

Automated research platforms

Ali R. Vahdati^{1, *}, Charles N. de Santana^{2, *}, and Carlos J. Melián^{3, *}

1 Department of Anthropology, University of Zurich, Switzerland

2 Institute of Evolutionary Biology and Environmental Studies, University of Zurich, Switzerland

3 Department of Fish Ecology and Evolution, Center for Ecology, Evolution and Biogeochemistry, EAWAG, Kastanienbaum, Switzerland

* These authors contributed equally to this work.

Contents

1	Abstract	3
2	Introduction	4
3	The structure of automated research platforms	5
3.1	Data Integration	5
3.2	Complexity Reduction	5
3.3	Pattern-Process Inference	5
3.4	Validation	5
3.5	Visualization	5
3.6	An example with <i>ROBHOOT</i>	5
3.6.1	Data Integration (<i>DAADI</i>)	6
3.6.2	Complexity Reduction (<i>GOCORE</i>)	6
3.6.3	Pattern-Process Inference (<i>PROPENCE</i>)	6
3.6.4	Validation (<i>VATION</i>)	6
3.6.5	Visualization (<i>VITION</i>)	6
4	Discussion	6
5	Acknowledgments	7
6	Tables	9
7	Figures	10

1 Abstract

High-resolution data coming from many sources is a standard in science, engineering and investment landscapes. Yet, automated inference providing insightful patterns and processes integrating databases with analytical frameworks remain challenging. In this work we review and discuss the challenges for automated workflows to integrate data and pattern-process-based inference accounting for many sources of uncertainty. Our results suggest that automated research platforms will help to do science of science and to take better informed decisions in research, management and investment landscapes.

Keywords: data integration, multilayer networks, approximate Bayesian computation, process-based inference.

2 Introduction

We are in a era of massive data accumulation, integration and pattern detection. Yet, obtaining insights from such an integration accounting for reproducibility, inference and prediction power is at a very incipient stage (Ioannidis, 2005). There are many challenges when aiming to integrate data, inference and prediction. For example, sampling design and experiments (Voelkl et al., 2018), randomizations to achieve solid statistics , and process- or pattern-based model selection and inference just to name a few. They all require many intermediate decisions that make the scientific process difficult to reproduce and challenging to improve. Currently, there are many protocols and platforms automatizing partial steps of the scientific cycle (Table 1), but integrative and automated research platforms with the ability to automatize the whole scientific cycle is currently lacking (Figure 1).

Open automated research platforms might play a leading role in addressing these challenges. Open platforms might allow for maximizing reproducibility of pattern-processes detection along the many paths in the scientific enterprise (Figure 1). These platforms might also help to make the scientific method such a data, pattern, process and insight integration less hard to follow.

The design or research platforms is still in its infancy. Many factors are involved in research platforms: the programming language, the number of packages, its efficiency and the, its functionality in parsing ...

Many questions in science strongly depend on our own bias, lab inertia in the methods and data explored. Therefore, exploring new paths would require new efforts to learn new methods or new collaborations...

Reproducibility and robustness across the different stages of a research platform are two of the desired properties. Reproducibility guarantees the future improvement of the results in future analysis. Many programming languages have tools to facilitate reproducibility (notebooks) and notebooks implementing many languages are already available (jupyter...). Automated research platforms track the explored paths (i.e., the within and between layer interactions) and outlines how close each path is to the empirical patterns accounting for

Sampling design and experiments...

Randomizations to achieve solid statistics...

One of the most discussed challenges nowadays is how to balance pattern and process inference. Many problems might not require a mechanistic understanding to make predictions. Recent examples are AI algorithms playing chess and go. They do not require a theory of mind to win. On the other side, there can be problems that might require a solid mechanistic understanding to make accurate predictions. Examples of these problems can be global warming or astrophysics. Therefore platforms that learn to combine AI and process-based methods

3 The structure of automated research platforms

In this section we outline the steps to develop a research platform. We introduce each of the layers outlined in Figure 1, data collection and integration (DC), complexity reduction, pattern-process inference, validation and visualization. The second part introduces a simple example using the *ROBHOOT* package.

For any given question, there are different methods within each layer that can complete the task. Ideally, one should be able to choose the best method from each layer and connect them to reach insightful patterns and predictions from the data. How many paths are there? Which of these minimize bias? Which topology within and between layers give the best response to our question?

3.1 Data Integration

Data access platforms within and qacross disciplines are highly scattered across the web¹. Researchers have to deal with a highly complex set of intermediate stages and regulations before having access to the raw data. Having “easy” access to the information in a “perfectly informed market” should be simple and efficient, but unfortunately, it is not. Data integration in research platforms is rapidly evolving and there are many platforms that can have access and deliver real time data plots (Table 1).

