificial Intelligence. SOMs have been used to identify biologically important differences in subsets of proteins helping identify potential drug targets [5].
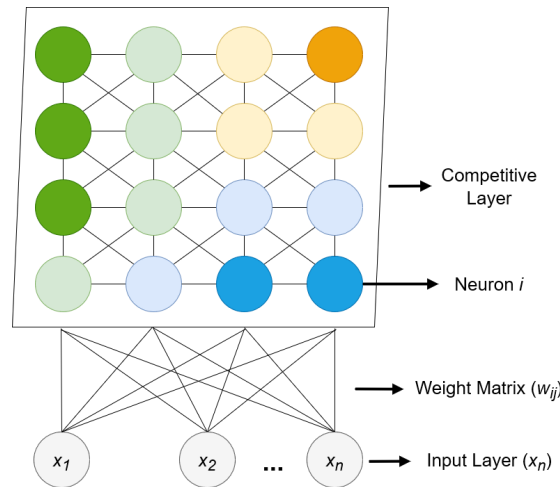


Fig. 2: Self Organizing Map Architecture Illustration.

## Goals

This work tackles the task of glycoproteogenomics integrative analysis in colorectal cancer (CRC) using deep learning. Two major approaches are explored: i) Analysis of individual omics data (conventional approach); ii) Integrative glycoproteogenomics data analysis.

In this context, this work pursues 2 major goals:

1. Identification and validation of **new molecular signatures** associated with CRC and neoantigens holding potential for precise cancer targeting;

2. Patients' **phenotyping and stratification**, as well as CRC subtype identification.

## Tasks

The Clinical Proteomic Tumor Analysis Consortium (CPTAC) Data Portal is a public repository of proteomic sequence **datasets** containing data for more than 450 colorectal cancer cases. This data will be the base of our work, but we'll actively look for more sources.

The **exploration and preprocessing** of the data will be a crucial step for the success of this work. Enabling more complex tasks, like combining the various omics layers and training deep learning models.

**Unsupervised deep learning** will be the domain of Artificial Intelligence employed by our methodology. We plan to explore at least two major techniques, Autoencoders and Self Organizing Maps (SOM):

- **Autoencoders** are a modular and efficient way to learn an encoding. Their potential can be explored to support the integrative component, capable of feature extraction and dimensionality reduction. Furthermore, **Variational Autoencoders** arise as a generative alternative, adding the possibility of training the model to generate clustered instances.

- **SOMs** have also been described as powerful tools, and are recently being tested in pair with other methods. The plan is to test them as a **standalone** method, and later in combination with an Autoencoder to process its encoded layer, **combining** the quality feature extraction capability of Autoencoders with the dimensionality reduction, clustering, and visualization capabilities of SOMs.
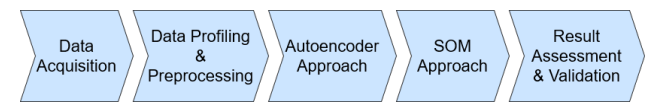


Fig. 3: Work steps pipeline.

## References

[1] Melina Arnold et al. "Global patterns and trends in colorectal cancer incidence and mortality". In: *Gut* 66.4 (2017), pp. 683–691.

[2] Nupur Biswas and Saikat Chakrabarti. "Artificial Intelligence (AI)-Based Systems Biology Approaches in Multi-Omics Data Analysis of Cancer". In: *Frontiers in Oncology* 10 (2020).

[3] Travers Ching et al. "Opportunities and obstacles for deep learning in biology and medicine". In: *Journal of The Royal Society Interface* 15.141 (2018), p. 20170387.

[4] Justin Guinney et al. "The consensus molecular subtypes of colorectal cancer". In: *Nature medicine* 21.11 (2015), p. 1350.

[5] Clara Higuera, Katheleen J Gardiner, and Krzysztof J Cios. "Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome". In: *PloS one* 10.6 (2015), e0129126.

[6] Muta Tah Hira et al. "Integrated multi-omics analysis of ovarian cancer using variational autoencoders". In: *Scientific reports* 11.1 (2021), pp. 1–16.

[7] Diederik P Kingma and Max Welling. "An introduction to variational autoencoders". In: *arXiv preprint arXiv:1906.02691* (2019).

