
Expert-driven approximation of MPE

Subtitle: Reinventing the World

Master's Thesis submitted to the
Faculty of Informatics of the *Università della Svizzera Italiana*
in partial fulfillment of the requirements for the degree of
Master of Science in Informatics
Artificial Intelligence

presented by
Thomas Francesco Tiotto

under the supervision of
Prof. Alessandro Facchini
co-supervised by
Prof. Alessandro Antonucci

September 2019

I certify that except where due acknowledgement has been given, the work presented in this thesis is that of the author alone; the work has not been submitted previously, in whole or in part, to qualify for any other academic award; and the content of the thesis is the result of work which has been carried out since the official commencement date of the approved research program.

Thomas Francesco Tiotto
Lugano, Yesterday September 2019

To my beloved

Someone said ...

Someone

Abstract

This is a very abstract abstract.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et

vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras

ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Contents

Contents	xi
List of Figures	xiii
List of Tables	xv
1 Introduction	1
1.1 Context	1
1.2 Problem and Significance	3
1.3 Response	4
2 Literature review	5
2.1 Explainability	5
2.2 Importance of Explainability	5
2.3 Explainability in Bayesian Networks	5
2.4 Notable Works	5
2.5 Explaining the Most Probable Explanation	5
2.6 A Progressive Explanation of Inference in “Hybrid” Bayesian Networks for Supporting Clinical Decision Making	5
3 Mathematical Background	7
3.1 Probability Theory	7
3.1.1 Probability distributions	7
3.1.2 Random Variables	8
3.1.3 Conditional Probabilities	8
3.1.4 Independence	8
3.1.5 Correlation	9
3.2 Information Theory	9
3.2.1 Entropy	10
3.2.2 Normalised Entropy	10
3.2.3 Mutual Information	11
3.3 Graph Theory	11
3.3.1 Graphs	11
3.3.2 Polytrees	12
3.3.3 D-separation	12
3.4 Bayesian Networks	13

3.4.1	Bayesian Networks Structure Learning	15
3.4.2	Bayesian Networks Updating	16
3.5	Summary	17
4	Methodology	19
4.1	Data set	19
4.1.1	Istituto Cantonale di Patologia	19
4.1.2	Motivation	20
4.1.3	Provided Data Set	20
4.2	Methods	21
4.2.1	Libraries	21
4.2.2	Algorithms	26
4.3	Novel contributions	27
4.3.1	Theory	27
4.3.2	Algorithms	27
4.4	Validation	37
4.4.1	Clinical validation	38
4.4.2	Explainability validation	39
4.5	Summary	39
5	Results	41
6	Conclusions	43
6.1	Future developms	43
	Glossary	45
	Bibliography	47

Figures

3.1	D-Separations in the provided data set (see Sec. 4.1) generated using the <code>plot_model</code> function (see Subsec. ??)	14
4.1	Example output of <code>plot</code> (Schreiber [b])	24
4.2	Example output during the first round of the d-separation-aware variant of <code>dialogue</code> . The variable “mut17q21” is the initial evidence.	29
4.3	Example output during the second round of the d-separation-aware variant of <code>dialogue</code> . “morfologia” is added to the evidence set and this makes a large part of the network redundant.	30
4.4	Example output during the d-separation-aware variant of <code>dialogue</code> . The tuple (“FISH”, “Aampl”) was proposed but the expert refused it and accepted the alternative (“c erbB 2”, “2+”). The main MPE Branch has ID 1 while the “what-if” one has ID 2.	33
4.5	Example output of the Pseudo-MPE from random evidence algorithm.	36

Tables

4.1	Data set variables	21
4.2	Data set distribution before pre-processing	22
4.3	Data set preprocessing steps	23
4.4	Probability quantifiers in natural language	37

Chapter 1

Introduction

1.1 Context

state the general topic and give background of what your reader needs to know to understand the problem outline the current situation evaluate the current situation (advantages/ disadvantages) and identify the gap

In genere mancano referenze

While neural networks and Artificial Intelligence (AI) - as a field - have existed for nearly seventy years, the concept of artificial intelligence dates back to at least Ancient Greece. In ancient times, artificial intelligence embodied in mechanical men was part of the domain of myth; in the twentieth century, of that of science. During this last decade, Artificial Intelligence can't anymore be described by a limited set of terms, as it has materialised out of Man's imagination, broken out of laboratories and has been given lease to act in the world at large.

What do you mean? Why "reductive"?

No sector of our economy has been left untouched by the recent and rapid rise of machine learning that has been enabled by the rediscovery of deep neural networks, the availability of Big Data and cheap parallel computing power. Fields as diverse and as critical as are government, healthcare, finance and bioinformatics have been revolutionised and the possibility has been set for new ones - such as self-driving vehicles - to be born. The ever increasing reliance of our society on ever more complex machine learning-driven algorithms can only make us worry ever more about the ethical dilemmas posed by such a situation. Our society has only very recently been confronted with the dilemma of assigning blame when a driverless car causes the death of a person but this moral problem is only the tip of the iceberg, even when focusing only on the automotive industry. For example, how should a self-driving car behave when confronted with a real-world analogous of the classic Trolley Problem - a situation where each course of action is liable to cause harm? On what basis should a person be denied a mortgage, access to university or a job interview? How can we be sure that there is no bias in the system? How do we even define if the system is behaving morally? Would it currently be feasible for a person that feels they have been harmed by such a decision to appeal it, as prescribed by the recent EU General Data Protection Regulation (GDPR)? As more and more decisions are made in an automated way, with many of them significantly impacting both individuals and society at large, it comes natural to stop and wonder what are the characteristics we would want the systems making these decisions to have.

Explainable AI (xAI) is the sub-field of AI that rests at the intersection between Computer

Science, Social Sciences and Philosophy and whose aim is to define our desiderata of artificially intelligent systems and machine learning algorithms from the point of view of their explainability. The basic idea is that the prerequisite for the evaluation of the ethical and moral implications of a machine's decision is for the system to be "interpretable" or "explainable". Within the xAI community, there is currently no unanimously agreed upon definition of which these desiderata should be or of the best way to implement them in real systems. There is also no common, agreed-upon, definition of what is meant by the phrase "understanding a system": some authors equate it to having a *functional understanding*, void of the low-level details, while others decline it into the concepts of *interpretation* and *explanation*, the former indicating the output of a format that a human user can comprehend and the latter a set of features that have contributed to generating the system's decision.

The difficulties start even in trying to define what interpretability really is. Does it mean to gain the trust the system's user? Of type of user in particular? Does trust stem from some property of the decisions the system makes or from some other inherent characteristic of the machine? A common approach to solving the difficulty in defining interpretability is to try and define it post-hoc by categorising systems into ontologies, based on their perceived interpretability; unfortunately this seems like a circular way of approaching the problem: the classification of system models is being done utilising the same criterion that is trying to be uncovered by doing so. For reference, a commonly used classification is the following:

- **Opaque systems:** these are systems that offer no insight into the mapping between inputs and outputs; all closed-source algorithms fall under this definition;
- **Interpretable systems:** this is the vastest category, as the characteristic of these systems is *transparency* i.e. their inner workings are accessible but the onus of comprehensibility falls completely onto the user. The classical example is that of neural networks where the mapping from inputs to outputs (the *weights*) is inspectable by the user who can, theoretically and depending on her skill, interpret them;
- **Comprehensible systems:** systems falling into this category emit additional symbols together with their outputs with the explicit intent of giving the user the means to interpret and understand the automated decisions; the additional symbols may be visualisations, natural-language text or any other means of demystifying the output. These extra symbols would need to be graded based on the user's expertise, as comprehension is a property that involves both man and machine but materialises on the human side.

Some authors propose to classify systems as *non-interpretable*, *ante-hoc interpretable/transparent* and *post-hoc interpretable*; this roughly corresponds to the ontology presented above.

What I hope can be gleamed from this brief introduction to the field of Explainable Artificial Intelligence, is that many of the problems it aims to tackle are hard *per-se* and may not have a unique optimal solution. This is because these issues are not only engineering problems, but exist at the intersection between man and machine and as such can't be tackled using only the methods of Computer Science. There is no way to satisfactorily investigate the human element of the situation without resorting to the well-established methods of the Social Sciences. There is little hope to know in which direction to proceed without the guiding force that can only come from philosophy, because of its millennia-long tradition in thinking about ethical and high-level issues. It should be clear that when the human - and particularly the ethical - domain are part of the equation, it is impossible *by definition* to find an optimal and unique solution.

references to articles ?

explain better what follows here, since it is important

what does functional mean?

what are low level details

??? explain better. gives examples, help the reader

does not sound well expressed here. also v has trust to do with interpretability? explain intuitions, motivations

why are you speaking about ontologies here not simply of categories etc?

in general. explain what a system, input and output are

References!!!

for instance? which methods? example?

you have to link what you just said with what you are gonna do later. for instance, illustrate what you are saying with examples. use one close to what you are going to do. also, you speak about NN but never about graphical methods, such as BN. you have to speak about them here too.

1.2 Problem and Significance

identify the importance of the proposed research - how does it address the gap? state the research problem/ questions state the research aims and/or research objectives state the hypotheses

again, references!

AI has a trust problem. The bigger problem with AI is not anymore its utility, as that has mostly been solved by deep neural networks, but its capacity to elicit the trust of the users. To be truly useful, an automated system should be able to make itself be trusted in a manner proportional to the criticality of its application. Unfortunately, the explainability and, by extension, the “trustability” of machine learning models are inversely proportional. There are many examples of modern methods - such as boosted trees, random forests, bagged trees, kernelized-SVMs - that show this tendency, but it is best exemplified by *deep neural networks* (DNN). Deep Neural Networks are machine learning models constructed by stacking many layers of artificial neurons, these systems are currently state of the art on a variety of tasks but are among the least easily interpretable systems due to the fact that they represent information in an implicit and distributed manner among their network weights. Some older methods, like decision trees or rule-based methods, are inherently more interpretable due to their simplicity and the fact that they can explicit state their reasoning steps, but are less accurate and flexible than more modern techniques.

ok but come on, there is not only NN in the world

The runaway success obtained by modern Machine Learning in a variety of domains, on a spectrum that goes from engineering to social work, has created the desire to also start applying these methods to mission-critical and traditionally more entrenched fields. A perfect example of a field exhibiting both these characteristics is that of medicine. The first successful artificially intelligent systems date back to the 1970s and '80s and were based on *symbolic methods* integrated with *knowledge-bases*. These systems were by design capable of providing an explanation for their reasoning and were thus accepted by the medical community in an implementation known as *expert systems* that aimed to perform functions similar to those of a human expert. The deficiency of modern AI methods in being able to provide causal links for their reasoning process has held back their acceptance in the field of medicine, regardless of their superior performance and accuracy.

In a high-stakes domain such as the medical one, it would be unthinkable for a doctor to trust the predictions of an AI system a priori; any decision with profound moral implications - such as prescribing or interrupting the treatment of a patient - would have to first be validated by a human. The possibility of carrying out this validation and its quality are dependent on the degree of interpretability of the model that made the decision. Unfortunately, as has been repeated many times, the best performing models are often also the most opaque to inspection.

