University of Liège

Faculty of Applied Sciences Department of Electrical Engineering & Computer Science

PhD dissertation The properties of the properti



Prof. PIERRE GEURTS Advisors:

Prof. Louis Wehenkel

June 2019

JURY MEMBERS

GILLES LOUPPE, Professor at the Université de Liège (President);

PIERRE GEURTS, Professor at the Université de Liège (Advisor);

LOUIS WEHENKEL, Professor at the Université de Liège (Co-advisor);

BENOÎT FRÉNAY, Professor at the Université de Namur;

ROBIN GENUER, Professor at the Université de Bordeaux (France);

PATRICK MEYER, Professor at the Université de Liège;

ERWAN SCORNET, Professor at Ecole Polytechnique (France).

ACKNOWLEDGMENTS

To all the family members, friends and colleagues that helped me through the accomplishment of this thesis.

T•H•A•N•K • Y•O•U

Nowadays new technologies, and especially artificial intelligence, are more and more established in our society. Big data analysis and machine learning, two subfields of artificial intelligence, are at the core of many recent breakthroughs in many application fields (e.g., medicine, communication, finance, ...), including some that are strongly related to our day-to-day life (e.g., social networks, computers, smartphones, ...). In machine learning, significant improvements are usually achieved at the price of an increasing computational complexity and thanks to bigger datasets. Currently, cutting-edge models built by the most advanced machine learning algorithms typically became simultaneously very efficient and profitable but also extremely complex. Their complexity is to such an extent that these models are commonly seen as black-boxes providing a prediction or a decision which can not be interpreted or justified. Nevertheless, whether these models are used autonomously or as a simple decision-making support tool, they are already being used in machine learning applications where health and human life are at stake. Therefore, it appears to be an obvious necessity not to blindly believe everything coming out of those models without a detailed understanding of their predictions or decisions.

Accordingly, this thesis aims at improving the interpretability of models built by a specific family of machine learning algorithms, the so-called tree-based methods. Several mechanisms have been proposed to interpret these models and we aim along this thesis to improve their understanding, study their properties, and define their limitations.

The first part of this thesis introduces the techniques used to build these models, i.e. decision tree and ensemble of randomised trees induction algorithms. It also presents the basis of feature selection, a data analysis method aiming at identifying the essential features of a model and allowing to improve the model performances and/or its interpretability.

The second part of this thesis focuses on the two most popular importance measures, aiming at measuring the relative importance of features in the model, derived from tree-based methods. Our contribution in this part is two-fold. On one hand, we review the main literature on that topic, with a focus on theoretical analyses. On the other hand, we improve the theoretical characterisation of one subclass of these importance measures, known as the Mean Decrease of Impurity (MDI), and study it in greater details, both theoretically and practically.

The last part of this thesis is a collection of several works addressing some limitations of existing importance measures in some specific applications. We thus propose an extension of the MDI importance measure that can take into account different contexts in which the problem can be put, so as to provide further insight into the feature importances. We also study a new tree-based method that yields an efficient feature selection even in presence of large datasets and/or under memory constraints. Lastly we discuss the strengths and weaknesses of a solution to the network inference problem based on a tree-based importance measure, and propose a non tree-based method that we have designed as part of a network inference challenge that we eventually won.

De nos jours, les nouvelles technologies, et tout particulièrement l'intelligence artificielle, sont toujours plus ancrées dans notre société. L'analyse de grands volumes de données et l'apprentissage automatique, deux sous-domaines de l'intelligence artificielle, sont au centre des plus récentes percées dans de nombreux domaines (e.g., la médecine, la communication, la finance, ...), et en particulier des applications intimement liées à notre vie quotidienne (réseaux sociaux, ordinateurs, smartphones, ...). En apprentissage automatique, les améliorations significatives sont souvent obtenues au prix d'une plus grande complexité computationelle et grâce à des quantités de données toujours plus grandes. A l'heure actuelle, les modèles de pointe obtenus par les algorithmes d'apprentissage automatique les plus sophistiqués sont généralement à la fois très efficaces et extrêmement complexes. Leur complexité est telle qu'ils sont souvent vus comme des «boîtes noires» fournissant une prédiction ou une décision qui ne peut ni être interprétée ni être justifiée. Néanmoins, que ces modèles soient considérés de manière autonome ou comme de simples outils d'aide à la décision, ils sont déjà utilisés dans des applications d'apprentissage automatique desquelles dépendent la santé et des vies humaines. Par conséquent, il apparait comme une évidente nécessité de ne pas croire les prédictions de ces modèles aveuglément, sans les avoir comprises.

Dans ce contexte, cette thèse a pour but d'améliorer l'interprétation qui peut être faite de modèles construits par une famille particulière d'algorithmes d'apprentissage automatique basées sur les arbres de décision. Plusieurs mécanismes ont été mis en œuvre pour interpréter ces modèles et nous visons tout au long de cette thèse à améliorer leur compréhension, à étudier leurs propriétés et à en définir les limites.

La première partie de cette thèse introduit les techniques de construction de ces modèles, à savoir les arbres de décision et les ensembles d'arbres aléatoires. Elle présente également les bases de la sélection de variables, méthode d'analyse de données qui a pour but d'identifier les variables essentielles d'un problème permettant à la fois d'améliorer les performances des modèles et leur interprétabilité.

La seconde partie de cette thèse se concentre sur les deux mesures d'importance les plus populaires, visant à déterminer l'importance relative des variables dans le modèle, dérivées des méthodes à base d'arbres. Notre contribution dans cette partie est double. D'une part, nous examinons la littérature traitant ce sujet, avec une attention toute particulière pour les analyses théoriques. D'autre part, nous améliorons la caractérisation théorique d'une sous-classe de mesures d'importance, à savoir celle basée sur la réduction d'impureté (MDI), et nous l'étudions de manière détaillée théoriquement et pratiquement.

La dernière partie de cette thèse est une collection de plusieurs travaux qui se concentrent sur certaines limitations des mesures d'importance existantes dans des applications spécifiques. Ainsi, nous proposons une extension de la mesure d'importance MDI capable de prendre en compte les différents contextes dans lesquels le problème peut être placé, et cela de manière à fournir une connaissance approfondie sur l'importance des variables. Nous étudions également une nouvelle méthode à base d'arbres capable de fournir une sélection de variables performante, et ce, même en présence de grands volumes de données et/ou en cas de contraintes de mémoire. Et enfin, nous discutons les forces et les faiblesses d'une so-

lution à ce problème au d'inférence de réseaux utilisant les mesures d'importances dérivées d'arbres de décision. Nous proposons également une méthode développée lors d'une compétition d'inférence de réseaux, qui ne fait pas intervenir les arbres de décision mais qui nous a permis de remporter cette compétition.

CONTENTS

1	1.1		ıtion	
	1.2	Outline	e of the manuscript	 17
	1.3	Public	ations	 17
	5.4.0			10
		KGROL		19
2			LEARNING AND FEATURE SELECTION	21
	2.1		ne learning vs Artificial Intelligence	
	2.2		s data?	
			Nature of features	
			Interactions between features	
	2.3	•	vised learning	
			Predictions	
			Model assessment and selection	
			Other forms of learning	
	2.4		e selection for supervised learning	
		2.4.1	Relevance of features	
			2.4.1.1 On the quantitative measure of irrelevance	
			Markov boundary	
		2.4.3	Redundancy	
			2.4.3.1 Yu and Liu [2004]'s redundancy	
			2.4.3.2 Total redundancy	
			2.4.3.3 Asymmetric and partial redundancies	
			2.4.3.4 Quantitative measure of redundancy	
			2.4.3.5 Redundancy and relevance	
			2.4.3.6 Redundancy and correlation	
			Feature selection problems	
			Feature selection methods	
		2.4.6	Feature subset search algorithms	 50
		2.4.7	Discussion	 52
3	DEC	ISION	TREES AND ENSEMBLE METHODS	57
	3.1	Introdu	uction	 57
	3.2	Super	vised Learning with Decision trees	 58
		3.2.1	Semantics of tree based prediction models	 58
			3.2.1.1 From graph theory to decision tree terminology	 58
			3.2.1.2 A tree structure shaped by the features	 59
			3.2.1.3 Decision tree models	 60
		3.2.2	Learning a decision tree model from data	 63
			3.2.2.1 Splitting rules	 64
			3.2.2.2 Stopping rules and pruning	 68
			3.2.2.3 Labeling the leaves	 69
		3.2.3	Interpretability of decision tree models	 70
	3.3	Tree-b	ased ensembles	 71
		3.3.1	Random forest type of methods	 73