Data Integration and standarization – Size effects – N labs vs N samplings per lab: Accuracy and uncertainty: How do initial distributions change accuracy and uncertainty? Trade-offs experimental vs big data

3.2 Complexity Reduction

3.3 Pattern-Process Inference

Outline classical variance-covariance matrices, AI algorithms and process-based methods.

3.4 Validation

Describe briefly Bayesian Inference, Approximate Bayesian computation, AIC and BIC model comparison methods. Gibbs sampling – Bayes factors

3.5 Visualization

3.6 An example with *ROBHOOT*

In this section, we illustrate a semi-automated tool combining access to data from both centralized and decentralized platforms and integrating the datasets to infer insights and predictions obtained from analyzing patterns in the datasets (Figure 1). We aim to develop

¹<https://github.com/melian009/Robhoot/blob/master/resources/databases.md>

ROBHOOT in two stages. The first stage will be to develop the free-access platform to have access to integrated databases. The second stage will be to run it automatically to produce insights and pattern inference given specific questions (Figure 1).

3.6.1 Data Integration (*DAADI*)

3.6.2 Complexity Reduction (*GOCORE*)

3.6.3 Pattern-Process Inference (*PROPENCE*)

3.6.4 Validation (*VATION*)

3.6.5 Visualization (*VISION*)

4 Discussion

5 Acknowledgments

References

- Ioannidis, J. P. A., 2005. Why most published research findings are false. PLoS medicine 2:e124. URL <http://www.ncbi.nlm.nih.gov/pubmed/16060722>.
- Voelkl, B., L. Vogt, E. S. Sena, and H. Würbel, 2018. Reproducibility of preclinical animal research improves with heterogeneity of study samples. PLoS Biology 16.

6 Tables

Table 1	
Data platforms	Webpage
Nakamoto Terminal	https://www.nterminal.com
BigQuery	https://cloud.google.com/bigquery/

7 Figures

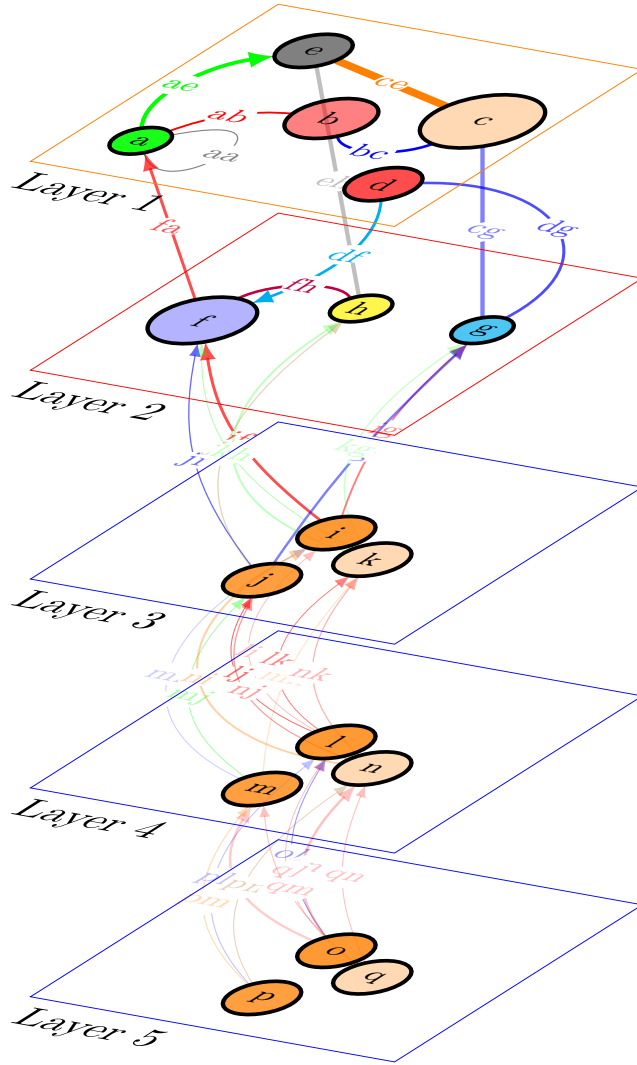


Figure 1: A five layer research platform: Data Integration (Layer 1), Complexity reduction (Layer 2), Pattern-process inference (Layer 3), Validation (Layer 4) and Visualization (Layer 5). Research platforms might play a leading role in accounting for bias and reproducibility in the pattern-process detection enterprise. a) A fully connected 5-layer research platform, and b) A specific path representing the best solution to solve a task.