Explainability is not a necessary condition only for the verification of the system which, as we have just discussed, is a presupposition for it to be applied in mission-critical domains, but also for the extraction of knowledge from data. The amount of information that a machine learning model can process is many orders of magnitude greater than that inspectable by any human; this may let a computer spot new patterns in the data that aren't immediately apparent or are latent given only a moderate amount of samples. Being able to turn this information into

not clear. you were speaking about "validation" before. is this taken as a synonymous of "verification"? if yes, why two words for the same concept? if not, then explain better

why?

new knowledge implies the system having the ability to output human-interpretable symbols that are capable of communicating it in a comprehensible and effective way.

so, what is knowledge? why explainability
essary to create knowledge? you do not ex
this here, imho

There has recently been much research carried out on trying to explain and extract knowledge from deep neural networks together with attempts to marry the connectionist and symbolic approaches to artificial intelligence - a subfield known as *neuro-symbolic computation* while also reconsidering mixed approaches such as *Bayesian Networks*. A Bayesian Network is a graphical and computationally efficient way of representing dependencies between random variables. The graphical component is immediate as in the model each random variable is represented by a node of a Directed Acyclic Graph (DAG), with the edges connecting them standing for their dependencies. The efficiency stems from the fact that the graph structure imposes a factorisation of the joint probability space and thus lets each variable be calculated using only the values of its parents.

again, a lot of nice talk about NN, but then little about BN and thus what you are gonna do

1.3 Response

outline the methodology used - outline the order of information in the thesis - a roadmap - Maximum 2500 words.

The work carried out in this thesis concentrates on explainability in the medical domain and presents both a practical part, with the implementation of a Bayesian network-based system focused on *knowledge-extraction*, as defined at the end of the previous section, and a theoretical one, regarding the definition and validation of desiderata for an artificially intelligent system using the aforementioned system.

The implemented system aims at supporting medical decision making through the instauration of a dialogue with the user/domain expert. To this end, the information implicit in the data is used as basis for a constructive dialogue with the user; this starts with the expert informing the system of which knowledge is certain i.e. a variable's value that has been observed in a specific patient, and continues via a process where the next most probable (*variable, state*) pair is proposed, with the expert having the choice of accepting it or refusing it, if she believes that the variable under examination doesn't adequately explain the accumulated evidence. Each accepted variable is added to the evidence set, as the system gives priority to the domain expert's judgement. The result of the dialogue is an *explanation tree* whose nodes represent (*variable, state*) pairs and are organised into branches, depending on the flow of the dialogue; more specifically, there will always be a *main branch* corresponding to the choices of the user and none or more *alternative branches* whose role is to inform the expert of the possible alternative outcomes to his decisions.

This software system was developed and tested tested in collaboration with *Istituto Cantonale di Patologia*, a medical institute in Locarno, Ticino, Switzerland that specialises in the analysis of tissue samples received from hospitals, clinics and private doctors. Its main activity is to characterise the samples by using *histo-cytopathologic techniques*, with particular focus on the diagnosis of cancer and tumoural diseases in general.

The theoretical part of this thesis aims to understand how an

Chapter 2

Literature review

What is explainability? How is it defined? By whom? When? Why is it important? Notable works in the field

Suggestions when collecting references:

Copy everything in one document (with references !), but do not use it directly when writing the text

Copy it in the document (with references !) and color it , but do not use it directly when writing the text

What is the main point of the sentence / paragraph / article ?

What do I want to convey ?

Read a few references / paragraphs and only then write it down

2.1 Explainability

2.2 Importance of Explainability

2.3 Explainability in Bayesian Networks

2.4 Notable Works

2.5 Explaining the Most Probable Explanation

Butz et al. [2018]

controllare se qualcuno ha lavorato nell'estate
i metodi del paper

2.6 A Progressive Explanation of Inference in “Hybrid” Bayesian Networks for Supporting Clinical Decision Making

Chapter 3

Mathematical Background

Bayesian Networks (BN) are a class of Probabilistic Graphical Models that are used to represent systems under conditions of uncertainty. To give a formal definition we will first need a few basic concepts from probability and graph theory.

3.1 Probability Theory

3.1.1 Probability distributions

Definition 3.1 A probability distribution is a function $\mathbb{P} : \mathcal{S} \rightarrow \mathbb{R}$ with \mathcal{S} a set of events of interest. To be a valid probability distribution \mathbb{P} must satisfy:

- $\mathbb{P}(\sigma) \geq 0 \quad \forall \sigma \in \mathcal{S}$
- $\sum_{\sigma} \mathbb{P}(\sigma) = 1 \quad \forall \sigma \in \mathcal{S}$
- $\alpha, \beta \in \mathcal{S} \wedge \alpha \cap \beta = \emptyset \Rightarrow \mathbb{P}(\alpha \cup \beta) = \mathbb{P}(\alpha) + \mathbb{P}(\beta)$

Each event $\sigma \in \mathcal{S}$ must have a probability $\mathbb{P}(\sigma) \in [0, 1]$ and the sum of all these must equal 1. An event with $\mathbb{P}(\sigma) = 0$ is deemed *impossible* while one with $\mathbb{P}(\sigma) = 1$ is *certain*.

There is some discord regarding how to actually *interpret* the probability of an event. What I believe to be the initially commonly held view is the *frequentist* one, that views the probability of an event as the ratio of times it would occur over a great number of trials. So, for example, saying that obtaining a heads has probability 0.5 when tossing a coin would mean that over repeated throws we would observe heads half the time.

Another, commonly held view is the *Bayesian* (from the 18th century mathematician Thomas Bayes) one in which probabilities are viewed as the *subjective* degree of belief attributable regarding the manifestation of an event. In this interpretation, stating that a coin has 0.5 probability of landing on heads simply means that the person making the claim personally believes that the chances of seeing heads or tails are the same. This is obviously a “softer” definition compared to the frequentist one but it is nonetheless useful in that it lets one characterise certain events that haven’t come about yet or are liable to happen only once or a few times.

Philosophically, Bayesian inference assigns a probability to a hypothesis (a *prior*) while the frequentist method tests a raw hypothesis empirically before assigning it any probability. As

Bayesian inference naturally embraces and deals with uncertainty, it is an enormously useful tool to model and reason about the real, stochastic world we live in.

3.1.2 Random Variables

Definition 3.2 A random variable is a function that associates every outcome in \mathcal{S} with a value.

Random variables are a way of bringing to the fore the attributes of interest of events while dealing with them in a clean, mathematical way. The values that a random variable can take are a function of the events in sample space \mathcal{S} , each of these is assigned a value by the random variable function. I will only be dealing with *categorical random variables* i.e. those whose codomain is a discrete set of values. Every random variable has a probability distribution induced by the cardinality of the subsets of its values; in the case of categorical-valued one, such a distribution is *multinomial*.

If we were to take a Bayesian point of view, we would consider a random variable as simply representing the subjective degree of belief we would have over a set of outcomes we believed possible.

3.1.3 Conditional Probabilities

After having defined the basic notion of probability, we can construe one the basic building blocks of Bayesian Networks: the concept of *conditional probability*

Definition 3.3 The conditional probability, “the probability of event β given event α ” is:

$$\mathbb{P}(\beta \mid \alpha) = \frac{\mathbb{P}(\beta \cap \alpha)}{\mathbb{P}(\alpha)} \quad (3.1)$$

That is, the relative proportion of event β compared to event α ; this intuitively represents the probability of β *knowing* that α has already occurred.

Equation 3.1 can be easily manipulated to obtain another basic element of Bayesian Networks: what is called the *chain rule of conditional probabilities*:

$$\mathbb{P}(\beta \cap \alpha) = \mathbb{P}(\beta \mid \alpha)\mathbb{P}(\alpha) \quad (3.2)$$

This can be generalised to any number of events:

$$\mathbb{P}(\alpha_1 \cap \dots \cap \alpha_n) = \mathbb{P}(\alpha_n \mid \alpha_1 \cap \dots \cap \alpha_{n-1}) \dots \mathbb{P}(\alpha_1 \mid \alpha_2) \mathbb{P}(\alpha_1) \quad (3.3)$$

Intuitively, it means that we can decompose joint probabilities as products of conditional probabilities. As we’ll see, this is how the values in a Bayesian Network are calculated.

3.1.4 Independence

Now, we have just seen in Equation 3.1 that, in general, $\mathbb{P}(\beta \mid \alpha) \neq \mathbb{P}(\beta)$ because $\mathbb{P}(\beta \cap \alpha) \neq \mathbb{P}(\beta)\mathbb{P}(\alpha)$

Definition 3.4 Two events α and β are unconditionally independent $A \perp B$ - or simply independent - when:

$$\mathbb{P}(\beta \mid \alpha) = \mathbb{P}(\beta) \Leftrightarrow \beta \perp \alpha \quad (3.4)$$

This means that knowing that α took place doesn't change our beliefs around β happening. In the real world it is hard, or actually impossible if we consider existence at a fine-enough level to involve Chaos Theory, to find two such perfectly non-interacting events. Thus, a more useful concept is that of *conditional independence* where two previously dependent event become independent when also conditioned on a third one.

Definition 3.5 Two events α and β are conditionally independent ($\beta \perp \alpha \mid \gamma$) when:

$$\mathbb{P}(\beta \mid \alpha \cap \gamma) = \mathbb{P}(\beta \mid \gamma) \Leftrightarrow (\beta \perp \alpha \mid \gamma) \quad (3.5)$$

3.1.5 Correlation

Correlation, as defined by Stolp et al. [2006], is a measure of the degree to which two random variables are linearly dependent. The most used measure of such a dependence is the *Pearson Correlation Coefficient* or *bivariate correlation*.

Definition 3.6 The correlation coefficient ρ of random variables X and Y is given by:

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad (3.6)$$

$$= \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \quad (3.7)$$

$$= \frac{\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} (x - \mu_X)(y - \mu_Y) p_{XY}(x, y)}{\sqrt{\sum_{x \in \mathcal{X}} (x - \mu_X)^2 p_X(x)} \sqrt{\sum_{y \in \mathcal{Y}} (y - \mu_Y)^2 p_Y(y)}} \quad (3.8)$$

with p_{XY} the joint probability mass function of X and Y and p_X and p_Y the marginal distributions of X and Y , respectively.

That is, the bivariate correlation coefficient for random variables X and Y is given by the covariance of X and Y divided by the product of their standard deviations.

The covariance is the *first centred moment* of the *joint distribution* of X and Y while the *standard deviation* is the square root of the *second centred moment* of the marginals.

ρ_{XY} is normalised so its values vary in the interval $[-1, 1]$; the correlation coefficient represents the degree of linear association between the two variables with $\rho_{XY} = -1$ being called *perfect anticorrelation* and $\rho_{XY} = +1$ *perfect correlation*. The two correspond to the cases where the linear equation perfectly describes the relationship between X and Y ; the sign indicates the slope of the regression line describing the relationship i.e. if an increase in one of the two variables corresponds to an increase in the other in the pair, or viceversa. The closer ρ_{XY} tends to 0, the feebler the relationship between X and Y with the case $\rho_{XY} = 0$ indicating that the two variables are *independent*.

3.2 Information Theory

The birth of the field of *information theory* is usually traced back to the seminal paper “A Mathematical Theory of Communication” (Shannon et al. [1949]) where Claude Shannon set the mathematical basis for the quantification of the amount of *information* transmissible over a noisy channel. In his words “The fundamental problem of communication is that of reproducing at one point, either exactly or approximately, a message selected at another point.” The

concepts of field are broad enough to have influenced practically every other scientific discipline and deep enough to have enabled the “digital age”, for example by enabling the creation of ever more complicated coding schemes for the compression, reconstruction and obfuscation of digital data.