	3.3.2 Random Forests and Extra-Trees: parameters, properties, in-
	terpretability
	3.3.2.1 Parameters
	3.3.2.2 Variable importances
	3.3.2.3 Out-of-bag samples and estimates
	3.3.2.4 Proximity measure 8
П	CHARACTERISATION OF IMPORTANCE MEASURES 83
4	A SURVEY OF THE LITERATURE ABOUT TREE-BASED FEATURE IM-
	PORTANCE MEASURES 85
	4.1 Contribution of a feature to a tree-based model
	4.2 MDI and MDA feature importance measures
	4.2.1 MDI importance measure
	4.2.2 MDA importance measure
	4.2.3 Discussion of MDI versus MDA
	4.3 Theoretical analyses
	4.3.1 Asymptotic properties of MDA
	4.3.2 Asymptotic properties of MDI
	4.3.3 Correlated and redundant features
	4.3.3.1 MDA
	4.3.3.2 MDI
	4.4 Empirical analyses
	4.4.1 Soundness
	4.4.2 Split randomisation parameter K
	4.4.3 Feature ranking stability and number of trees
	4.4.4 Importance measures vs prediction performances
	4.4.5 Biases
	4.4.5.1 Bias due to masking effect
	4.4.5.2 Bias due to correlation
	4.4.5.3 Bias due to number of categories and scale of mea-
	surement
	4.4.5.4 Bias due to the category frequencies
	4.4.5.5 Bias due to empirical impurity estimations
	4.4.5.6 Bias due to binary splits and split value selection 107
	4.4.5.7 Bias due to bootstrapping
	4.5 Meaningful thresholds on feature importances
	4.6 Extensions and derivations
	4.7 Other importance measures
	4.8 Some applications exploiting feature importances
5	CHARACTERISATION OF MDI IMPORTANCE MEASURE 117
	5.1 Degree of relevant variables
	5.2 Totally randomised and totally developed trees
	5.3 Importances of relevant and irrelevant variables
	5.4 Impact of redundant variables
	5.5 Non-totally randomised trees
	5.6 Non-fully developed trees
	5.7 Binary trees
	5.7.1 Binary splits
	5.7.2 Relevance in hinary trees

		5.7.3 Importance scores in binary trees	. 136
	5.8	In non-asymptotic conditions	
		5.8.1 With a finite number of trees	. 138
		5.8.2 With a finite number of samples	. 140
	5.9	Result summary	. 141
Ш	EX ⁻	TENSIONS AND DERIVATIONS OF IMPORTANCE MEASURES	143
6		H A CONTEXTUAL EFFECT	145
_	6.1	Motivation	
	6.2	Context-dependent feature selection and characterization	
		Context analysis with random forests	
		6.3.1 Variable importances ¹	
		6.3.2 Identifying context-dependent variables	
		6.3.3 Characterizing context-dependent variables	. 152
		6.3.4 In practice	
		6.3.5 Generalization to other impurity measures	. 154
	6.4	Experiments	. 156
	6.5	Conclusions and future work	. 160
ΑP	PENE	DICES	161
	6.A	Details of Example 6.1	. 161
	6.B	Proof of Theorem 6.1	. 161
	6.C	Proof of Theorem 6.2	. 162
	6.D	Proof of Theorem 6.3	. 163
	6.E	Results for Problem 3	. 163
7	IN V	ERY HIGH DIMENSIONS	165
	7.1	Motivation	. 165
	7.2	Sequential random subspace	. 166
	7.3	Theoretical analysis	. 167
		7.3.1 Soundness	
		7.3.2 Convergence	
	7.4	1	
	7.5	Conclusions and future work	. 176
ΑP		DICES	177
		Proof of Theorem 7.2	
	7.B	Convergence analysis	
		7.B.1 Simplifying assumptions	
		7.B.2 Average times	
		7.B.3 Markov chain interpretation	
	7.C	Details for Section 7.4	. 181
		7.C.1 On the use of a random probe to distinguish relevant features	
		from irrelevant features	
		7.C.2 On the datasets and on the protocol	
		7.C.3 Detailed results	
8		WORK INFERENCE AND CONNECTOMICS CHALLENGE	
		Motivation	
	8.2	Tree-based network inference based on variable importances	
		8.2.1 GENIE3	. 187

¹This section is a reminder of the MDI importance measure and its asymptotic characterisation. See Section 5.2 for more details.

		8.2.2	Direct interaction	187
		8.2.3	Edge orientation	188
	8.3	Conne	ctomics challenge	189
		8.3.1	Preamble	189
		8.3.2	Connectome inference	190
		8.3.3	Signal processing	190
		8.3.4	Connectome inference from partial correlation statistics	192
		8.3.5	Experiments	193
		8.3.6	Conclusion for connectome inference	195
ΑP	PEND	ICES		197
	8.A	Descri	ption of the "Full method"	197
		8.A.1	Signal processing	197
		8.A.2	Weighted average of partial correlation statistics	198
		8.A.3	Prediction of edge orientation	198
		8.A.4	Experiments	199
	8.B	Supple	ementary results	199
	8.C	On the	selection of the number of principal components	200
9	CON	ICLUSI	ON	203
	9.1	Main fi	indings	203
	9.2	Limitat	ions and future work	205
	9.3	Open	research questions	206
IV	API	PENDI	CES	231
Α	NOT	ATION	S AND SYMBOLS	233
В	NOT	ATION	S, AND DEFINITIONS OF ENTROPIES AND MUTUAL IN-	
	FOR	MATIO	N	239
С	DIG	IT REC	OGNITION PROBLEM	241

INTRODUCTION

1.1 MOTIVATION

From Alan Turing and Claude Shannon in the 1940's and the birth of computer science to the recent breakthroughs in the Internet of Things and in Artificial Intelligence (AI), the scientific and technological worlds of data collection and computing have tremendously evolved. In the last 20 years, this phenomenon has been accelerating significantly. There was a real "boom" in terms of new discoveries and breakthroughs. Among those recent and popular successes, many were made in the field of Machine Learning (ML). This field unifies all researches that aim at equipping machines (high performance computer grids, robots, cars, smart-phones, etc.) with the ability of learning a new task, and then improving their performances, by the mere fact of exploiting more data. Let us mention for example the famous softwares of Google, AlphaGo and AlphaGoZero, that learned how to play and even become champion of the game of Go as well as several other highly complex boardgames. Progresses in ML are either dedicated to help researchers to exploit growing empirical datasets in their fields (e.g., physics, medicine, environmental sciences, social sciences, linguistics...) or to improve day-to-day life. ML applications include sorting incoming e-mails, translating text (e.g., Google translate, DeepL), understanding and producing spoken language (e.g., Siri from Apple, Ok Google, Alexa from Amazon), and even self-driving cars and autonomous robots.

Following the main trend of the ML domain, those applications are constantly improved with the avowed goal of always achieving better performances and reducing the costs. In machine learning, significant improvements are usually achieved at the price of an increasing computational complexity and thanks to bigger datasets. Currently, cutting-edge models built by the most advanced machine learning algorithms are commonly seen as black-boxes because they are either too complex to be comprehensible, or because they are kept secret by their owners.

In the future, there will be countless new ML applications in which human health and life are at stake. Making a diagnosis (i.e., identification of a disease), estimating a prognosis (i.e., predicting the expected development of a disease), personalising a medical treatment and so many other medical decisions are already available or currently developed. It is obvious that one will not blindly believe everything coming out of those machines. The failure of Google Flu (predicting flu pandemics) illustrates that machines are not always infallible, but they may be of great help. To gain trust in machine learning based solutions, it is and will remain crucial to understand these algorithms and the reasons behind their decisions or predictions in a given application, e.g., examine choices of (military) autonomous drones and self-driving cars and knowing why Google Death forecasts someone's near death. That is one of the reasons motivating a second trend in ML focusing on the interpretability of models rather than on their mere predictive and computational performances only (see, e.g., [Lipton, 2016; Doshi-Velez and Kim, 2017]). An interpretable model means that one understands the problem that is modelled and apprehends the underlying inference mechanism. Therefore, in some circumstances, the preference is for an

interpretable model, that is not necessarily the most accurate or the fastest one but that manages to extract relevant knowledge from the data.