[8] Teuvo Kohonen. "Self-organized formation of topologically correct feature maps". In: *Biological cybernetics* 43.1 (1982), pp. 59–69.

[9] Suhas Vasaikar et al. "Proteogenomic analysis of human colon cancer reveals new therapeutic opportunities". In: *Cell* 177.4 (2019), pp. 1035–1049.

[10] Bing Zhang et al. "Proteogenomic characterization of human colon and rectal cancer". In: *Nature* 513.7518 (2014), pp. 382–387.

## Acknowledgements

# Learning prognostic biomarkers from three-dimensional biomedical data of psychiatric disorders

Leonardo Alexandre[1,3,4] <leonardoalexandre@tecnico.ulisboa.pt>, Rafael S. Costa [2,3], Rui Henriques [1,4]

**1** Instituto Superior Técnico, Universidade de Lisboa, Lisboa, Portugal, **2** LAQV, NOVA School of Science and Technology, Universidade NOVA de Lisboa, Portugal, **3** IDMEC, Instituto Superior Técnico, Universidade de Lisboa, Portugal, **4** INESC-ID, Lisboa, Portugal

## Goals

The proposed PhD research topic aims at developing **machine learning approaches to identify prognostic biomarkers** in psychiatric disorders and support therapeutic choices from cohorts with available **social behavior**, **neuroimaging**, **psychocognitive tests**, **omics**, **and undertaken therapeutics data**. These biomarkers can be represented by **three-dimensional patterns**, with the dimensions being *patients-variable-time*, extracted from **heterogenous biomedical data** present in neurophysiological longitudinal studies. Focus will be placed on the **discovery**, **statistical assessment** and **validation** of biomarkers in order to guarantee the **usability** by medical professionals. The proposed work will create a **novel computational approach** able to pre-process heterogenous biomedical data, **exhaustively search the three-dimensional space**, and **guarantee the statistical significance of the biomarkers** discovered. After validating the biomarkers discovered we will **extend the developed approaches** towards predictive setting **to support diagnostics and therapeutic choices of new patients**.
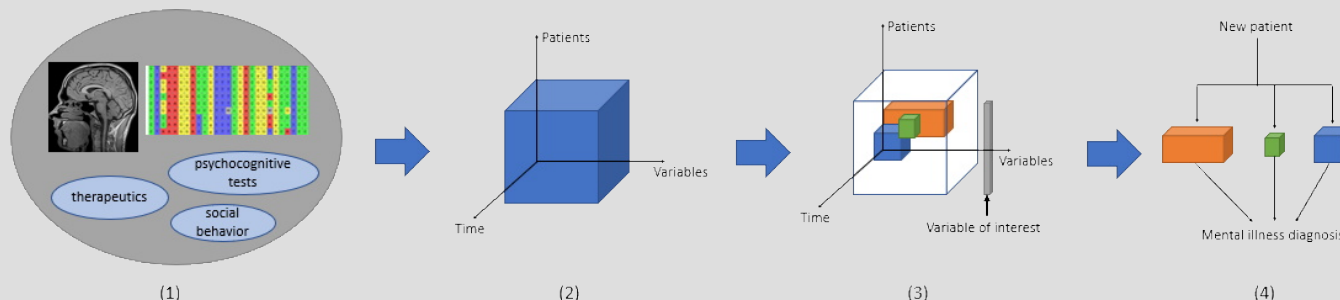
## Detailed Description



Figure 1: **(1)** Different types of data are collected such as neuroimages, omics, undertaken therapeutics, psychocognitive tests, and social behavior. **(2)** Data is preprocessed and structured with symbolic techniques transforming it to the format of *patient-variable-time*. **(3)** Computational approaches such as triclustering are applied searching for pattern with discriminative power towards a mental illness or therapeutic outcome (variable of interest). **(4)** A predictive model which employs patterns to help diagnose and treat new patients.