3.2.1 Entropy

In classical mechanical statistics, entropy can be seen as a measure of the uncertainty, or randomness, of a physical system. This concept was reapplied by Shannon to measure the amount of randomness in a random variable.

Definition 3.7 *Given a random variable X with probability distribution $\mathbb{P}(X)$, its entropy $H(X)$ is defined as the expected amount of information content carried by X (Schneider [2005]):*

$$H(X) = \mathbb{E}(I(X)) = \mathbb{E}(-\log(\mathbb{P}(X))) = -\sum_{i=1}^n P(x_i) \log_b P(x_i) \quad (3.9)$$

The base b of the logarithm defines the unit of measure. Shannon used $b = 2$ as he was dealing with the transmission of digital, binary-coded data; in this case the unit of measure are *bits*.

The simplest example of how information entropy characterises a random variable X , is in imagining X to model a coin and the task being to predict the probability of the outcome of a throw being heads. If the coin is fair, we will not be any more surprised to see the outcome being heads than tails; the entropy is maximum as there is maximum uncertainty regarding the outcome. However, if the coin is not fair and tails is more probable then we will be more surprised than not to see the outcome being heads. The entropy is sub-maximal because there is less uncertainty regarding the outcome: tails is more probable than heads. If one of the outcomes is impossible, for example if the coin has two heads, then the entropy of the coin is 0 as there is no uncertainty regarding the result of a toss.

3.2.2 Normalised Entropy

Plain entropy is not a good choice when trying to characterise random variables with different cardinalities of their sample space. Let us suppose that the objective is to find the variable with the least “entropic” distribution and we suppose that their values have all been generated by the same process, say Gaussian. Simply calculating their entropies and ordering them according to this criterion will bias the selection process towards the variables with smallest cardinality. This is because we supposed them to be distributed in the same way so there will naturally be less uncertainty when there are fewer possible outcomes. This can easily be understood by imagining the distributions to all be random uniform.

To obviate to this problem we need to *normalise* the entropy so that different-sized variables can be directly compared to each other. To achieve this, we can look at a measure of *normalised entropy* or *efficiency*:

$$\eta(X) = -\sum_{i=1}^n \frac{p(x_i) \log_b(p(x_i))}{\log_b(n)} \quad (3.10)$$

From Eq. 3.10 it can be seen that $\eta(X) \in [0, 1]$; it is thus normalised and comparable among distributions. This ratio expresses the amount of entropy found in the distribution compared

to the maximum possible entropy when using n symbols, corresponding to the uniform distribution:

$$H\left(\underbrace{\frac{1}{n}, \dots, \frac{1}{n}}_n\right) = -\sum_{i=1}^n \frac{1}{n} \log_b\left(\frac{1}{n}\right) = -n \cdot \frac{1}{n} \log_b\left(\frac{1}{n}\right) = -\log_b\left(\frac{1}{n}\right) = \log_b(n) \quad (3.11)$$

3.2.3 Mutual Information

Another way of characterising the interrelatedness of two variables is through the concept of *mutual information*, as defined in Cover and Thomas [2006], that is closely linked to entropy (Eq. 3.9).

Definition 3.8 *The mutual information of two random variables X and Y is given by:*

$$I(X, Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{XY}(x, y) \log\left(\frac{p_{XY}(x, y)}{p_X(x)p_Y(y)}\right) \quad (3.12)$$

NB: In Information Theory, the convention is that $0 \log(0) = 0$.

I_{XY} , intuitively, measures the amount of information that X and Y share that can also be seen as the degree to which one variable is informative of the other. If X and Y are independent then they share no mutual information and knowing one of the two gives no information about the other. This can be immediately understood by rewriting Eq. 3.12 as:

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (3.13)$$

The mutual information $I(X, Y)$ is the reduction in uncertainty of one of the variables given the knowledge of the other. if X and Y are perfectly correlated ($\rho_{XY} = \pm 1$) then they both convey the same amount of information and I_{XY} is equal to the entropy $H(X) = H(Y)$.

3.3 Graph Theory

Many problems in Machine Learning (ML) don't involve classification or prediction of single data points in isolation, but of set of entities that may present a more, or less, complex relation with each other. Most real-world phenomena fit into the latter framework. Graphs are one of the most powerful tools for the modelling of this class of problems, as their structure naturally captures the wide variety of relations that may exist between entities. These range from the atomical structure of a molecule to a social network of friends. In both these examples graphs help in reasoning, visualising and making inferences and predictions.

3.3.1 Graphs

Definition 3.9 *A graph is a tuple*

$$\mathcal{G} = (\mathcal{V}, \mathcal{E}) \quad (3.14)$$

with $\mathcal{V} = \{v_1 \dots v_n\}$ the set of vertices and $\mathcal{E} = \mathcal{V} \times \mathcal{V}$ the set of edges.

For our scopes, we will only be considering the case where every element in \mathcal{E} is a pair either of the form (v_i, v_j) or (v_j, v_i) with $i \neq j$. That is to say that the class of graphs presently of interest for us are those where there can be at most a single directed edge between any node in \mathcal{V} and no self-loops. We are also interested in enforcing that there be no cycles in the graph, i.e. sequences of nodes of the form $v_i \rightarrow v_j \rightarrow \dots \rightarrow v_i$. The resulting graph possessing only directed edges and no cycles is commonly called a *directed acyclic graph*, or DAG for short. This data structure is of paramount importance as it's the fundamental graphical representation used for Bayesian Networks.

3.3.2 Polytrees

We now have all elements to be able to formally define a Bayesian Network. I will also define polytrees and trees because these are a fundamental concept for the work carried out in this thesis.

Definition 3.10 A loop is a trace $v_i, v_j \dots v_i$ of nodes obtained by following edges regardless of their direction

Definition 3.11 A directed graph containing no such loops is called a polytree.

Definition 3.12 A tree is a particular case of polytree where each node has at most one parent.

3.3.3 D-separation

Dependence-separation or *d-separation*, as the name entails, is a concept relating to the conditional dependence between variables. It was first presented by Pearl and Dechter [1988] To define it, we first have to clarify when two sets of nodes X and Y are causally connected. This is so if $Z = \emptyset$ and they are part of one of the following three structures, called *v-structures* in this context:

- $X \rightarrow Z \rightarrow Y$
- $X \leftarrow Z \leftarrow Y$
- $X \leftarrow Z \rightarrow Y$

This means that knowing something about X also tells us something new about Y . X and Y are causally independent if they appear in the following v-structure:

- $X \rightarrow Z \leftarrow Y$

Such a configuration is called a *collider* and it blocks the flow of information from X to Y . If $Z \neq \emptyset$ the cases are reversed so colliders are open and the other three structures are blocked.

Definition 3.13 Given disjoint subsets $X, Y, Z \subset \mathcal{X}$, X and Y are *d-separated* if:

- $Z \neq \emptyset$: no path between X and Y presents a collider
- $Z = \emptyset$: there is a collider on every path between X and Y

The independencies between variables are encoded in the structure of the DAG so every distribution whose BN has the same connections between nodes also has the same independencies, regardless of the values of the variables.

An example using on the real network used in this thesis is shown in Fig. 3.1. Here we can see how the network's topology and the nodes chosen to be in the observed set Z define the resulting separations. In this case $X = \{\text{mut17q21}\}$, $Z = \{\text{differenziazione, pN}\}$ and $Y = V \setminus X \setminus Z$; the resulting set of nodes that are d-separated from X contains $\{\text{c erb 2, FISH, loss 17, situ, pT, M, lateralita}\}$. Basically, we are asking for the set of all nodes in the DAG that are d-separated from X , given evidence Z . The reason for this can easily be given by enumerating all paths through the v-structures in the network and applying the definitions for causal connections given above:

- $\text{mut17q21} \rightarrow \text{eta arrotondata}$: are directly connected so the cannot be independent
- $\text{mut17q21} \rightarrow \text{recettori estrogeni}$: like above
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{morfologia}$: $\text{recettori estrogeni}$ is not in Z so the path is open and mut17q21 is d-connected to morfologia
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{recettori progesterinici}$: like above
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{ki67}$: as above
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione}$: as above
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{c erbB 2}$: differenziazione is in Z so this path is closed, as there are no other paths between mut17q21 and c erbB 2 , the two nodes are d-separated
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{FISH}$: the only path to mut17q21 is through c erbB 2 , that is already d-separated from mut17q21 , so the two nodes are d-separated
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{FISH} \rightarrow \text{loss 17}$: as above
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{pN}$: as $\text{recettori estrogeni}$, differenziazione and pN are arranged in a *collider* and differenziazione is in Z , then the path is open
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{pN} \rightarrow \text{situ}$: pN is in Z and differenziazione , pN , situ are arranged as a *fork*, thus the path is closed
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{pN} \rightarrow \text{pT}$: as above, M is also d-separated from mut17q21 so there are no open paths between mut17q21 and pN
- $\text{mut17q21} \rightarrow \text{recettori estrogeni} \rightarrow \text{differenziazione} \rightarrow \text{pN} \rightarrow \text{M}$: differenziazione , pN , M are arranged as a *chain* and pN is in Z , the path is closed
- lateralita : lateralita is not connected to any other node and is thus d-separated from all others a priori

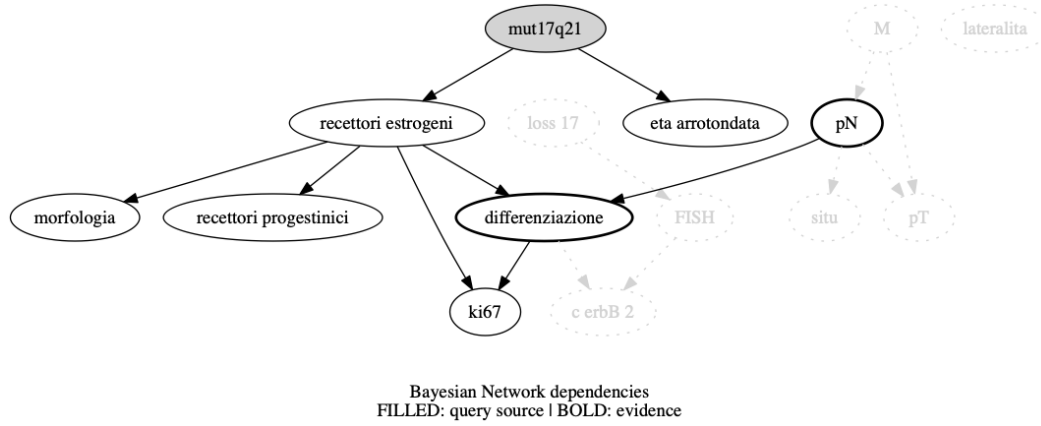


Figure 3.1. D-Separations in the provided data set (see Sec. 4.1) generated using the `plot_model` function (see Subsec. ??)

3.4 Bayesian Networks

Definition 3.14 A Bayesian Network (BN) is a probabilistic graphical model represented by a DAG where each vertex corresponds to a random variable X_i and the edges model the dependencies among these.