Performances and interpretability are typically not concomitant and a trade-off between those two properties is usually a desirable feature for a ML method. Works are then made to improve the interpretability of existing black-box approaches while others focus on boosting the performances of already interpretable models.

Among the broad set of existing machine learning methods, this thesis only considers tree-based models. Within that kind of methods, single decision trees are very popular method and considered as highly interpretable. The model takes the form of a tree-structured graph representing a sequential reasoning to take a complex decision. The interest of the model is that it follows the reasoning everyone can make to handle difficult problems. However, this approach often provides highly variable models (because of the greedy nature of the approach) which in turn leads to rather modest levels of accuracy. In this thesis, special attention will be given to tree-based ensemble methods, also known as random forests (RF). While improving significantly the accuracy with respect to single trees, they unfortunately provide also much less interpretable models. With an ensemble of trees, many different explanations for a single decision are aggregated and interpreting the resulting prediction is not possible any more.

As mentioned, some efforts are usually made to interpret accurate models and, in this case, to recover some of the interpretability of a single decision tree. This can be done by identifying the constitutive elements (variables) of the model and their relative importances. For example, trying to predict someone's wine taste, we could determine that the wine colour is quite important and plays a decisive role in the wine taste discovery. In the literature on random forests, several different so-called 'feature-importance' measures have been proposed in order to restore some interpretability, and also in order to help selecting relevant subsets of features, whenever this is useful.

Despite their success, RF methods and in particular importance measures derived from these models still contain some grey areas:

- (a) Parameters of the methods have been usually studied with the scope of maximising the model performances. How do these parameters impact the quality of importance measures? Are optimal values for performances similar to those providing the best understanding of the problem?
- (b) What is actually measured by an importance measure? Is it its usefulness in the model? Does the importance evaluate the contribution of the variable in the model? How is defined the contribution of a variable?
- (c) Are those importance measures consistent? Are all variables equally treated when their importance is evaluated?
- (d) For a given importance measure, one can retrieve a numerical score for each variable. Is this sufficient to interpret all kinds of data structures, such as interacting features?

Along this thesis, we focus on answering some of those important questions in the light of our own work and of major contributions from the literature. We also propose some improvements to respond to some of the main limitations of the importance measures.

1.2 OUTLINE OF THE MANUSCRIPT

The first part of this manuscript aims at summarising important notions about supervised learning, feature selection and tree-based methods. In particular, Chapter 2 describes the different natures and roles of variables and how they may interact together to form complex structures. Then, in the context of supervised learning, the interest of a variable is formalised by various notions of relevance and redundancy. This chapter is concluded by a description of feature selection problems and methods, that aim at using a dataset to find the most relevant features in order to improve performances of machine learning models and/or to improve their interpretability. Chapter 3 introduces tree-based models: from the single decision tree algorithm to state-of-the-art tree-based ensemble methods. Some key points or methods are highlighted for a better understanding of the subsequent chapters.

The second part of this manuscript is dedicated to the most popular importance measures derived from tree-based ensembles. In particular, Chapter 4 reviews the main literature on that topic, with a focus on theoretical analyses. Chapter 5 then focuses on one subclass of these importance measures (known as the Mean Decrease of Impurity (MDI)) and studies it in greater details, both theoretically and practically.

The third and last part collects several contributions made in order to improve existing importance measures and/or in the context of some specific applications. In particular, Chapter 6 proposes an extension of the MDI importance measure to take into account different contexts in which the problem can be put, so as to provide further insight into the feature importances. Chapter 7 describes a new method using tree-based ensembles to perform feature selection under memory constraints. Finally, Chapter 8 considers the network inference application. Its first part describes a tree-based solution and highlights some of the limitations of the method facing some challenges of network inference. The second part focuses on a network inference challenge and a non tree-based approach that we have designed in order to win the competition.

1.3 **PUBLICATIONS**

This dissertation summarises several contributions to tree-based importance measures. Publications that are directly related to this work include:

 G. Louppe, L. Wehenkel, A. Sutera, and P. Geurts. Understanding variable importances in forests of randomized trees. In Advances in neural information processing systems, pages 431-439, 2013

This publication is of interest in Chapters 4 and 5.

 A. Sutera, A. Joly, V. François-Lavet, A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In Neural Connectomics Workshop, pages 23-35, 2015

Chapter 8 is the result of that publication.

- A. Sutera, G. Louppe, V. A. Huynh-Thu, L. Wehenkel, and P. Geurts. Contextdependent feature analysis with random forests. In Uncertainty In Artificial Intelligence: Proceedings of the Thirty-Second Conference, 2016 Chapter 6 is the result of that publication.
- A. Sutera, A. Joly, V. François-Lavet, Z. A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In Neural Connectomics Challenge, pages 23-36. Springer, 2017

Second version of [Sutera et al., 2015].

 A. Sutera, C. Châtel, G. Louppe, L. Wehenkel, and P. Geurts. Random subspace with trees for feature selection under memory constraints. In A. Storkey and F. Perez-Cruz, editors, Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics, volume 84 of Proceedings of Machine Learning Research, pages 929-937, Playa Blanca, Lanzarote, Canary Islands, 09-11 Apr 2018. PMLR. URL http://proceedings.mlr.press/ v84/sutera18a.html

Chapter 7 is the result of the methodological part of that publication. Theoretical results of that publication are also of interest in Chapters 4 and 5.

During the course of this thesis, several fruitful collaborations have also led to the following publications. These are not discussed within this dissertation.

- D. Taralla, Z. Qiu, A. Sutera, R. Fonteneau, and D. Ernst. Decision making from confidence measurement on the reward growth using supervised learning: A study intended for large-scale video games. In Proceedings of the 8th International Conference on Agents and Artificial Intelligence (ICAART 2016)-*Volume 2*, pages 264–271, 2016
- F. Olivier, A. Sutera, P. Geurts, R. Fonteneau, and D. Ernst. Phase identification of smart meters by clustering voltage measurements. In *Proceedings of* the 20th Power Systems Computation Conference (PSCC 2018), 2018
- M. Wehenkel, A. Sutera, C. Bastin, P. Geurts, and C. Phillips. Random forests based group importance scores and their statistical interpretation: application for alzheimer's disease. Frontiers in Neuroscience - Brain Imaging Methods, 2018

Part I BACKGROUND

Overview

The goal of this chapter is to provide some general background in supervised machine learning and feature selection. We start with a short motivation about the role of machine learning in the context of artificial intelligence. Then we discuss data structures and supervised learning problems. The bulk of the chapter focuses on feature selection methods. Along the way, we also introduce terminology, and some related mathematical notions and notations.

"Can machines think?"

— Alan Turing, 1950

2.1 MACHINE LEARNING VS ARTIFICIAL INTELLIGENCE

By studying the possibility of a machine to think, which led to his famous test to establish human level intelligence of a machine, Turing [1950] laid the foundation stone for a new field of research, called Artificial Intelligence (AI). Since then, in their quest to give a sort of intelligence to machines (and most prominently to computers), scientists have developed theories and algorithms to enable computers to learn from examples. This topic forms a sub-domain of AI called machine learning (ML). The goal of machine learning is to allow a machine to progressively improve its ability to solve some tasks by exploiting some relevant data collected over time. This contrasts with the habit of classical programming that implements computer programs based on a frozen set of human-based knowledge. Learning algorithms may actually allow a machine to discover knowledge that was missed by human experts or that is too complex to be discovered by them. Thus, the purpose of ML methods is dual. On the one hand, ML methods aim at producing models derived from data that allow for accurate predictions, e.g. to take decisions or to guess not yet observed values. On the other hand, those models need to be interpretable in order to help humans to explore data and understand complex systems. Both goals however equally require the same thing: (a lot of) data. That is why the next section presents the notion of data and its constituent elements known as observations and features.