The **core contribution** of the prospective PhD research will be **defining prognostic biomarkers in psychiatric disorders**. The patient **being diagnosed** and treated with appropriate medication according to their neurobiological profile is **critical to prevent morbidity** and in some cases **mortality** [11], but, due to many **symptoms overlapping** between mental disorders [8,10], this task becomes **hampered**. This profile can be extracted from **three-dimensional patterns**, in the form of *patient-variable-time*, contained in a wide range of **heterogeneous longitudinal data** from multiple visits such as **social behavior**, **psychocognitive tests**, **neuroimaging**, **omics**, and **undertaken therapeutics**. The data from psychiatric disorders cohorts will be obtained from **public databases**, like OpenNeuro (https://openneuro.org/) and Schizoconnect (http://schizconnect.org/). The time component of the patterns allows it to characterize the neurobiological changes from psychiatric follow-up and accompany the response to medication. In **schizophrenia patients** these would present as a **deficit in the socio-cognitive skill** [19], an **interplay between receptors in genes** [9], or **imaging markers** [1].

**To achieve the aforementioned goal, state-of-the-art computational approaches will be extended** to address the following challenges: 1) the **heterogeneous nature of the input data**, 2) the need to scale and **exhaustively search the three-dimensional space** outputting patterns with different homogeneity criteria and intrinsic noise, 3) the importance of assessing of **statistical significance of the discovered patterns**.

To deal with the heterogeneous **data inputted, symbolic techniques** that consider border values will be implemented. To **explore the three-dimensional space** to search for pattern with different homogeneity, multiple approaches to exhaustive **triclustering algorithms** will be implemented and improved upon[4,6,7]. Given the difficulty of promoting the **scalability** of such algorithms, **heuristics** can be employed while still **guaranteeing optimal results**. **Statistical tests to assess the statistical significance and discriminative power** of data associations have been studied and applied when searching for patterns in two-dimensional data and can be extended to triclustering [3].

**The research** done in the context of my PhD is aimed to be **submitted to top conferences** and **create more discussion** on the important topic that is mental illness. The work will be developed so that **results will outpour to actual healthcare providers** and will explore synergies with **ongoing national and international R&D projects**, establishing a close contact with their research teams and contributing to address their open challenges.

## State-of-the-art

The **number of patients diagnosed** with depression, attention deficit and hyperactivity disorder (ADHD), anxiety, bipolar disorder, and/or schizophrenia **is considerably rising** (https://www.nimh.nih.gov/health/statistics/mental-illness.shtml).

In the **biomedical domain, three-dimensional data can describe a patient's profile** or multiple patients profile. Approaches such as **triclustering**, the grouping of objects across three dimensions, are **essential to better delimit diseases, understand disease progression**, and **responses to stimuli and drugs**. Henriques and Sara [2], provide an extensive structured view on the triclustering problem and existing contributions. They **stress the need for upcoming contributions**, such as integrative approaches able to combine the potentialities of different algorithms, new triclustering searches for sparse 3D data, and new principles to guarantee the statistical significance of temporal patterns.

**Many contributions** have been made in gene expression data **in various biomedical domains using a triclustering approach**. **Siska** *et al.* [16] **used OPTricluster**, a triclustering algorithm designed to analyze 3D short time series data, on gene expression data from four yellow fever patients after vaccination **with the purpose of identifying genes that have the same pattern of change in expression across time and experimental conditions**, and **Sari** *et al.* [14] **used TimesVector**, a triclustering algorithm designed to find similarly and differentially expressed patterns in gene-sample-time data, **on gene expression data of medulloblastoma disease**, with the time dimension representing the effect of doxycycline on the medulloblastoma cells. **Both previous works** evaluate the discovered triclusters in terms of overall quality and **present no statistical significance**. In neurodegenerative diseases such as Amyotrophic Lateral Sclerosis, **Soares** *et al.* [17] **proposed a triclustering-based classification**, using triclustering to find disease progression patterns in three-way clinical data and **classifying new patients to determine if they will need NIV in the next 90 days**.

In the **psychiatric domain**, works such as Ikezawa *et al.* [5] Sen *et al.* [15] **show the importance of searching for biomarkers** in hemodynamic responses and MRI's as a way to classify mental disorders such as Schizophrenia, ADHD and autism. But **the use of triclustering** in order to unravel relevant domain knowledge and support predictive tasks **is still scarce**. Rahaman *et al.* [12] **used biclustering**, the search for bi-dimensional patterns, applied **to neuroimaging of patients suffering from Schizophrenia, and more recently** Rahaman *et al.* [13] **proposed a method for finding triclusters in resting-state fMRI data** of both healthy patient's and suffering Schizophrenia with positive results. **Tamminga** *et al.* [18] **present what the Bipolar-Schizophrenia Network for Intermediate Phenotype has learned**, using multiple computational approaches, to better understand psychosis in patients with Schizophrenia, Schizoaffective, and Psychotic Bipolar. **They stress the need for these computational approaches to be able to deal with large datasets and provide insights on complex targets like the brain**.