Such a model is basically a way of compactly representing an explicit joint distribution $\mathbb{P}(X_1 \cap \dots \cap X_n) = \mathbb{P}(X_1) \dots \mathbb{P}(X_n)$, that is factorised into $\mathbb{P}(X_n | X_1 \cap \dots \cap X_{n-1}) \dots \mathbb{P}(X_2 | X_1) \mathbb{P}(X_1)$. The way this compactness is achieved is in exploiting the independencies that exist among the random variables:

$$\forall X_i : (X_i \perp \neg Desc(X_i) | Pa(X_i)) \quad (3.15)$$

with $Pa(X_i)$ the set of nodes that are parents of X_i and $Desc(X_i)$ the nodes that are not *descendants* of X_i . That is to say, every random variable X_i , given its parent nodes, is independent of all other nodes in the Bayesian Network that are not descended from it. Also, a BN gives the flexibility to drop the many weak dependencies that are bound to exist between variables thus leading to an even simpler model. A full probability table for a joint distribution of random variables obscures the independencies and requires an exponential number of entries for the representation. A Bayesian Network on the other hand can represent the same distribution using only a linear number of parameters. The way that Bayesian Networks can be used to reduce the storage requirements for uncertain information is by taking advantage of the conditional independencies embedded in the underlying distribution being modelled. The power of BNs comes from the additional information encoded in their structure and this was first explicitly described in its entirety by Pearl and Dechter [1988] who defined the concept of dependence separation (see Subsec. 3.3.3) and applied it to Bayesian Networks.

One nice characteristic of BNs is that they very naturally model the type of mixed causal and stochastic processes that we find in all of Nature. Imagine we want to represent the process modelled by joint distribution $\mathbb{P}(B, A) = \mathbb{P}(B) \mathbb{P}(A)$; using the chain rule for conditional probabilities (Eq. 3.3) we can write this as $\mathbb{P}(B | A) \mathbb{P}(A)$. A BN modelling this process would be composed of two nodes A and B with an edge from the former to the latter $A \rightarrow B$, A is called the “parent” of B . Each of these two nodes would have its own probability table, with $\mathbb{P}(A)$

representing the *prior* distribution over A and $\mathbb{P}(B | A)$ the *conditional probability distribution* of B given A .

We can now see why these types of models are named *Bayesian Networks*: the inference process is based in a given prior distribution/belief and evolves through a parent \rightarrow child relationship to constantly yield an updated *posterior* belief. The BN DAG encodes a generative sampling where each variable's value is determined stochastically by Nature, based on the value of its parents. This process is also highly compatible with our view of causality and this is one of the reason that makes BNs highly interpretable. The prior $\mathbb{P}(A)$ can be seen as the result of some stochastic process caused by a series of latent (unmodelled) variables while the posterior $\mathbb{P}(B | A)$ is stochastically, causally determined by A . As I have mentioned in the previous paragraphs, there are probably no truly “prior” distributions in the Universe, at the modelling scale we are usually interested in. Only on arriving on the quantum particle level may we find “pure” stochastic, uncaused processes due to quantum collapse.

A good example of how BNs are well compatible with our notion of causality may be to imagine A as the random variable modelling the predisposition to having a certain disease and B to actually developing the symptoms for it. *First*, genetic and epigenetic factors such as the environment stochastically contributed to having the predisposition and *then* the development of the symptoms was stochastically determined by the degree of predisposition. Adding an extra time dimension certainly helps us in dealing with this class of models.

3.4.1 Bayesian Networks Structure Learning

In many probabilistic models initialisation is fast but then fitting the data is slow (ex. k-means). For Bayesian Networks the converse is true: fitting is fast as only sums of the counts in the data are needed but identifying the correct graph structure can take super-exponential time. Learning the Bayesian Network structure from data is commonly known as the Bayesian Network Structure Learning (BNSL) problem. The methods to solve this problem can be roughly categorised into one of three types.

Search and Score

This is the most naïve method as it does a brute force search over all the possible graph structure space - i.e. all DAGs with the same number of variables as the input data - and scores all these depending on some cost function. This process is super-exponential but though the use of dynamic programming and heuristic search algorithms it can become sub-exponential. Nonetheless, solving the exact BNSL is only feasible up to 30 variables.

Constraint Learning

Methods of this type calculate some measure of correlation to identify the presence and direction of edges between nodes. A typical test is to iterate over all triplets while testing for conditional independencies. Thanks to the d-separation properties outlined in Sub-Section. ??, this test is able to identify the correct edges. The algorithm is quadratic in time in the number of vertices.

Approximations

Several heuristical approaches have been developed to be able to find good network structures in an efficient manner. Examples of these are:

- Chow-Liu, that builds a tree approximation of the probability distribution
- Greedy Hill-Climbing, that adds/removes/flips an edge at a time
- Optimal Reinsertion, that iteratively calculates the optimal *Markovblanket* (the subset of all nodes that are sufficient to determine the value of another subset) of an ever-smaller subset of nodes

3.4.2 Bayesian Networks Updating

All the types of inference presented are instances of *diagnostic reasoning*, also known as *abductive reasoning*. This type of explanation can either be modelled as a conditional probability or a MAP query and is of fundamental importance in many important problems of machine learning including medical diagnosis, that is of particular interest to us.

Conditional Probability Query

The *updating* problem is the process of updating the probabilities of nodes in the BN based on the observation of the values of other vertices. This process of conditioning on observed information is also called *data propagation*.

The following algorithm was described by Normand and Tritchler [1992] and applies to our case where the random variables follow a multinomial distribution. What we want, is to calculate the conditioned probability $\mathbb{P}(B \mid D)$ i.e. the updated probability of node B based on observed evidence E .

Definition 3.15 *The conditional probability query for variable B given evidence E is:*

$$\mathbb{P}(B \mid E) = \alpha \pi(B) \lambda(B) \quad (3.16)$$

with $\pi(B) \lambda(B)$ analogous to the prior and likelihood of B , respectively.

The likelihood of B depends only on the weighted likelihoods of its children C_1, \dots, C_k :

$$\lambda(B) = \prod_l \lambda_{C_l}(B) \quad (3.17)$$

$$\lambda_{C_l}(B) = \sum_{C_l} \lambda(C_l) P(C_l \mid B) \quad (3.18)$$

and its prior similarly depends only on the information received from its parents A :

$$\pi(B) = \sum_A P(B \mid A) \pi_B(A) \quad (3.19)$$

$$\pi_B(A) = \alpha \pi(A) \prod_{S_B} \lambda_{S_B}(A) \quad (3.20)$$

The information is propagated down if any variable observed is above B while up if any variable observed lives in the tree rooted in B . Initially all leaf nodes' likelihoods are set at 1 and the priors of root nodes are assumed to be observable.

Maximum a Posteriori Query

Another common type of question we might ask a BN is the following: “given evidence E which is the most likely assignment of a subset of variables Y ?”. This is known as *Maximum a posteriori (MAP)* inference and is a much harder problem than a conditional probability query. We are trying to solve the optimisation problem.

Definition 3.16 *As defined by Koller et al. [2009a]. Given evidence/observed variables $E = e$, $E \subseteq \mathcal{X}$ and sets $Y \subseteq \mathcal{X} - E$ and $Z = \mathcal{X} - E - Y$, with \mathcal{X} the set of all variables in the BN, the MAP query for Y is the assignment of values $Y = y$ that has maximum probability:*

$$\text{MAP}(Y = y \mid E = e) = \underset{y}{\operatorname{argmax}} \sum_z \mathbb{P}(Y = y, Z = z \mid E = e) \quad (3.21)$$

The MAP problem is hard to solve efficiently; that is it is part of the *NP-hard* complexity class, as proved by Shimony [1994]. Calculating it in a brute-force way would mean enumerating all the possible variable-value tuples and computing their joint probabilities; as these are exponential in the number of variables, the problem is evidently untractable. Moreover, this is true even in a Bayesian Network. Such a model may possess a linear number of parameters but the underlying distribution is still exponential. Explicitly calculating the MAP defeats the very purpose of the BN, that is computational efficiency. For this reason, there exist a host of approaches to optimising MAP: elimination algorithms, gradient methods, simulated annealing and other stochastic local searches, belief propagation and integer linear programming.

A very important thing to note is that the greedy assignment where each variable picks its most likely value can be very different from the most likely joint assignment of all variables.

Most Probable Explanation Query

A special case of MAP is the *Most probable explanation (MPE)* that,

Definition 3.17 *As defined by Koller et al. [2009a]. Given evidence/observed variables $E = e$, $E \subseteq \mathcal{X}$ and $W = \mathcal{X} - E$, the MPE query for W is the assignment of values $W = w$ that has maximum probability:*

$$\text{MPE}(W = w \mid E = e) = \underset{w}{\operatorname{argmax}} \mathbb{P}(W = w \mid E = e) \quad (3.22)$$

This is an easier problem than MAP, as can be seen by comparing Eq. 3.21 with Eq. 3.22; MAP presents both a summation and a maximisation and as such is part conditional probability query, part MPE query. All algorithms for the computation of MAP obviously apply to MPE too, but there exist efficient approximate algorithms for MPE that do not generalise to MAP such as Loopy Belief Propagation [Pearl and Dechter [1988]] and Stochastic Local Search [Kask and Dechter [1999]].

3.5 Summary

Chapter 4

Methodology

The inspiration for the work carried out in this thesis was the paper “Explaining the Most Probable Explanation” by Butz et al. [2018], that has been presented in detail in Sec. 2.5. This paper proposed a system that would build a Bayesian Network modelling a medical data set and, through the interaction with a medical expert, distill an explanation tree. This tree, deemed to represent the solution to the MPE query, could then be used to generate a natural language explanation that the authors claim would lead to the extraction of extra knowledge from the original data set.

The driving hypothesis of the paper was that Bayesian Networks and the solution to the MPE problem would be a powerful tool in helping medical experts gain insights into data. Unfortunately, the paper did not provide any indication that a such a system had ever been built and any validation of the method was left for future work. As of the finalisation of this thesis (July 14, 2019), there has been no work done in substantiating the conclusions by Butz et al. [2018]. As discussed in Chap. 1 and Chap. 2 there is an ever greater need for Machine Learning models and systems to be explainable, especially in mission-critical domains as is healthcare. Current machine learning systems are for the most part opaque and there is confusion regarding even what would constitute a good explanation of their working.

For these reason, I believe that building a proof-of-concept system whose logic was inspired by the method presented in the aforementioned paper and validating it with real medical experts would be an important step forwards in the direction of answering the following questions: Can the method of the paper be corroborated? Are Bayesian Networks a good ML model to bootstrap an explanation from? How good an explanation does the proposed method give, as validated by a domain expert? What improvements are there to be made?

controllare altri autori

migliorare dopo aver fatto literature review
perche avro piu idee

4.1 Data set

As anticipated in Chap. 1 the work carried out in this thesis had a certain degree of collaboration with a third party, Istituto Cantonale di Patologia.

4.1.1 Istituto Cantonale di Patologia

Istituto Cantonale di Patologia (ICP) is an institute based in Locarno that is specialised in the histological analysis of tissue samples received from private patients, clinics and hospitals, mainly

in support of cancer diagnosis. Its Laboratory of Medical Diagnostics supports pathologists in the diagnosis of neoplastic diseases through the application of cytogenetic techniques; that is, the focus is on understanding how chromosomes relate to cell behaviour, particularly during mitosis and meiosis. One of the techniques used is Fluorescence in Situ Hybridization (FISH) that is able to localise the presence or absence of specific DNA sequences in chromosomes. These tests are aimed at identifying the precise profile of the cancer cells and thus inform the clinician on the best treatment for the specific patient.

In addition to its clinical support activities, the Istituto also carries out scientific research aimed at better understanding certain types of cancers at a basic level. In the last ten years, the ICP has published more than 200 peer-reviewed papers and more than 100 works in non-peer reviewed journals and is active at a national and international level.