2.2 WHAT IS DATA?

In the context of this thesis, a *dataset* \mathbf{D} is a collection of data and is organised as a set of N observations $\{\mathbf{o}^i\}_{i=1}^N$. An *observation* \mathbf{o}^i , also called *sample* or *example*, is a (line)vector of p values $\mathbf{o}^i = (o_1^i, \ldots, o_p^i)$, where the element o_j^i corresponds to

the value of the feature j. A feature (or equivalently a variable¹) is a function taking as argument an object (belonging to some underlying set of possible objects) and whose values belong to a certain domain.

A dataset of N observations described by p features is usually represented by a matrix of size $N \times p$.

A large dataset refers to a dataset where N is very large while a small dataset refers to a dataset where N is small. A high-dimensional (respectively low-dimensional) dataset corresponds to the case where p is very large (respectively small), while a big dataset corresponds to the case where $N \times p$ is very large. From a statistical viewpoint, the number N of samples should ideally be (much) larger than the number p of features in order to cover sufficiently well all possible combinations of features values. In practice, datasets with N \ll p are often encountered and they indeed raise important challenges in the learning process [Kuncheva and Rodríguez, 2018].

2.2.1 Nature of features

In machine learning, a feature encodes some observed information by taking a value from its domain. The number of possible values and the relationship between them allow to define several types of features, listed hereunder.

CONTINUOUS A feature is continuous if it can take any value within an interval of R. This results in an uncountable number of possibilities, and one can always find a new value between two other ones as close as they can be. A continuous feature is also *ordered*: its values are inherently numerical and hence they are (logically) ordered.

A few examples of continuous features are height (domain is \mathbb{R}^+), weight (\mathbb{R}^+) , time (\mathbb{R}^+) , speed (\mathbb{R}) , flow (\mathbb{R}) , correlation score ([-1,1]), error rate ([0,1]).

RESCALING OF CONTINUOUS FEATURE VALUES

A continuous feature may be rescaled without loss of information by mapping its domain to [-1, 1] or [0, 1] for instance. In the same machine learning application, ranges of different continuous features may vary widely from each other and some machine learning algorithms (e.g., artificial neural networks or support vector machines) might require to rescale all continuous features to the same range to work properly (e.g., by helping or speeding up optimisation) or to compare features with each others (e.g., in k-nearest-neighbours so that all features can contribute equally).

DISCRETE A feature is discrete when it takes its values in a set of at most a countable (and usually finite) number of values. Its values can either be numerical or categorical, ordered or not. The number of possible values defines the cardinality of such a feature. A m-ary feature (i.e., a feature of cardinality m) can take $\mathfrak m$ different values. In particular, a feature of cardinality two is a binary feature and its set of possible values is typically represented as $\{0,1\}$ or $\{-1,1\}$.

¹Both terms will be used in this thesis without distinction.

Usually, discrete features are divided in three sub-types:

- A numerical discrete feature takes on numerical values from a countable or finite subset of \mathbb{N} or \mathbb{R} . Its values are thus naturally *ordered*.
 - Examples of numerical discrete features usually refer to counts or proportions of indivisible elements: the number of children, the number of passengers, the proportion of expensive of cars, etc.
- An ordinal discrete feature takes on values that are not numerical but are still following a logical order.
 - Examples of ordinal discrete features usually refer to a scale, a degree of magnitude and can often straightforwardly be replaced by numerical values if necessary: position {first, second, third}, the degree of severity (of a car accident, a disease) {low, intermediate, high}, the coffee strength {mild, strong}, etc.
- A categorical discrete feature (also known as nominal discrete feature) takes on values from an finite set of elements without logical order. Values, referred to as *classes* or *categories*, are *unordered*.
 - Examples are eye colour taking values in {blue, brown, green}, mood ∈ {happy, sad}, etc.. While they are not ordered, they can however be encoded as numerical values if necessary (see side note on page 24).

When only the existence of a logical order between the feature values is of interest, continuous, numerical or ordinal discrete features are united as ordered features and, conversely, categorical discrete features are unordered features.

2.2.2 Interactions between features

Beyond their individual natures, the relations between features may also play a key role. Indeed, features can be seen as individual entities that carry some information (e.g., a value), but to consider features to their full extent, they need to be seen in the context of other features possibly interacting with them.

In what follows, we first define a model of interacting features and then focus on the interactions between variables.

A (causal) model of interacting features

Following Pearl [2009a]'s definition, a (causal) model is a triple M = (U, V, F) [White et al., 2011], where

- \bullet U is a set of background variables $\{u_1,\ldots,u_m\}$ that are determined outside the model. Such variables are also called exogenous.
- V is a set of variables $\{v_1, \dots, v_n\}$ that are determined within the model. Such variables are also called endogenous.
- F is a set of functions $\{f_1, \ldots, f_n\}$ specifying how each endogenous variable is determined by other variables of the model. More precisely, each fi provides the value of v_i given the values of a subset of all other variables $U \cup V^{-i}$ where V^{-i} is the set V without the variable v_i (i.e., $V^{-i} = V \setminus \{v_i\}$).



Numerical encoding of categorical features

Some methods (e.g., neural networks, support vector machines) are not able to handle features with non-numerical values (i.e., ordinal and categorical discrete features). Values of such features thus need to be encoded, converted into numerical values.

With an ordered feature, one can easily attribute a numerical value to each possible values while respecting the logical order between them (e.g., {low, middle, high} into $\{1,2,3\}$ and low < high is preserved through 1 < 3). Similarly, numerical values can be assigned to each class of a categorical feature. For example, let us take a categorical feature representing the eye colour with possible classes {blue, brown, green}. A classical numerical encoding would give $\{blue = 1, brown = 2, green = 3\}$. However, this introduces an order between the classes that was not originally there. Having blue eyes is not "lower" than having brown eyes but assigned numerical values (1 and 3) induce a spurious ordering.

Another encoding consists in replacing a categorical feature by several binary features B. Two binary variables are enough to perfectly encode a variable with four different classes (x binary variables give up to 2^x combinations). However, all binary variables are required to unambiguously retrieve the value. This is the binary equivalent of the classical encoding.

One-hot encoding associates one binary feature hi to each possible value of the original feature such that the binary value is equal to 1 only if the original feature has the corresponding class (e.g., h_1 corresponding to blue). In this case, a larger number of binary features are required to represent all possible values of the original feature but there is no ordering implied by this encoding. A summary is made in Table 2.1.

Eye colour	Classical Enc.	Binary Enc.		One-Hot End		
		b ₁	b_2	h ₁	h ₂	h ₃
blue	1	0	0	1	0	0
brown	2	0	1	0	1	0
green	3	1	1	0	0	1

Table 2.1: Example of different encodings of a categorical variable

The structure of such a model may be represented in the form of a directed graph, where each vertex corresponds to one of the (exogenous or endogenous) variables, and where for each endogenous variable v_i there is an edge pointing to its vertex from each one of the vertices corresponding to the other variables actually intervening in the function f_i. More details about the associated graph and uniqueness are given in [Pearl, 2009a].



EXTENDING THE MODEL

Some exogenous variables may become endogenous if one extends the (causal) model by adding new features ($\notin V \cup U$). In some way, the characterisation associated to one feature will depend on the considered model.

Based on this characterisation, endogenous and exogenous variables are particularly interesting in terms of interactions between variables. In the following section, we characterise some of those interactions.

Direct, indirect and confounded interactions between variables

From the previous section, it appears that variables may interact with each others. An endogenous variable v_i is determined by (potentially) all other variables in V^{-i} . It means that some variables in V^{-i} interact with v_i to determine its value. Let us notice that exogenous variables interact with endogenous variables asymmetrically. Indeed, they can influence the value of variables in V but their values can not be determined, as defined, by variables in V. On the other side, interactions implying endogenous variables can be symmetrical because one endogenous variable v_i may influence and be influenced by the value of another endogenous variable v_i .

Let us extend the characterisation of interacting variables to include indirect influences of variables.