## References

[1] Andreou, Christina, and Stefan Borgwardt. "Structural and functional imaging markers for susceptibility to psychosis." Molecular psychiatry 25.11 (2020): 2773-2785.

[2] Henriques, Rui, and Sara C. Madeira. "Triclustering algorithms for three-dimensional data analysis: a comprehensive survey." ACM Computing Surveys (CSUR) 51.5 (2018): 1-43.

[3] Henriques, Rui, and Sara C. Madeira. "BSig: evaluating the statistical significance of biclustering solutions." Data Mining and Knowledge Discovery 32.1 (2018): 124-161.

[4] Hu, Zhen, and Raj Bhatnagar. "Discovery of versatile temporal subspace patterns in 3-D datasets." 2011 IEEE 11th International Conference on Data Mining. IEEE, 2011.

[5] Ikezawa, Koji, et al. "Impaired regional hemodynamic response in schizophrenia during multipleprefrontal activation tasks: a two-channel near-infrared spectroscopy study." Schizophrenia research 108.1-3 (2009): 93-103.

[6] Ji, Liping, Kian-Lee Tan, and Anthony KH Tung. "Mining frequent closed cubes in 3D datasets." Proceedings of the 32nd international conference on very large data bases. 2006.

[7] Jung, Inuk, et al. "TimesVector: a vectorized clustering approach to the analysis of time series transcriptome data from multiple phenotypes." Bioinformatics 33.23 (2017): 3827-3835.

[8] Konstantareas, M. Mary, and Terri Hewitt. "Autistic disorder and schizophrenia: diagnostic overlaps." Journal of autism and developmental disorders 31.1 (2001): 19-28.

[9] Lang, Undine E., et al. "Molecular mechanisms of schizophrenia." Cellular Physiology and Biochemistry 20.6 (2007): 687-702.

[10] Pearlson, Godfrey D. "Etiologic, phenomenologic, and endophenotypic overlap of schizophrenia and bipolar disorder." Annual review of clinical psychology 11 (2015): 251-281.

[11] Procyshyn, Ric M., et al. "Medication errors in psychiatry." CNS drugs 24.7 (2010): 595-609.

[12] Rahaman, Md Abdur, et al. "N-BiC: A method for multi-component and symptom biclustering of structural MRI data: Application to schizophrenia." IEEE Transactions on Biomedical Engineering 67.1 (2019): 110-121.

[13] Rahaman, Md Abdur, et al. "A novel method for tri-clustering dynamic functional network connectivity (dFNC) identifies significant schizophrenia effects across multiple states in distinct subgroups of individuals." bioRxiv (2020).

[14] Sari, Ika Marta, et al. "Mining Biological Information from 3D Medulloblastoma Cancerous Gene Expression Data Using TimesVector Triclustering Method." 2020 4th International Conference on Informatics and Computational Sciences (ICICoS). IEEE, 2020.

[15] Sen, Bhaskar, et al. "A general prediction model for the detection of ADHD and Autism using structural and functional MRI." PloS one 13.4 (2018): e0194856.

[16] Siska, Dea, et al. "Triclustering Algorithm for 3D Gene Expression Data Analysis using Order Preserving Triclustering (OPTricluster)." 2020 4th International Conference on Informatics and Computational Sciences (ICICoS). IEEE, 2020.

[17] Soares, Diogo, et al. "Towards Triclustering-Based Classification of Three-Way Clinical Data: A Case Study on Predicting Non-invasive Ventilation in ALS." International Conference on Practical Applications of Computational Biology Bioinformatics. Springer, Cham, 2020.

[18] Tamminga, Carol A., et al. "Biotyping in psychosis: using multiple computational approaches with one data set." Neuropsychopharmacology 46.1 (2021): 143-155.

[19] Tripathi, Adarsh, Sujita Kumar Kar, and Rashmi Shukla. "Cognitive deficits in schizophrenia: understanding the biological correlates and remediation strategies." Clinical Psychopharmacology and Neuroscience 16.1 (2018): 7.