4.1.2 Motivation

My first contact with the ICP was during a meeting with Dr. Vittoria Martin (Martin [2012]), molecular cytogeneticist, in date 28/01/2019. The institute had expressed interest in bringing machine learning into their workflow in order to both augment their profiling capabilities for patients and to be able to extract new knowledge from their existing data. This knowledge-extraction may lead towards the confirmation of current scientific theories or may be the first step towards the formulation of novel ones.

magari togliere CV e mettere descrizione

chiedere a Vittoria cosa si aspettavano/aspettano

My interest in collaborating with the Istituto stemmed from the desire to apply the methods described in Sec. 4.3 to a real-world case. Being the theoretical work being carried out in this thesis an expert-driven MPE approximation, collaboration with the institute has also provided the opportunity to implement a proof of concept using real histological data. The doctors and researchers of the Istituto have been able to validate, from an Explainable AI and clinical relevance points of view, the model software that I have developed. That is to say, they have validated the capacity of the developed software both in its capacity to support clinical decision making and to surface clarifying explanations of the data set and in its adherence to established medical literature.

da confermare con Vittoria piu avanti

4.1.3 Provided Data Set

The data set I was provided with was created by *Registro Tumori Ticino* (Locarno, Ticino) in order to highlight possible new relations between clinical, histopathological and molecular features, as well as to potentially discover novel biomarkers for the progression of the disease. It consists of the histological records over 38 recorded variables of 3218 breast cancer patients who have been diagnosed between the years 2005 and 2014 within the Ticino canton of Switzerland. The data set had been pre-processed by collaborators of IDSIA in agreement with the ICT with 13 of the variables being dropped, because not considered relevant. In particular, all variables relating to the patient post-treatment were discarded as well as those recording the diagnosis date. In Tab. 4.1 is a description of the measured variables, together with their clinical meaning. The value distribution of the data set is shown in Tab. 4.2.

The indications from Dr. Martin on how to further preprocess the data are shown in Tab. 4.3. Note that some variable names were simplified.

Table 4.1. Data set variables

Variable	Clinical meaning
Codice globale	Unique patient identifier
mut17q21	Mutation of chromosome 17
loss 17	Loss of chromosome 17
età arrotondata	The age of the patient at diagnosis
Lateralità	The affected breast
Situ SUBGROUP MZ	The primary site code of the tumour
Morfologia SUBGROUP MZ	The morphology classification of the tumour
pT SUBGROUP MZ	Primary tumour in the TNM classification for breast cancer
pN SUBGROUP MZ	Pathologic in the TNM classification for breast cancer
M 8.2.96	Distant metastasis in the TNM classification for breast cancer
Differenziazione	Tumour grade
Recettori estrogeni percento 1.1.2003	Expression of estrogen receptors
Recettori progestinici percento 1.1.2003	Expression of progestin receptors
c erbB 2 cod percento 1.1.2003	ErbB2 marker expression
Ki67 cod percento	Tumoural proliferation index
FISHRatio	FISH analysis result

4.2 Methods

4.2.1 Libraries

Pomegranate

pomegranate (Schreiber [a]) is an open-source probabilistic models package for python. Its core philosophy is that every probabilistic model, from Hidden Markov to Bayesian Network, can be seen as a probability distribution and, as such, can be flexibly composed into hierarchical mixture models (Schreiber [2017]). The package implements:

- Probability Distributions
- General Mixture Models
- Hidden Markov Models
- Bayes Classifiers and Naïve Bayes
- Markov Chains
- **Bayesian Networks**
- Factor Graphs

Table 4.2. Data set distribution before pre-processing
















Variable	Unique values	Distribution
mut17q21	2	
loss 17	3	
eta arrotondata	74	
ateralità	3	
situ	5	
morfologia	5	
pT	23	
pN	6	
M	3	
differenziazione	5	
recettori estrogeni	40	
recettori progestinici	40	
c erbB 2	4	
Ki67	52	
FISH	5	

Table 4.3. Data set preprocessing steps

Variable	Action
Codice globale	Remove variable
mut17q21	Remove blanks
loss 17	Remove blanks
eta arrotondata	Bin into “< 40”, “40 – 50”, “≥ 50”
lateralita	Remove blanks and “sconosciuta”
situ	Remove blanks
morfologia	Remove blanks and “unuseful” if performance on classification is subpar
pT	Remove blanks and “unuseful”
pN	Remove blanks and bin into “0” and “≠ 0”
M	Remove blanks
differenziazione	Remove blanks and “Sconosciuto o non applicabile”
recettori estrogeni	Remove blanks and bin into “negativo” if ≤ 10 , “debolmente positivo” if ≤ 50 , “fortemente positivo” if > 50
recettori progestinici	Remove blanks and bin into “negativo” if ≤ 10 , “debolmente positivo” if ≤ 50 , “fortemente positivo” if > 50
c erbB 2	Remove blanks
ki67	Remove blanks and bin into “<14”, “14-20”, “20-30”, “>30”
FISH	Remove blanks

This package was chosen among others for its good implementation of Bayesian Networks, its clear API and its performance. The package is written in cython and natively supports multi-core parallelism and out-of-core learning. Network structure learning from data, described in 3.4.1, appears to be particularly efficient, thanks to the implementation of prior knowledge into the graph selection process as described by Schreiber and Noble [2017]. The claim of this novel selection process is that it possesses the speed of a heuristic approach while yielding a far better quality estimate.

Pomegranate currently only supports Discrete Bayesian Networks so the random variable of each node must have a categorical distribution.

Structure learning from data is achieved using the `from_samples` method of the `BayesianNetwork` class, with the default algorithm being the novel one described by Schreiber and Noble [2017]. The *probability* of a sample is calculated using the *probability* function of an object of `BayesianNetwork` type; the `predict_proba` function is used to return the probability of each variable in the model given some evidence. *Predictions* (described in detail in Subsec. 3.4.2) are run by passing to the `predict` function of an object a matrix with `None` as placeholders for missing values. *Fitting* is done through the `fit` function that uses MLE estimates to update each node's distribution in the model based on the input data.

A `BayesianNetwork` object can also be displayed graphically by calling its `plot` function. The output is a DOT file that is generated using the PyGraphviz package (PyGraphviz developer team), that is a python interface to the famous Graphviz (gra) graph visualisation software. An example of such an output is shown in Fig. 4.1.

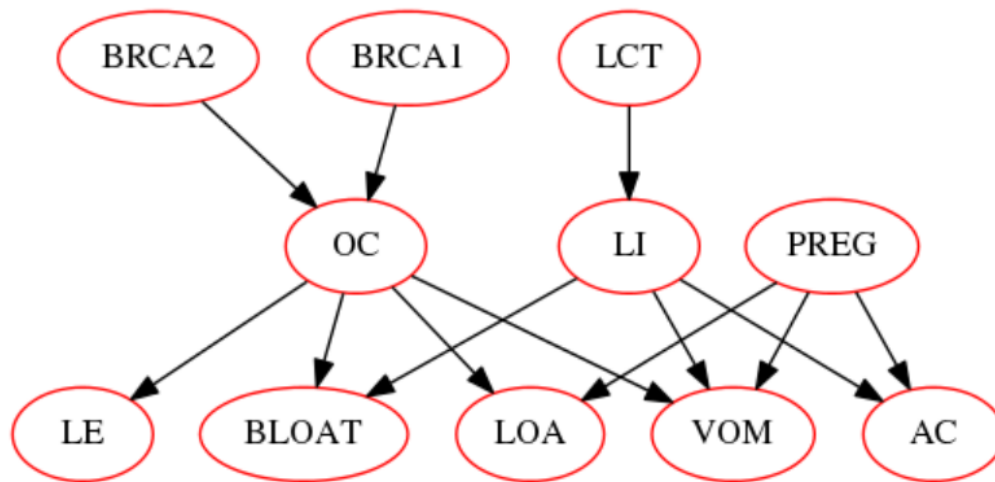


Figure 4.1. Example output of `plot` (Schreiber [b])

Pgmpy

Pgmpy (pgm) is, like pomegranate, another recent probabilistic graphical model package for Python. Unlike pomegranate, it natively implements various exact and approximate inference algorithms, like Variable Elimination, Belief Propagation and Max-Product Linear Programming.

The reason that I elected to use two different probabilistic graphical model libraries is because there is still no Python package that offers all the needed functionality. Pomegranate

implements a novel structure learning algorithm, that I wanted to use, but is severely lacking in functionality in many other areas. Pgmpy, on the other hand, has a very good API as regards inference.

Gurobi

Gurobi (gur) is a closed-source mathematical programming solver for Linear Programming, Quadratic Programming and Mixed-Integer Programming optimisation problems. It claims to be the fastest solver available for these classes of problems. Gurobi offers object-oriented and matrix-oriented interfaces to, among others, Python, MATLAB and Excel.

Pandas

pandas (pan) is an extremely widely-used open-source python library that provides data structures and methods to aid in data analysis. The package excels in the manipulation of tabular data in the form of `DataFrame`, that is the analogous of R's `data.frame`. A `DataFrame` can be seen as a “general 2D, size-mutable structure with potentially heterogeneously-typed columns”. The syntax for slicing is very close to R's as are many other functionalities; this is because one of Pandas' explicit goals is to offer all of CRAN's functionalities and be easily approachable by anyone already knowing the other language.

Pandas the default choice for this thesis' implementation because it is the *de facto* standard in data analysis applications. Its flexibility in reading Excel spreadsheets (the format the data set the project was built on, see Sec. 4.1) and in then manipulating the data confirmed that this was a good choice. Note that to read files in the Excel formats the additional `xlrd` package is needed.

Scikit-learn

scikit-learn (sci) aims at providing a unified API for basic Machine Learning; it does not include advanced paradigms such as Reinforcement Learning or graphical models for structured learning. The latter omission was the reason that lead me to select pomegranate as the basis for the implementation of a Bayesian Network. What is included are a stack supervised and unsupervised ML tools to prepare data sets, define machine learning models ranging from spectral analysis-based to ensemble methods to clustering and multiple evaluation and model selection utilities.

NumPy

NumPy (num) is another *de facto* standard package when doing scientific computing with python. Most scientific packages (including pandas, scikit-learn and TensorFlow) depend on NumPy for low-level operations; this is because NumPy is contains fast implementation of n-dimensional array objects together with powerful manipulation functions. In addition to this, NumPy implements linear algebra operations, Fourier Transform and random number generation. The closest parallel to NumPy - as R was for pandas, is MATLAB.

Networkx

NetworkX (net) is another widely-used package that is specialised in the creation and manipulation of graph-structured data. The main use for this package was in building the MPEGraph data structure that I used to build a polytree representing the dialogue with the expert.

4.2.2 Algorithms

Model construction

The data was given in .xlsx format and was imported using panda's `read_excel` function that returned a `DataFrame` object. The so imported data was then preprocessed by dropping unwanted records and binning the remaining ones as described in Tab. 4.3 . The actual BN representation is learned at runtime by calling the `from_samples` method of pomegranate's `BayesianNetwork` to solve the structure learning problem (defined in Subsec. 3.4.1).

d-separation

A naïve implementation according to the definition (presented in Subsec. 3.3.3) to check for d-separation between node X and Y would have a complexity in the order of the number of trails between X and Y ; this leads to an exponential in the size of the graph running time. Luckily, Koller et al. [2009b] present what is a linear time algorithm to solve the problem.