INTERMEDIATE VARIABLE is a variable providing a (causal²) link between two other variables³. Let us consider two variables x and y. There may be a (causal) path going from x (a cause) directly to y (an effect), or indirectly through some intermediate step(s). A variable (i.e., the intermediate step) on the pathway from x to y is an *intermediate variable*. An intermediate variable mediates the effect of x on y. Figure 2.1 shows an example of model with an intermediate variable z between x and y. Practically, the starting point (the source) x may be a treatment or an exposure and the ending point y may be a survival status or a disease [Deng et al., 2013]. For example, let us associate x with a certain drug that affects the heartbeat, y with the survival status of a patient. One may observe that the drug have a positive effect on the survival of the patient. However, the drug does not directly modify the survival status. Actually, the drug helps to regulate the heartbeat which in turn may improve the survival expectation of the patient. In this example, the heartbeat is an intermediate variable between the treatment and the outcome [Deng et al., 2013].

²Causality is not specifically addressed in this thesis (see reference text book [Pearl, 2009a] for more details on causality). Many scientific fields, such as medicine or economy, are however interested in causal mechanisms and study the effect of intermediary variables and confounders (see, e.g., [Pearl, 2001, 2009b; Deng et al., 2013; Ananth and Schisterman, 2017]).

³Such a variable is also known as an intervening, mediating or intermediary variable.

A more trivial example is the relationship between the income and the life expectancy. One can not actually "buy" a longer life but money can contribute to better medical care that help to live longer. In this case, the quality of medical care is the intermediate variable.

From that, we can define the *direct effect* as the influence of x on y that is not mediated by other variables [Pearl, 2001]. Conversely, the indirect effect is the influence of x on y that is mediated by other variables.



Figure 2.1: Example of model with an intermediary variable z in the pathway from the cause x to the effect y.

DIRECT AND INDIRECT EFFECTS SIMULTANEOUSLY

There may be several paths from x to y and so x may have simultaneously direct and indirect (through intermediates variables) effects on y. Figure 2.2 illustrates two paths: a direct one and another that goes through an intermediate variable z. In this case, the indirect effect is meant to quantify the influence xthrough indirect paths only. One may notice that this is not practically possible to block paths (i.e., holding a set of variables constant) such that the direct pathway would be circumvented. More thorough definitions of direct and indirect effects are given in [Pearl, 2001].



Figure 2.2: Framework where there is a direct path from x to y and an indirect path from x going through z to y.

CONFOUDING VARIABLE OR CONFOUNDER is a (unstudied, exogenous) variable, say z, which influences two other variables x and y (conditionally or not to x), and tends to confound our reading of the effect of x on y [Pearl, 2009b; Li et al., 2011]. Figure 2.3 gives a possible model where x and y are confounded by a third variable z that influences both x and y (conditionally or not to x). As illustrative example, let us examine an example proposed by Kamangar [2012]: the risk of Down's syndrome⁴ for a newborn baby. Let us associate x with the parity (i.e., mother's number of pregnancies), y with the Down's syndrome (i.e., whether or not the baby is affected by the syndrome), and z with the maternal age (i.e., mother's age when giving birth to the baby). Researches that only consider parity and the risk of Down's syndrome tend to show that the risk for a baby to be affected is associated with the number of his/her mother's pregnancies. For instance, the first-born has lower risk to be affected by the Down's syndrome than the fifth one. However, one needs to take the maternal age into account to determine the real association between the parity and

⁴The Down's syndrome is a genetic disorder caused by the presence of an extra copy of human chromosome 21 [Patterson, 2009].

the risk of Down's syndrome. The fifth children of a young 30-year-old mother has actually lower risk of getting affected than the first baby of a 40-year-old mother. In this case, the mother's age is a confounder⁵ that accentuates the effect of parity on the risk of being affected by the syndrome. Many studies (e.g., in bioinformatics [Li et al., 2011], in ecology [Ewers and Didham, 2006], in medicine [Møller et al., 2000; Del Campo et al., 2012; Ananth and Schisterman, 2017; Wu et al., 2018]) focus on the effect of confouding factors as a way of taking another look at previous observations.

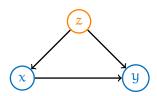


Figure 2.3: Example of model with a confounding variable z for x and y.

When the confounding bias comes from contextual elements (e.g., the specific conditions in which an experiment is made), these circumstances are assumed to be encoded by a specific context variable, further referred as a contextual variable. When taking into account the context, some feature dependencies may be accentuated or toned down while other may be unchanged being non-contextual (see side note about Simpson's paradox on page 28).

2.3 SUPERVISED LEARNING

In all generality, machine learning consists in learning models from data. This learning can be supervised when used data is labelled, i.e., where each sample is associated with a label or a specific value. Supervised learning thus focuses on learning a model from a learning set (i.e., labelled data) that can be used to predict the label of new (unseen) objects.

A learning set LS is a collection of input-output pairs [Liu and Wu, 2012; Schrynemackers, 2015]

$$LS = \{(x^1, y^1), \dots, (x^N, y^N)\} \in (\mathfrak{X} \times \mathfrak{Y})^N$$

where \mathcal{X} and \mathcal{Y} are respectively the input and output spaces, $\mathbf{x}^i = \{x_1^i, \dots, x_p^i\} \in \mathcal{X}$ is the vector of the i^{th} sample made of p input variable values and $y^i \in \mathcal{Y}$ is the corresponding output⁶ (label).

From a learning set, a supervised learning algorithm A aims at finding a function $f: \mathcal{X} \to \mathcal{Y}$ that expresses the relationship between the inputs and the output. Such a model is able to provide a prediction $f(x) = \hat{y}$ approximating the true value y for a new input vector x.

Section 2.3.1 focuses on the prediction of a supervised model. Section 2.3.2 defines the relevant notions of error for model assessment and selection. Section 2.3.3

 $^{^5}$ Let us note that a confounder is not on the path and can not be an intermediate variable. The number of pregnancies of a woman does not influence her age.

⁶Typically, there is only one output to predict as it will be the case in this thesis. However, sometimes applications require to predict several outputs simultaneously (e.g., the full state of a system in power system management). Learning with more than one output is called multi-output learning (see, e.g., Joly [2017]).

SIMPSON'S PARADOX [SIMPSON, 1951; PEARL, 2009B]

The Simpson's paradox [Simpson, 1951] refers to a setting where there is a trend in a given population p, and, at the same time, this trend disappears or reverses in every subpopulation of p. Pearl [2009b] formalises^a it as follows:

"An event C increases the probability of E in a given population p and, at the same time, decreases the probability of E in every subpopulation of p. In other words, if F and \neg F are two complementary properties describing two subpopulations^b, we might well encounter the inequalities:

$$P(E|C) > P(E|\neg C) \tag{2.1}$$

$$P(E|C,F) < P(E|\neg C,F) \tag{2.2}$$

$$P(E|C, \neg F) < P(E|\neg C, \neg F)$$
(2.3)

[...] For example, if we associate C with taking a certain drug, E with recovery, and F with being a female then - under the causal interpretation of Equations 2.2 and 2.3 - the drug seems to be harmful to both males and females yet beneficial to the population as a whole (Equation 2.1). Intuition deems such a result impossible, and correctly so."

Such paradoxical setting - yet surprising - shows that this is possible to have a certain effect (or no effect in case of equality) without considering an external factor (here, F) and opposite effects when taking into account this factor (see Chapter 6 in which variable F will refer to some contextual conditions, i.e., a contextual variable).

^aPearl [2009b] consciously chooses letters C and E to connote with *cause* and *effect*.

 $[^]b$ Symbol \neg is the logical *not* operator. \neg F refers to the complementary value of F, i.e., not F.

briefly presents other forms of learning but only supervised learning is considered in the rest of this thesis.

2.3.1 **Predictions**

The output variable, also known as a target variable, can be either continuous or discrete and the learning algorithm must take this nature into account. The learnt model thus differs depending on the nature of the variable to predict. The performance of the model (i.e., the quality of its predictions) is usually measured by means of a loss function L: $y \times y \to \mathbb{R}^+$ (see Section 2.3.2). It provides a numerical score based on the comparison of the predictions with the targeted (actual) values.