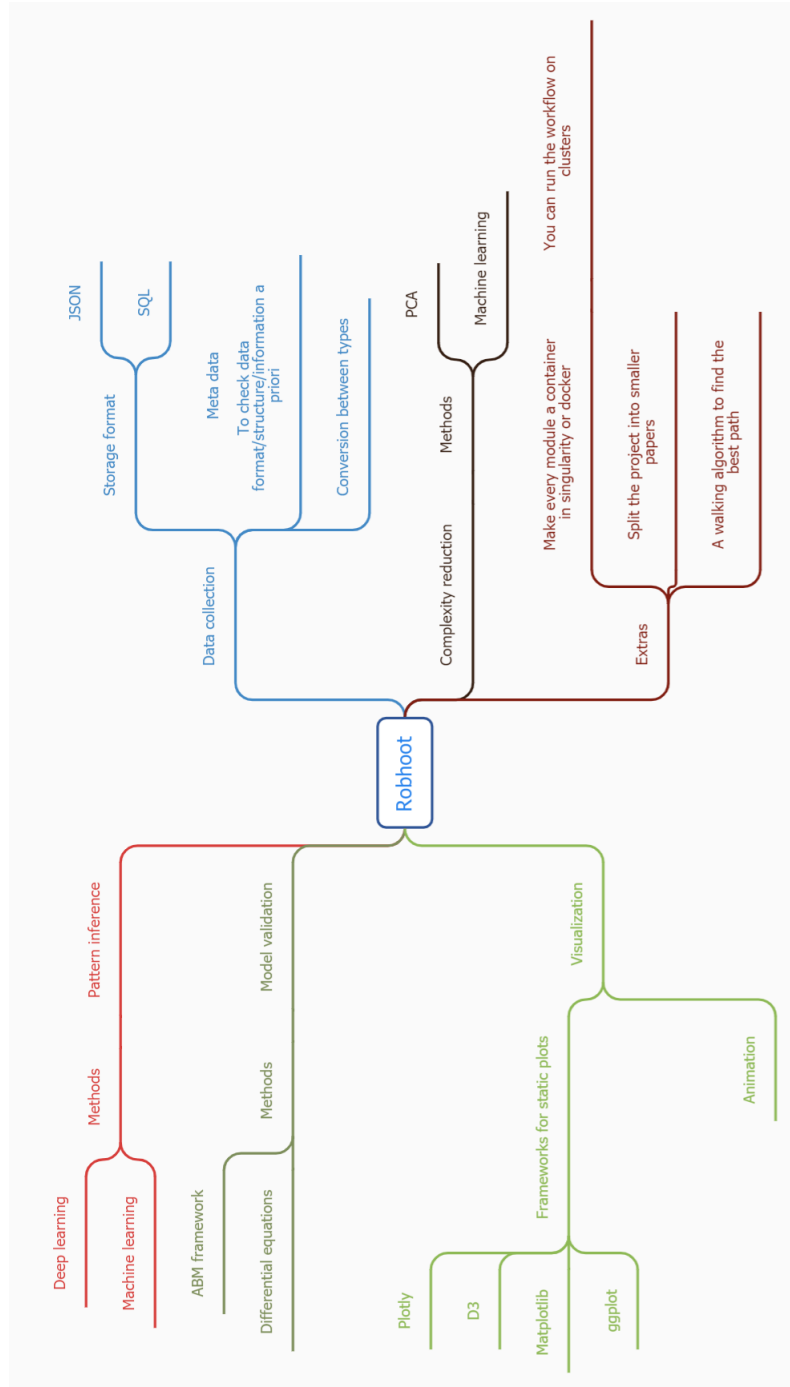


Figure 2: A mind map outlining the different methods to be used within each layer.

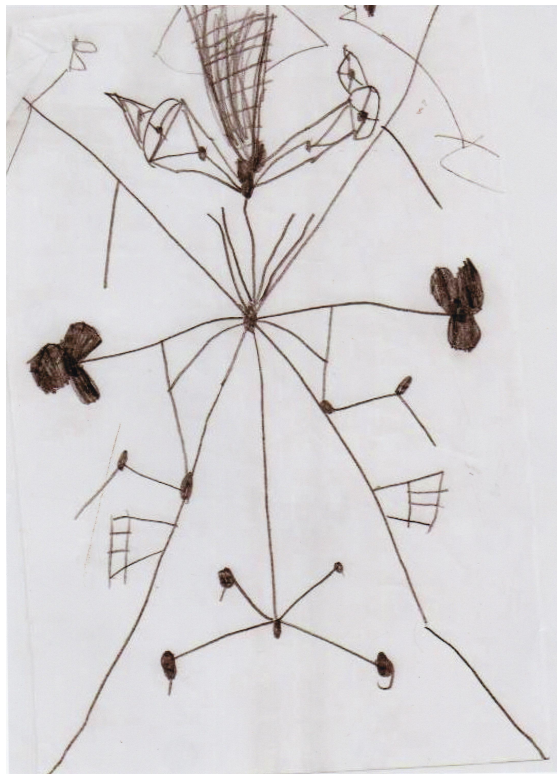


Figure 3: *ROBHOOT*

- An open multilayer platform for data integration, inference and prediction

Index

Open automated research platforms, 4

robust algorithms, 4

robust experiments, 4

robust inference, 4