The `reachable` procedure takes as input the DAG representing the Bayesian Network \mathcal{G} , a source variable X and a set of observed variables Z ; on exit it returns the set of variables R that are reachable from X . The procedure runs in two phases, traversing the graph twice: first bottom-up from leaves to roots, then viceversa. During the first run, the algorithm finds all nodes A that are ancestors of the evidence set Z . During the second, the procedure distinguishes the direction it visits each node in order to determine if it is traversable or not. Any node Y that is not in the evidence set is marked as reachable; if it is being visited in direction “up” it can be traversed as the v-structure is a chain. All the parents of Y are marked to be visited in the “up” direction (i.e. from below) and the converse is done for Y 's children. If Y is being visited in the “down” direction its children are again added to be visited in the “down” direction, because Y is traversable. Additionally, if Y happened to be in the set A , found in the first step, then Y 's parents are marked to be visited in the “up” direction then the collider is active and Y can be traversed (a collider is open iff. the central node or any of its descendants are observed).

The full procedure can be found in Koller et al. [2009b]; my implementation shown in Alg. 1 follows this pseudocode very closely but the procedure `d-separated`, instead of finding all nodes R that are d-connected to the input X , tests if a given target Y is d-separated from X or not. This gives some extra flexibility in how the function can be used. To find the set S of all nodes d-separated from X I simply iterate the `d-separated` test over all nodes V in the graph representing the BN.

MPE

thresholded dialogue aggiungere

quanti dettagli implementativi devo mettere
ch'Al ho un sacco di funzioni di visualizzaz
etc

descrivere come ottenere MPE con Gurobi:
porter, UAI ecc...

Algorithm 1 d-separation algorithm

```

1: separated_list =  $\emptyset$ 
2: for target  $Y \in V$  do
3:   append  $d\text{-separated}(X, Y, E)$  to separated_list            $\triangleright$  will return true or false
4: end for

```

Other Machine Learning methods

4.3 Novel contributions

tengo qui o sposto nel cap 4?

4.3.1 Theory

Selection based on entropy

la discussione con l'esperto ci da' belief revision dei dati

counterfactual explanations

magari togliere theory e mettere tutto assieme

utilizzo entropia

4.3.2 Algorithms

An important part of my work was developing the algorithms needed to adapt the ideas presented in the paper “Explaining the Most Probable Explanation” by Butz et al. [2018] and “A Progressive Explanation of Inference in ‘Hybrid’ Bayesian” by ?. From the former, the construction of the probability tree through a constructive dialogue with the domain expert, the building of counterfactual explanation branches, the automatic generation of the most probable probability tree from initial evidence. From the latter, the generation of an “Inverse explanation”. Finally, a simple procedure to output a natural language explanation was developed.

“Pseudo-MPE”

The so-called “Pseudo-MPE” algorithm is inherently wrapped up with the concept of *dialogue* and is central to the explanatory powers of the system being developed in this thesis. The algorithm was developed as a way of implementing the “MPE branch” of the “Argumentative Probability Tree” hypothesised by Butz et al. [2018]. It was termed “Pseudo-MPE” because there are no guarantees of it returning the MPE solution (see Subsec. 3.4.2), as noted by Koller et al. [2009a] in their definition of the MAP problem.

dove critico il fatto che non calcolano veramente l'MPE come dicevano? qui o in cap.4 on ca

At a lower level of detail, the algorithm may be broken into:

- a dialogical part, that interfaces with the expert user through the use of natural language, menus and visualisations
- the part responsible for constructing the “MPE branch”

The former process was informed and shaped by the results obtained by the methods described in Subsec. 4.4.2 and, as such, presented substantial elements of novelty.

The latter process is, at its core, a greedy procedure that aims at selecting the “best” next $(state, value)$ tuple at each step, based on some measure of optimality and on the variables already in the evidence set. In the actual implemented system the two parts are intertwined, given their close inter-dependence.

The dialogue procedure starts by asking the user to select a subset of variables and their relative values to add as initial evidence. This initial evidence is used to radicate the MPE Branch. It should be noted than in the description given by Butz et al. [2018], the Argumentative Probability Tree is a real tree as each node is guaranteed to have at most one parent. My application, on the other hand, constructs an *Argumentative Probability Polytree* (see Subsec. 3.3.2) because, as will better be described in Chap. 5, it was seen early on that the users much preferred to be able to start from a set of initial evidences and not be limited to a single one. The algorithm then proceeds to call the `next_most_probable_states` subroutine that is tasked with returning an ordered list of $(state, value)$ pairs. It does this by calculating the posterior distribution given evidence of all the states not already in the evidence, then calculating the efficiency (see 3.2.2) and the maximally probable symbol of each state’s distribution and finally returning the $(state, value)$ tuples ordered according to their normalised entropy (most efficient/least entropic at the head). The $(state, value)$ pair at the head of the list is proposed to the user who has the faculty to accept the system’s evaluation or refuse it. If the user accepts, the state is added to the evidence set and to the MPE Branch under construction. Thus, the evidence set’s cardinality increases by one each time a user accepts a proposal. The updated evidence will be used to calculate the new list of $(state, value)$ pairs at the following round. If the expert chooses to refuse, then she is iteratively presented with the remaining $(state, value)$ items, in order of decreasing efficiency. Once she accepts one of the explanations given by the system, the `generate_alternative_branch` subroutine is called to automatically generate a maximally probable MPE Branch, radicated in the newest $(state, value)$ node of the MPE Branch. The proposal loop for alternative states runs until there are increasingly less probable elements in the list and exits with a partial solution if the user refuses all of them at a given step. Thus, the Pseudo-MPE solution is constructed only if the user runs through all variables proposed, accepting each one at least once.

Three slightly different operational modes of the algorithm were implemented. This was done for research purposes, in order to understand which of the three, if any or if a combination of their distinctive features, the expert users would find the most useful from a usability, comprehensibility and explainability standpoint:

- **exhaustive:** In the basic dialogue type, the set of variables under consideration monotonically decreases by one every time the user accepts one of the system’s proposals and the dialogue terminates only when the user has accepted all variables at least once or refused all proposals at a given step. In the first case the user will have the Pseudo-MPE solution while in the second she will be left with a partial assignment to some of the variables not present in the initial evidence. The pseudocode is shown in Alg. 2.
- **d-separated:** In the second variant, the set of variables under considered at each step is dynamic and depends on the separation properties of the underlying Bayesian Network’s DAG and the evidence set constructed by the user’s choices. Differently from the first type of dialogue, an additional `evidence_d_separation` subroutine is called before `next_most_probable_states` to calculate the set of variables that are d-separated from the evidence set, up to that step of the dialogue. `next_most_probable_states` is then executed but the variables that the previous function found to be separated from the ev-

idence, are removed from the returned list. This way, variables that can have no effect given the current evidence are not proposed. As the d-separation operation is not monotonic, adding new nodes to the evidence set can both increment or diminish the number of nodes that will be proposed at each step. The user is shown an updated view of the independencies of the graph at each step; an example of such an output is shown in Fig. 4.2 and 4.3. The pseudocode is shown in Alg. 3.

- **thresholded:** The final variant of the algorithm prunes the set of variables using a different strategy from the previously presented one. In this case, the $(state, value)$ pairs in the list returned by `next_most_probable_states` are dropped automatically based on their probability. Pairs whose probability is below a user-defined threshold or are “worse than random” (for ex. a $(state, value)$ tuple will be discarded if $state$ is binary and the probability of $value$ is lower than 0.5) are removed and not proposed to the user. This thresholding strategy based on the probability of the tuples is paired with one where there is a user-defined threshold on the number of times that the expert can refuse a particular $(state, value)$. In the general dialogue, tuples can be proposed multiple times, with an ever lower probability, if the user has previously refused them; in the thresholded scheme a $(state, value)$ pair can only be proposed a maximum number of times before being permanently discarded.

The underlying Bayesian Network that represents the data set is learned and queried through the pomegranate (see Subsec. 4.2.1), API but the great majority of all the code is completely custom-written. This was necessary because pomegranate, while having a powerful backend, was found to be severely lacking in the breadth and flexibility of its API. Many basic operations, such as the calculation of a joint distribution, were not available so the only way was to implement lower-level workarounds while still using pomegranate for the most basic operations, as the calculation of a posterior distribution. In particular, dialogue is implemented with the only direct calls to the API being when learning the network and when calling `predict_proba`, that queries the BayesianNetwork object to calculate the posterior distribution of the states given the current evidence. D-separation, in the second variant of the algorithm, is calculated via the `evidence_d_separation` procedure that implements the pseudocode presented in Alg. 1.

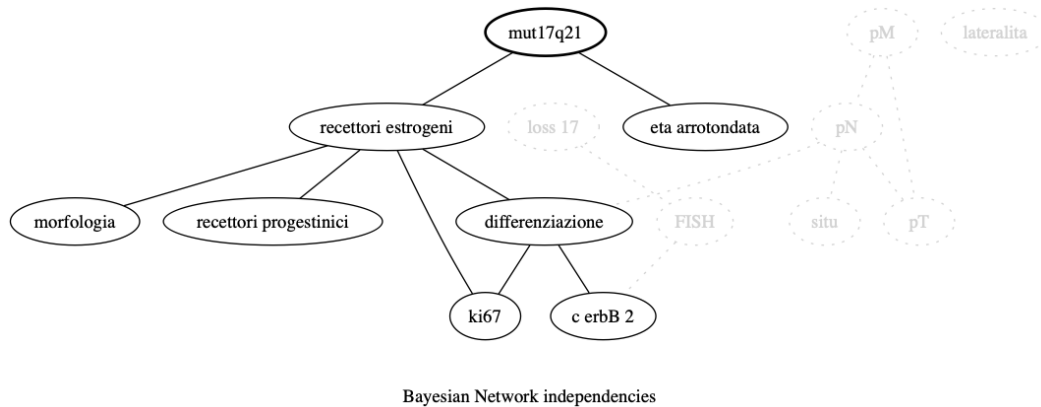


Figure 4.2. Example output during the first round of the d-separation-aware variant of dialogue. The variable “mut17q21” is the initial evidence.

Algorithm 2 Exhaustive pseudo-MPE algorithm

```

1: evidence = user selected (state, value) tuples
2: MPE_polytree = MPE Polytree rooted in evidence
3: while True do
4:   mpe_states = next_most_probable_states(evidence)
5:   if mpe_states is not empty then
6:     next_state = head of mpe_states
7:     propose next_state to user                                     ▷ the least entropic state
8:     if the user refuses next_state then
9:       for alternative_state in mpe_states  $\setminus$  next_state do
10:        propose alternative_state to user                         ▷ the next least entropic states
11:        if the user accepts alternative_state then
12:          call generate_alternative_branch() on MPE_polytree
13:          add alternative_state to MPE_polytree
14:          evidence = evidence  $\cup$  alternative_state
15:        else
16:          continue
17:        end if
18:      end for
19:    else
20:      add next_state to MPE_polytree
21:      evidence = evidence  $\cup$  next_state
22:    end if
23:  else
24:    return
25:  end if
26: end while

```

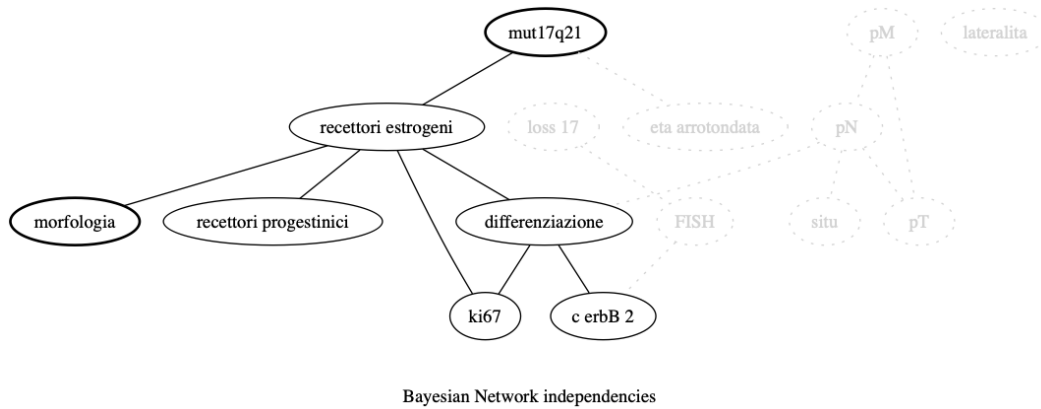


Figure 4.3. Example output during the second round of the d-separation-aware variant of dialogue. “morfologia” is added to the evidence set and this makes a large part of the network redundant.