Two kinds of models are defined:

A CLASSIFICATION MODEL predicts the value of a discrete output. This model typically chooses its prediction from a set of pre-defined values (e.g., usually output values in the learning set) and is thus unable to predict an unseen value (e.g., predicting yellow if only blue, brown and green have been observed in LS). A typical loss function for a classification model is the zero-one loss $L^{0-1}(f(\mathbf{x}), \mathbf{y}) = \mathbb{1}(f(\mathbf{x}) \neq \mathbf{y})$ which is equal to 1 if the condition is verified (i.e., if the prediction is wrong and differs from the real value) and otherwise equal to zero.

A REGRESSION MODEL predicts the value of a continuous output. This model is usually able to produce new output values different from those found in the learning set (e.g., by averaging subsets of these latter values). A typical loss function for regression is the squared error (SE) $L^{se}(f(x), y) = (y - f(x))^2$ which computes the difference between the prediction and the real values exaggerating large deviations by taking the square of the difference. Another common loss function is the absolute error $L^{ae}(f(\mathbf{x}), \mathbf{y}) = |\mathbf{y} - f(\mathbf{x})|$.

The model and the loss functions must be chosen accordingly with the considered application. Let us note that when trying to predict the value of an ordered discrete variable, one can also use a regression model. Given the logical order between the values, even an unseen predicted value can be related to the others.

2.3.2 Model assessment and selection

In this section, we focus on the assessment of the prediction performance of a model f. Let us consider a set of input variables $X = \{x_1, x_2, \dots, x_p\}$ and an output variable y. We denote $P_{x_1,x_2,...,x_p,y}$, or equivalently $P_{X,y}$, the joint probability density of variables x_1, x_2, \dots, x_p, y and $P_{y|x_1, x_2, \dots, x_p}$, or equivalently $P_{y|X}$, the conditional density of y given variables $x_1, x_2, ..., x_p$.

Given a loss function L (e.g., L^{0-1}, L^{se}, L^{ae}), the goal of supervised learning is to find a model f which minimises the prediction error over an independent test set (usually drawn from the same distribution than the learning set), and defined as follows:

Definition 2.1. The generalisation error (a.k.a., test error or expected prediction error) is the expected⁷ value of the loss function

$$Err(f) = \mathbb{E}_{X,y}\{L(f(X),y)\}$$
 (2.4)

over X and y randomly drawn from their joint distribution $P_{X,y}$.

Given a model \hat{f}_{LS} learnt from a learning set LS, its generalisation error is

$$\operatorname{Err}(\hat{\mathbf{f}}_{LS}) = \mathbb{E}_{X,u}\{L(\hat{\mathbf{f}}_{LS}(X), y)\}. \tag{2.5}$$

Another quantity of interest is the *expected generalisation error* $\mathbb{E}_{LS}\{\mathrm{Err}(\hat{\mathbf{f}}_{LS})\}$ over random learning sets of size N. Typically, $\mathrm{Err}(\hat{\mathbf{f}}_{LS})$ is used for model assessment and selection while $\mathbb{E}_{LS}\{\mathrm{Err}(\hat{\mathbf{f}}_{LS})\}$ is useful to characterise a learning algorithm.

From the distribution $P_{X,y}$ of a given problem and for a given loss function, it is actually possible the derive analytically and independently of any learning set the best possible model. First, let us rewrite the generalisation error by conditioning on X:

$$Err(f) = \mathbb{E}_{X,y}\{L(X), y\} = \mathbb{E}_{X}\{\mathbb{E}_{y|X}\{L(f(X), y)\}\}.$$
 (2.6)

From that, let us define the best possible model as follows:

Definition 2.2. The best possible model f_B , known as the **Bayes model**, that minimises $Err(f_B)$ is the one that minimises the inner expectation at each point x of the input space, that is:

$$f_B(X) = \underset{y' \in \mathcal{Y}}{\text{arg min}} \mathbb{E}_{y|X} \{ L(y', y) \}. \tag{2.7}$$

The generalisation error $Err(f_B)$ of the Bayes model is referred to as the *residual* error.

However, the joint distribution $P_{X,y}$ is usually unknown in practice and one needs to estimate the generalisation error from available data. Let us define the *average* prediction error as the average loss over a set LS' of N' observations (possibly different from the learning set LS used to learn \hat{f}_{LS}), that is,

$$\widehat{\mathsf{Err}}(\hat{\mathsf{f}}_{\mathsf{LS}}, \mathsf{LS}') = \frac{1}{\mathsf{N}'} \sum_{(\mathsf{x}^i, \mathsf{y}^i) \in \mathsf{LS}'} \mathsf{L}(\hat{\mathsf{f}}_{\mathsf{LS}}(\mathsf{x}^i), \mathsf{y}^i). \tag{2.8}$$

When \mathbf{LS}' is identical to the learning set \mathbf{LS} used to learn the model, $\widehat{\mathrm{Err}}(\widehat{\mathbf{f}}_{\mathbf{LS}}, \mathbf{LS})$ is known as the *training error* or *empirical risk*. Another approach, known as the *test set method*, consists in dividing the available learning set in two disjoint sets $\mathbf{LS}_{\mathrm{train}}$ (*training set*) and $\mathbf{LS}_{\mathrm{test}}$ (*test set*) that are respectively use to learn the model and estimate the generalisation error⁸. Similarly, the K *fold cross-validation*

$$\mathbb{E}_X\{f(X)\} = \sum_{x \in \mathcal{X}} P_X(x)f(x).$$

 $^{^7\}mathbb{E}_X\{f(X)\}$ denotes the expectation of a function $f(\cdot)$ with respect to the distribution P_X of a set of random variables X and defined as follows:

⁸Let us note that $\widehat{\text{Err}}(\widehat{f}_{LS_{\text{train}}}, LS_{\text{test}})$ estimates the generalisation error conditional on the learning set while other approaches such as *cross-validation* actually estimate the expected generalisation error.

(CV) consists in dividing the available learning set in K disjoint sets and learn in turn on K – 1 folds and estimate the error on the remaining fold. When the number of folds K corresponds to the number of samples, this method is then known as the leave-one-out cross validation.

2.3.3 Other forms of learning

Only one facet of machine learning is considered in this thesis, however many other forms of machine learning have been developed. This section is a brief summary of these other forms of learning.

differs from supervised learning by the absence of UNSUPERVISED LEARNING (labelled) outputs. Since, there are not outputs or targets to supervise the learning process, this part of machine learning focus on extracting informations from data (see, e.g., PCA, ICA, Gaussian mixture models). Gathering similar samples together by making clusters is one way to get some information from unlabelled data. Clustering is one of the most known unsupervised approaches and aims to gather similar samples into *clusters* (see, e.g., *k-means* and *k-metroids*).

SEMI-SUPERVISED LEARNING is halfway between supervised and unsupervised learnings. In this case, some of the samples in the training data are not labelled. Semi-supervised techniques aim at using those additional unlabelled data to better characterise the underlying data distribution than what could be done using only labelled data. Active learning is a particular case in which the learning algorithm can interact with the user in order to improve the quality of the learning process, e.g. by asking for a label.

differs from other kinds of learning by the fact that the TRANSFER LEARNING underlying distribution is not the same in the training data and in the testing data. Therefore, transfer learning mainly consists in learning a model and then apply it on a different but related application.

basically consists in transferring the information re-TRANSDUCTIVE LEARNING trieved from labelled examples to unlabelled ones (see [Bousquet, 2002] for details). The purpose is not to generate a model but only to label unlabelled samples. Transfer transductive learning is a particular case considering transfer learning in a transductive setting [Arnold et al., 2007; Rohrbach et al., 2013]. In this setting, the learning process can use labelled training data but the test set is unlabelled on the target domain (which is different than the training domain as in the transfer learning) but can be seen during training.

REINFORCEMENT LEARNING is apart from previously described forms of learning because it does not only rely on data. Indeed, the goal is not to discover an underlying distribution or mechanism but to determine an optimal control policy (i.e., the strategy that guides (future) chosen actions) from interaction with a system or from observations of a system [Ernst et al., 2005].

2.4 FEATURE SELECTION FOR SUPERVISED LEARNING

Machine learning problems in bioinformatics, neuroimaging, engineering, psychology (and many others) have in common that their typical dimensions have increased very significantly within the last two decades [Guyon and Elisseeff, 2003; Saeys et al., 2007]. Such applications usually go with high-dimensional datasets that are characterised by a large number of input features. Exploring the whole input space in such applications often requires to consider hundreds of thousands of variables. However, many supervised learning techniques were originally designed to cope with only a few tens or hundreds of variables. Furthermore, most practical supervised learning algorithms decrease in performances when facing many features that are not useful for the prediction of the output [Kohavi and John, 1997; Blum and Langley, 1997].

Therefore, reducing the input data dimension, e.g., by selecting a subset of the original features [Liu and Yu, 2005], has become a real prerequisite in such applications. In this context, the task of feature selection mainly consists in finding as small as possible subsets of features that are sufficient to build accurate predictors [Guyon and Elisseeff, 2003]), or alternatively in finding the subset of all informative features, i.e., all those that are somehow related to the output variable [Nilsson et al., 2007; Paja, 2018].

In addition to a dimensionality reduction, feature selection comes along with many potential benefits in terms of interpretability and performances.

Identifying and focusing on (the most) infor-IMPROVING INTERPRETABILITY mative or useful features gives insight of the features involved in the underlying mechanism behind the data and facilitates the data understanding and data visualisation [Guyon and Elisseeff, 2003; Saeys et al., 2007].

Unlike feature extraction or construction techniques (e.g., principal component analysis [Jolliffe, 2011] or partial least squares [Wold et al., 1984]), feature selection preserves original features and thus resulting selected subsets of features remain interpretable by a domain expert [Kohavi and John, 1997; Saeys et al., 2007; Wehenkel, 2018].

The dimensionality reduction helps to overcome **INCREASING PERFORMANCES** the curse of dimensionality and to avoid overfitting [Guyon and Elisseeff, 2003; Saeys et al., 2007]. Smaller data dimensions also reduce storage and computation requirements by providing faster and more cost-effective models [Guyon and Elisseeff, 2003; Saeys et al., 2007]. In presence of many input features that are not necessary for predicting the output, performances of most practical algorithms decrease [Kohavi and John, 1997] and this can be toned down by removing irrelevant features (i.e., not related at all with the output). For example, feature selection often increases the prediction accuracy in supervised learning and often improves the quality of clustering in the case of unsupervised learning [Saeys et al., 2007].

So far, feature selection has been summarized as finding a subset of features. In what follows, we refine this concept by first characterising the relevance of a feature which quantifies the amount of information provided about the target variable. Then we define the usefulness of a feature which is its contribution for a given learning algorithm in prediction accuracy and therefore allows one to define what would be

an optimal subset of features. Then we describe the two flavours of feature selection mentioned in this introduction, namely the *all-relevant* and the *minimal-optimal* problems. While the first problem consists in finding all relevant features in the sense of all features that are somehow related with the output variable, the second problem aims at identifying the smallest subset that yields similar (or better) accuracy performances than any other subset of features.

In the rest of this section, we review some concepts needed for our later developments while abstracting away from the fact that in practice we need to use a finite (and often small) learning set to identify suitable subsets of features for a given problem. We thus use concepts from probability theory and information theory, such as (conditional) independence, Markov boundary, and mutual information to characterize notions such as the relevance and optimality of input features and subsets of input features in the task of predicting the value of a particular output variable.

The notions of Markov boundary and redundancy motivate the fact that all relevant features are not necessary to capture all the information about the target output. Some particular settings that limits the feature selection (or the interpretation that can be retrieved from) will also be reviewed in this chapter such as the multiplicity of Markov boundaries, the difficulty to distinguish direct from indirect effects as well as contextual effects.

2.4.1 Relevance of features

Notational conventions

In the present and subsequent sections we use uppercase letters to denote both individual random variables and sets of random variables, and we reserve lower case letters to denote values of variables or configurations of subsets of variables. In order to lighten the presentation, we assume that all considered random variables are discrete unless explicitly specified differently. We denote the joint probability density of variables X, Y, Z by $P_{X,Y,Z}$ and its value for a combination of values of these variables by $P_{X,Y,Z}(x,y,z)$, and by $P_{X,Y|Z}$ (resp. $P_{X,Y|Z}(x,y|z)$) the conditional joint density of X and Y given Z (respectively its value).

Let us denote by V the set of all original input variables, with |V|=p, and by Y the target output variable. Let V^{-m} be the subset of V excluding the input feature $X_m \in V$ (i.e., $V^{-m}=V\setminus \{X_m\}$).

One facet of feature selection is concerned about the identification in V of the (most) relevant variables. Many definitions of relevance have been proposed in the literature over the years [Gennari et al., 1989; Almuallim and Dietterich, 1991b; Kohavi and John, 1997; Blum and Langley, 1997; Guyon and Elisseeff, 2006] (usually incompatible with each other [Kohavi and John, 1997; Kursa and Rudnicki, 2011]). A common and popular set of relevance notions that we retain has been proposed by Kohavi and John [1997] and is as follows:

Definition 2.3. A variable $X_m \in V$ is **relevant** with respect to the output Y iff there exists a subset $B \subset V^{-m}$ such that $X_m \not\perp Y|B$. A variable is **irrelevant** if it is not relevant.

In this definition the notation "X_m \(\mu \) Y|B" indicates (probabilistic) conditional dependence and is equivalent (in the case of discrete variables) to saying that

$$\exists b, x_{\mathfrak{m}}, y: \qquad \text{such that} \ P_B(b) > 0 \ \text{and} \\ P_{X_{\mathfrak{m}}, Y|B}(x_{\mathfrak{m}}, y|b) \neq P_{X_{\mathfrak{m}}|B}(x_{\mathfrak{m}}|b) P_{Y|B}(y|b).$$

When the subset B is empty, features are relevant by themselves:

Definition 2.4. A variable $X_m \in V$ is marginally relevant with respect to the output $Y iff X_m \not\perp \!\!\! \perp Y.$

Relevant variables can be further divided into two categories:

Definition 2.5. A variable X_m is **strongly relevant** with respect to the output Y iff $Y \not\perp \!\!\!\perp X_m | V^{-m}$.

Definition 2.6. A variable X_m is **weakly relevant** with respect to the output Y if it is relevant but not strongly relevant.

This definition is characterised by two degrees of relevance9 in order to cope with particular settings such as features that are relevant but not marginally (e.g., a XOR problem) [Nilsson et al., 2007]. Strongly relevant variables are thus variables that convey information about the output that no other variable (or combination of variables) in V conveys [Nilsson et al., 2007]. Figure 2.4 is a graphical representation of features in V according to the type of relevance with respect to Y. It shows that the subset of relevant features is made of all weakly relevant features and all strongly relevant ones. Let us note that a system can be constructed so that it contains relevant but no strongly relevant features [Kursa and Rudnicki, 2011].

Alternative, strictly equivalent, definitions of relevance can be formulated using the notion of conditional mutual informations¹⁰ (see [Meyer et al., 2008; Louppe et al., 2013]):

Definition 2.7. A variable $X_{\mathfrak{m}} \in V$ is **relevant** to Y iff there exists a subset $B \subset V$ such that $I(X_m; Y|B) > 0$. A variable is called **irrelevant** if it is not relevant.

Definition 2.8. A variable X_m is strongly relevant to Y iff $I(X_m; Y|V^{-m}) > 0$. A variable X_m is **weakly relevant** if it is relevant but not strongly relevant.

The equivalence between these definitions and Definitions 2.3, 2.5, and 2.6, follows from the equivalence between zero (conditional) mutual information and (conditional) independence¹¹.

⁹Kohavi and John [1997] showed that earlier definitions were not consistent to identify relevance in the case of a Correlated XOR problem (i.e., where the target Y is such that $Y = X_1 \oplus X_2$, where \oplus denotes a logical XOR) with five boolean features X_1, \ldots, X_5 and correlated/redundant features (X_2 and X_4 that are such that $X_4 = \overline{X_2}$ and that two degrees of relevance are required to achieve that. With respect to Y, X_1 is a strongly relevant feature, X_2 and X_4 are weakly relevant features due to their correlation/redundancy and X_3 and X_5 are irrelevant features.