Algorithm 3 Independencies pseudo-MPE algorithm

```

1: evidence = user selected (state, value) tuples
2: MPE_polytree = MPE Polytree rooted in evidence
3: while True do
4:   separated = evidence_d_separation(evidence)    ▷ based on evidence of previous
   step
5:   mpe_states = next_most_probable_states(evidence)
6:   mpe_states = mpe_states \ separated
7:   if mpe_states is not empty then
8:     next_state = head of mpe_states
9:     propose next_state to user                                ▷ the least entropic state
10:    if the user refuses next_state then
11:      for alternative_state in mpe_states \ next_state do
12:        propose alternative_state to user                    ▷ the next least entropic states
13:        if the user accepts alternative_state then
14:          call generate_alternative_branch() on MPE_polytree
15:          add alternative_state to MPE_polytree
16:          evidence = evidence ∪ alternative_state
17:        else
18:          continue                                           ▷ go to next proposal
19:        end if
20:      end for
21:    else
22:      add next_state to MPE_polytree
23:      evidence = evidence ∪ next_state
24:    end if
25:  else
26:    return
27:  end if
28: end while

```

Algorithm 4 Thresholded pseudo-MPE algorithm

```

1: user selected refuse_bound
2: refuse_thresholds =  $\emptyset$ 
3: lower_thresholds =  $\emptyset$ 
4: for  $v \in V$  do
5:   for value  $k$  of  $v$  do
6:     element [ $v, k$ ] = 0 in refuse_thresholds
7:   end for
8:   element [ $v$ ] =  $1/|v|$  in lower_thresholds ▷ refuse worse than random pairs
9: end for
10: evidence = user selected (state, value) tuples
11: MPE_polytree = MPE Polytree rooted in evidence
12: mpe_states = next_most_probable_states(evidence)
13: while True do
14:   for  $s \in mpe\_states$  do
15:     if probability of  $s < lower\_thresholds[s] \wedge refuse\_thresholds[s] > refuse\_bound$ 
16:       then
17:         remove  $s$  from mpe_states
18:       end if
19:   end for
20:   if mpe_states is not empty then
21:     next_state = head of mpe_states
22:     propose next_state to user ▷ the least entropic state
23:     if the user refuses next_state then
24:       increment refuse_thresholds[next_state]
25:       for alternative_state in mpe_states  $\setminus$  next_state do
26:         propose alternative_state to user ▷ the next least entropic states
27:         if the user accepts alternative_state then
28:           call generate_alternative_branch() on MPE_polytree
29:           add alternative_state to MPE_polytree
30:           evidence = evidence  $\cup$  alternative_state
31:         else
32:           increment refuse_thresholds[alternative_state]
33:           continue
34:         end if
35:       end for
36:     else
37:       add next_state to MPE_polytree
38:       evidence = evidence  $\cup$  next_state
39:     end if
40:   else
41:     return
42:   end if
43: end while

```

Alternative Explanation Branches

The function to generate alternative branches to the main MPE branch in the dialogue tree is called after the user refuses a $(state, value)$ in the dialogue and accepts one of the alternatives. The motivation is to present the user with a simple “what-if” analysis, replying to the question “Had I accepted the $(state, value)$ presented me by the system, what would have been the most probable configuration of the remaining $(state, value)$ pairs?”. The question is answered by generating a maximally probable, alternative MPE sub-branch rooted in the last node in the main MPE branch that is under construction through the dialogue.

The alternative branch is generated by what is essentially an automated version of dialogue that always accepts the first suggestion returned by `next_most_probable_states`. Given that `dialogue` and `generate_alternative_branch` are essentially one and the same, the latter inherits the same pruning strategies as the former. That is, `generate_alternative_branch` called during the exhaustive dialogue will generate a maximally likely assignment over all variables in $V \setminus E$ while when invoked from one of the other two variant of the dialogue algorithm, it will apply their same pruning strategies.

The implementation of the MPE Polytree is based on the NetworkX python package (see Subsec. 4.2.1). The creation of a chain of nodes is done by keeping a local pointer `alt_node` that refers to the last added node or set of nodes, if the node being added is the successor of multiple initial evidences.

An example of the output shown to the user at each step of the dialogue is shown in Fig. 4.4 while the pseudocode for the three variants are shown in Alg. 5, 6 and 7. Note that `alternative_evidence` is local to this algorithm and is separate from the `evidence` used in the main dialogue procedure.

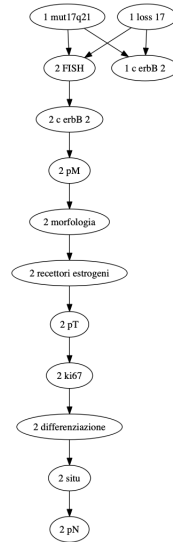


Figure 4.4. Example output during the d-separation-aware variant of `dialogue`. The tuple (“FISH”, “Aampl”) was proposed but the expert refused it and accepted the alternative (“c erbB 2”, “2+”). The main MPE Branch has ID 1 while the “what-if” one has ID 2.

Algorithm 5 Exhaustive alternative explanation branch algorithm

```

1: alternative_evidence = evidence
2: alt_node = last node in the main MPE Polytree
3: branch_id = branch_id + 1
4: while True do
5:   mpe_states = next_most_probable_states(alternative_evidence)
6:   if mpe_states is empty then
7:     return
8:   else
9:     next_state = head of mpe_states
10:    create next_state node, tag it with branch_id and make it son of alt_node
11:    update alt_node node to be next_state node
12:    alternative_evidence = alternative_evidence  $\cup$  next_state
13:   end if
14: end while

```

Algorithm 6 Independencies alternative explanation branch algorithm

```

1: alternative_evidence = evidence
2: alt_node = last node in the main MPE Polytree
3: branch_id = branch_id + 1
4: while True do
5:   separated = evidence_d_separation(alternative_evidence)
6:   mpe_states = next_most_probable_states(alternative_evidence)
7:   mpe_states = mpe_states  $\setminus$  separated
8:   if mpe_states is empty then
9:     return
10:  else
11:    next_state = head of mpe_states
12:    create next_state node, tag it with branch_id and make it son of alt_node
13:    update alt_node node to be next_state node
14:    alternative_evidence = alternative_evidence  $\cup$  next_state
15:  end if
16: end while

```

Algorithm 7 Thresholded alternative explanation branch algorithm

```

1: alternative_evidence = evidence
2: alt_node = last node in the main MPE Polytree
3: branch_id = branch_id + 1
4: while True do
5:   for  $s \in mpe\_states$  do
6:     if probability of  $s < lower\_thresholds[s] \wedge refuse\_thresholds[s] > refuse\_bound$ 
       then
7:       remove  $s$  from mpe_states
8:     end if
9:   end for
10:  mpe_states = next_most_probable_states(alternative_evidence)
11:  if mpe_states is empty then
12:    return
13:  else
14:    next_state = head of mpe_states
15:    create next_state node, tag it with branch_id and make it son of alt_node
16:    update alt_node node to be next_state node
17:    alternative_evidence = alternative_evidence  $\cup$  next_state
18:  end if
19: end while

```

“Pseudo-MPE” from Random Evidence

In order to compare the Pseudo-MPE output with the true MPE solution I implemented a simple algorithm that, starting from a random initial set of evidence, generates the relative pseudo-MPE and MPE solutions. The random initial evidence set is constructed by randomly choosing a number in the interval $k = [1, |V|]$, with V the set of vertices in the BN, and then randomly selecting k of the random variables in V to yield the set of variables E . The value of each variable is randomly chosen among the set of its values; as all variables are categorical, this can easily be done.

The implementation is based on the NetworkX python library (see Subsec. 4.2.1) package as what is being constructed is not a tree but a *polytree* (see Subsec. 3.3.2), as nodes may have multiple parents. Note that *alternative_evidence* is considered separate from the main *evidence* used in dialogue.

The pseudocode is shown in Alg. 8 while an example output can be seen in Fig. 4.5.

Inverse Explanation

da fare e trovare nome migliore

Natural Language Explanation

The probability of each proposed tuple is quantified in natural language based on the probability of the most probable value within the variable. These are shown in Tab. 4.4.

da fare e magari pensare anche a visual explanations?

Algorithm 8 Pseudo-MPE from random evidence algorithm

```

1: evidence =  $\emptyset$ 
2: generate random number  $k \in [1, |V|]$ 
3:  $S$  = choose  $k$  variables from  $V$ 
4: for  $s \in S$  do
5:   choose random  $v$  in the possible values of  $e$ 
6:   evidence = evidence  $\cup$  ( $s, v$ )
7: end for
8: MPE_polytree = MPE Polytree rooted in evidence
9: last_node = evidence
10: alternative_evidence = evidence
11: while True do
12:   mpe_states = next_most_probable_states(alternative_evidence)
13:   if mpe_states is empty then
14:     return
15:   else
16:     next_state = head of mpe_states
17:     create next_state node and make it son of last_node
18:     update alt_node node to be next_state node
19:     alternative_evidence = alternative_evidence  $\cup$  next_state
20:   end if
21: end while

```

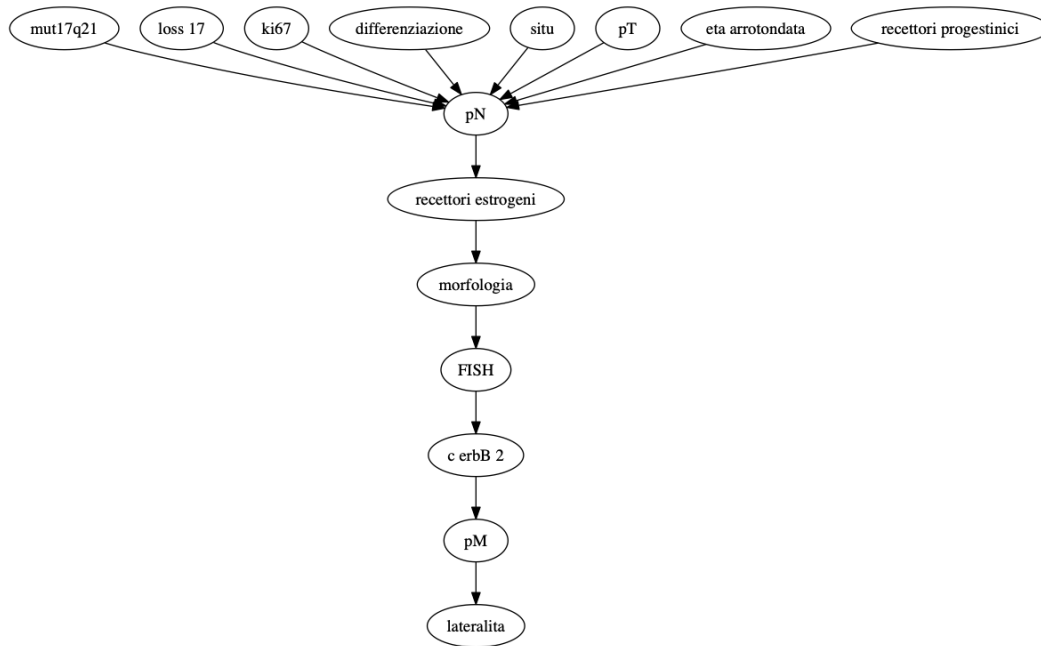


Figure 4.5. Example output of the Pseudo-MPE from random evidence algorithm.

Table 4.4. Probability quantifiers in natural language

Probability range	Natural language quantifier
(0, 0.2)	"highly unlikely"
(0.2, 0.3)	"very unlikely"
(0.3, 0.4)	"unlikely"
(0.4, 0.5)	"not plausibly"
(0.5, 0.6)	"plausibly"
(0.6, 0.7)	"possibly"
(0.7, 0.8)	"likely"
(0.8, 0.9)	"very likely"
(0.9, 1)	"highly likely"
(1)	"certain"

Pairwise Correlations

An interesting addition, in terms of both explainability and theory, is an algorithm to measure and graphically display the interrelatedness between pairs of variables. This is achieved by calculating the *conditional mutual information* between each pair of parent \rightarrow child variables. The definition of mutual information presented in Subsec. 3.2.3 is extended to account for the current state of the model i.e. the set of observed variables. What we are calculating is:

AA: e' giusto cos'?

$$I(X, Y|E = e) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{XY}(x, y|E = e) \log \left(\frac{p_{XY}(x, y|E = e)}{p_X(x|E = e)p_Y(y|E = e)} \right) \quad (4.1)$$

This is not done if the parent variable, in the $X \rightarrow Y$ tuple under consideration, is in the evidence set; this is because observing a variable in a BN conceptually corresponds to disconnecting it from its children.

The implementation takes advantage of pgmpy's (pgm) inference capabilities. To do this, I wrote a function to convert a pomegranate-based BN to an equivalent pgmpy-based one. The queries to the model are done using the variable Elimination algorithm, that is more than suitable for a small BN. The marginals for X and Y are calculated directly from the joint distributions, by marginalising the joint over one and then the other. The function exits with -1 if the parent variable is in the evidence set. As the mutual information $I(X, Y) \in [0, 1]$ this signals to the calling function to treat this couple X, Y differently.

The pseudocode for the algorithm is shown in Alg. 9.

4.4 Validation

discutere con Vittoria come validare al meglio

Having direct access to expert pathologists has not only helped in guiding research into the theoretical explainability properties of the system but also enabled their validation. There are two main validation streams to be addressed: from the clinical point of view and from the explainability one, with the results of the latter depending on those of the former.

Algorithm 9 Mutual information algorithm

```

1:  $X$  parent variable in the BN DAG
2:  $Y$  child variable in the BN DAG
3:  $E$  set of current evidence in the BN
4: if  $X \in E$  then
5:   return -1
6: end if
7:  $joint = p_{XY}(X, Y | E = e)$ 
8:  $Y\_marginal = \text{marginalise } joint \text{ over } X$ 
9:  $X\_marginal = \text{marginalise } joint \text{ over } Y$ 
10:  $mutual\_information = 0$ 
11: for  $y$  in  $Y\_marginal$  do
12:   for  $x$  in  $X\_marginal$  do
13:      $j = \text{entry in } joint \text{ corresponding to } y \text{ and } x$ 
14:     if  $j$  is 0 then
15:        $mutual\_information += 0$ 
16:     else
17:        $mutual\_information = j * \log(j / (y * x))$ 
18:     end if
19:   end for
20: end for

```

4.4.1 Clinical validation

A validation of the methods carried out in this thesis in their adherence to established clinical literature, is of paramount importance. A failure on the Bayesian Network's part in capturing the true relationships between the variables would hamper it in being able to give any meaningful representation of them. For the expert to even start to trust the system or to be able to make sense of its outputs, it is vital that there be as little cognitive dissonance between the experts' basic beliefs and expectations and those that he sees represented in the system.

For this reason, the initial validation phase with the Istituto Cantonale di Patologia concentrated on the clinical validation aspect. The methodology chosen to clinically validate the system was for a representative of the ICT to formulate a series of natural language queries; each one of these was annotated with the specific queried variable in the network and its value together with the known values of other variables. The ICT's delegate included the expected reply to the queries together with its probability, based on the latest medical literature and their personal expertise. So, each query asked for the probability truth value of the following:

$$(var_1 = a \wedge \dots \wedge var_n = z) \Rightarrow var_k = k \quad var_1 \dots var_n \in V, var_k \in V \setminus \{var_1 \dots var_n\} \quad (4.2)$$

This would translate into natural language into: "If var_1 is a and \dots and var_n is z , how probable is it that var_k k ?". The complete series of eight questions, together with the expected values, is shown in Tab. ???. It is quite evident that all these validation questions are instances of the Bayesian Network Updating problem, in particular they are Conditional Probability Queries (as defined in 3.4.2).

4.4.2 Explainability validation

Having validated the system in terms of its adherence to clinical literature, it could then also meaningfully be validated from an explainability point of view. The main question to be addressed is its capacity to relate to the expert user. Is the system able to engender the user's trust? In doing so, is she able to extract more knowledge from existing data when using the system than not? Especially in cases where there may be a dearth of data, can the expert maximise the benefit from the available information? Does the user subjectively feel that the system may positively impact her work? These are all hard questions to answer, as there is a very high degree of subjectivity involved. Thus to attempt to answer them, the chosen methods were borrowed from the social sciences.

In an earlier stage, the experts were introduced to the system in prototype form and instructed on the possible use cases it offered. Then, they were asked to keep track of times they felt they could have used the system, had it been available, during their day-to-day work and to report back to me with the query they would have submitted to the system. If the query were in the realm of possibilities offered by my program, I would submit the outputs back to the users together with the following questionnaire:

definire questionario da mandare con le domande

This process would both give feedback on the performance of the system and also inform its design.

In a later phase, having finalised the system's design, I provided it to the experts at the ICT for their use in their daily work. After **xx time** I followed up with them with an interview based on the following format:

definire formato e metriche intervista

This enabled me to quantify the performance of the system, as perceived by its users in a real setting, over an extended period of time.

4.5 Summary

Chapter 5

Results

Chapter 6

Conclusions

6.1 Future developments

Glossary

Bibliography

Graphviz. URL <https://www.graphviz.org>.

Gurobi. URL <https://www.gurobi.com>.

NetworkX. URL <https://networkx.github.io>.

NumPy. URL <http://numpy.org>.

Pandas. URL <https://pandas.pydata.org/about.html>.

Gurobi. URL <http://pgmpy.org>.

scikit-learn. URL <http://scikit-learn.github.io/stable>.

Raphaela Butz, Arjen Hommersom, and Marko van Eekelen. Explaining the Most Probable Explanation. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 11142 LNAI, pages 50–63. 2018. ISBN 9783030004606. doi: 10.1007/978-3-030-00461-3_4. URL http://link.springer.com/10.1007/978-3-030-00461-3_4.

Thomas Cover and Joy Thomas. *Elements of Information Theory*. 2006. URL <http://staff.ustc.edu.cn/~cgong821/Wiley.Interscience.Elements.of.Information.Theory.Jul.2006.eBook-DDU.pdf>.

Istituto Cantonale di Patologia. Istituto Cantonale di Patologia. URL <https://www4.ti.ch/dss/dsp/icp/chi-siamo/presentazione/>.

Kalev Kask and Rina Dechter. Stochastic local search for Bayesian networks. in *Proceedings Seventh International Workshop on Artificial Intelligence and Statistics*, 1999. URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.144.7172&rep=rep1&type=pdf>.

Daphne Koller, Nir Friedman, Sašo Džeroski, Charles Sutton, Andrew McCallum, Avi Pfeffer, Pieter Abbeel, Ming-Fai Wong, David Heckerman, Chris Meek, et al. *Probabilistic Graphical Models Principles and Techniques*, pages 26–27. MIT press, 2009a.

Daphne Koller, Nir Friedman, Sašo Džeroski, Charles Sutton, Andrew McCallum, Avi Pfeffer, Pieter Abbeel, Ming-Fai Wong, David Heckerman, Chris Meek, et al. *Probabilistic Graphical Models Principles and Techniques*, pages 74–76. MIT press, 2009b.

Vittoria Martin. Curriculum vitae. 5:13–15, 2012. URL https://m4.ti.ch/fileadmin/DSS/DSP/ICP/pdf/Collaboratori/VM{}_CV.pdf.

- Sharon Use Normand and David Tritchler. Parameter updating in a bayes network. *Journal of the American Statistical Association*, 87(420):1109–1115, 1992. ISSN 1537274X. doi: 10.1080/01621459.1992.10476266.
- Judea Pearl and Rina Dechter. Identifying Independencies in Causal Graphs with Feedback 2 BAYESIAN NETWORKS VS . *Intelligence*, (1):43, 1988. URL http://www.cs.ucla.edu/pub/stat/_ser/R243.pdf.
- PyGraphviz developer team. Pygraphviz. URL <https://pygraphviz.github.io>.
- Thomas D. Schneider. Molecular Information Theory Primer. (301), 2005. URL <http://users.fred.net/tds/lab/paper/primer/primer.pdf>.
- Jacob Schreiber. Pomegranate, a. URL <https://pomegranate.readthedocs.io/en/latest/>.
- Jacob Schreiber. Bayesian Network Structure Learning tutorial, b. URL https://github.com/jmschrei/pomegranate/blob/master/tutorials/B_Model_Tutorial_4b_Bayesian_Network_Structure_Learning.ipynb.
- Jacob Schreiber. Pomegranate: fast and flexible probabilistic modeling in python. 18:1–6, 2017. URL <http://arxiv.org/abs/1711.00137>.
- Jacob Schreiber and William Noble. Finding the optimal bayesian network given a constraint graph. 2017. doi: 10.7287/peerj-cs.122v0.2/reviews/2. URL <https://peerj.com/articles/cs-122/>.
- C. E. Shannon, W. Weaver, and R. E. Blahut. The mathematical theory of communication (TODO check jahr 48/49). *Urbana: University of Illinois press*, 117(April 1928):379–423, 1949. ISSN 07246811. doi: 10.2307/3611062. URL <http://math.harvard.edu/~ctm/home/text/others/shannon/entropy/entropy.pdf>.
- Solomon Eyal Shimony. Finding MAPs for belief networks is NP-hard. *Artificial Intelligence*, 68(2):399–410, 1994. ISSN 00043702. doi: 10.1016/0004-3702(94)90072-8. URL <http://www.cs.columbia.edu/~adrian/candp/Shi94-FindingMAPbeliefNetworksNP-hard.pdf>.
- Chandler Stolp, Shirley Dowdy, and Stanley Wearden. *Statistics for Research*, volume 3, page 637. 2006. ISBN 047126735X. doi: 10.2307/3324586.