¹⁰See Appendix B, for notations and definitions of several measures from information theory, including the conditional mutual information.

 $^{^{11}}X \perp \!\!\! \perp Y|Z$ and $X \perp \!\!\! \perp Y|Z$ are equivalent to I(X;Y|Z) > 0 and I(X;Y|Z) = 0 respectively [Cover and Thomas, 2012].

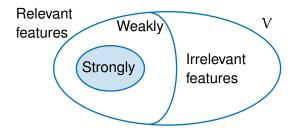


Figure 2.4: Graphical decomposition of the set of input variables V according to the feature relevance. The subset of relevant features can be further refined into two degrees of relevance: weak and strong relevances.

2.4.1.1 On the quantitative measure of irrelevance

In relation to Definition 2.7, several authors (eg., [Bell and Wang, 2000; Guyon and Elisseeff, 2006; Meyer et al., 2008]) proposed to use the notion of (conditional) mutual information to assess the level of relevance/irrelevance of a feature.

For example, Guyon and Elisseeff [2006] define a notion of "approximate irrelevance" as follows:

Definition 2.9. A variable X_m is approximately irrelevant at level ϵ if for all 12 subsets of features $B \subseteq V^{-m}$, $I(X_m; Y|B) \leq \epsilon$.

They further say that a variable $X_{\mathfrak{m}}$ is surely irrelevant if it is approximately **irrelevant at level** $\epsilon = 0$. Notice that this notion is equivalent to the previously introduced notion of irrelevance (Definition 2.7).

Let us finally mention that Guyon and Elisseeff [2006] claim that one single notion of (ir)relevance is enough if one simultaneously considers the notion of sufficient feature subset, while Kohavi and John [1997] preferred two degrees to characterise relevance. In addition to [Kohavi and John, 1997; Guyon and Elisseeff, 2006], we also further refer to [Bell and Wang, 2000] for a review on relevance.

2.4.2 Markov boundary

In this subsection, we introduce the notions of Markov blanket and Markov boundary that will be of interest in the rest of this chapter.

Let us consider a set of features V and a target variable Y, Markov blanket and Markov boundaries are defined as follows [Pearl, 1988; Tsamardinos and Aliferis, 2003; Statnikov et al., 2013]:

Definition 2.10. A *Markov blanket* of variable Y relative to V is a subset $M \subseteq V$ such $Y \perp \!\!\! \perp V \setminus M | M$.

Definition 2.11. A Markov boundary of variable Y relative to V is a Markov blanket of Y relative to V such that no proper subset of M is also a Markov blanket of Y relative to V.

Trivially, the set of all input features V is a Markov blanket of Y and a given Markov blanket can be arbitrarily extended by adding features (even irrelevant ones with respect to Y) [Statnikov et al., 2013]. That is why minimal Markov blankets - Markov

¹²including the empty subset and the set V^{-m} itself.

boundaries - are of greater interest in the context of feature selection¹³ [Margaritis and Thrun, 2000; Tsamardinos and Aliferis, 2003; Aliferis et al., 2003; Hardin et al., 2004; Nilsson et al., 2007; Statnikov et al., 2013]. Figure 2.5 shows how Markov boundaries relate with subsets of relevant features. As shown formally below, any Markov blanket (and hence any Markov boundary) includes all strongly relevant features, and no Markov boundary can contain any irrelevant feature. On the other hand, some weakly relevant features may belong to some Markov boundaries.

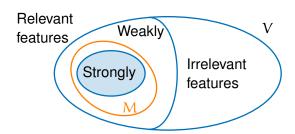


Figure 2.5: Graphical decomposition of the set of input variables V according to the feature relevance. A Markov boundary M in all generality gathers all strongly relevant features and some weakly relevant ones.

A target variable Y may have several Markov boundaries, for example because of redundancies between features [Statnikov and Aliferis, 2010; Geurts and Saeys, 2011; Statnikov et al., 2013]. However, the intersection of all Markov boundaries always includes the set of strongly relevant features.

Indeed we have the following property [Tsamardinos and Aliferis, 2003]:

Property 2.1. Let us consider a set V of input features and an output Y. If M is Markov blanket of Y, and X_m is a strongly relevant feature, then $X_m \in M$. Therefore, any Markov boundary of Y, as well as the intersection of all these Markov boundaries, contains all strongly relevant features.

Proof. Consider some subset M of V which is a Markov blanket of Y; thus

$$Y \perp \!\!\!\perp V \setminus M \mid M.$$
 (2.9)

Then consider some variable $X_m \in V \setminus M$; thus (2.9) may be rewritten as

$$Y \perp \!\!\! \perp (\{X_{\mathfrak{m}}\} \cup (V \setminus (M \cup \{X_{\mathfrak{m}}\}))) | M. \tag{2.10}$$

The weak union property $(X \perp \!\!\! \perp (Y \cup W) | Z \Rightarrow X \perp \!\!\! \perp Y | (Z \cup W)$, see side note on page 38) applied to (2.10) yields

$$Y \perp \!\!\!\perp X_m \mid M \cup (V \setminus (M \cup \{X_m\})), \text{ i.e. } Y \perp \!\!\!\perp X_m \mid V \setminus \{X_m\}.$$
 (2.11)

Therefore X_m is not strongly relevant.

Furthermore, a Markov boundary of Y never contains irrelevant features:

Property 2.2. Let us consider a set V of input features and an output Y. If M is a Markov boundary of Y, and X_i is an irrelevant input feature, then $X_i \notin M$.

¹³In computational biology, Markov boundaries are also known as (molecular) signatures, which are minimal subset of features that are of best interest to predict the value (i.e., the phenotypic response) of a target variables[Statnikov and Aliferis, 2010; Geurts and Saeys, 2011]. In this context, the non-uniqueness of Markov boundaries is known as signature multiplicity. Those two concepts are equivalent as it has been shown that maximally predictive and non-redundant molecular signatures are the Markov boundaries and vice-versa [Statnikov and Aliferis, 2010].

Proof. Consider a Markov blanket M of Y containing an irrelevant variable X_i. Then, rewriting M as $M^{-i} \cup \{X_i\}$, we have

$$Y \perp V \setminus (M^{-i} \cup \{X_i\}) | (M^{-i} \cup \{X_i\}). \tag{2.12}$$

Since X_i is irrelevant with respect to Y, we also have

$$Y \perp \!\!\! \perp X_i | M^{-i}. \tag{2.13}$$

Using the contraction property (i.e., $X \perp\!\!\!\perp Y|Z$ and $X \perp\!\!\!\perp W|(Z \cup Y) \Rightarrow X \perp\!\!\!\perp (Y \cup W)|Z$, see side note on page 38) between Equations 2.13 and 2.12, we thus have

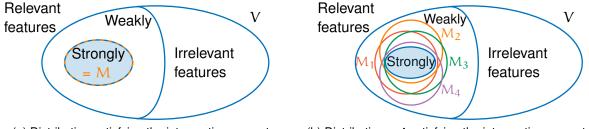
$$Y \perp \{X_i\} \cup (V \setminus (M^{-i} \cup \{X_i\}))|M^{-i}$$
(2.14)

$$\Leftrightarrow \qquad Y \perp \!\!\!\perp V \setminus M^{-i} \mid M^{-i}. \tag{2.15}$$

Equation 2.15 implies that M^{-i} is also a Markov blanket of Y, so that M can not be a Markov boundary of Y.

Following [Nilsson et al., 2007], let us define a strictly positive density Pv over the full set of input variables V as a density such that $P_V(v) > 0$ for all configurations ν of the variables in V. When P_V is strictly positive ¹⁴ (see side note on page 38), the Markov boundary of Y is unique and it contains only strongly relevant features [Tsamardinos and Aliferis, 2003; Hardin et al., 2004; Nilsson et al., 2007; Sutera et al., 2018].

Figure 2.6 illustrates the relation between the concept of Markov boundary and relevance. Figure 2.6a shows that the (unique) Markov boundary coincides with the set of strongly relevant features when the distribution verifies the intersection property (proof in [Nilsson et al., 2007, Theorem 10]). Figure 2.6b illustrates the fact that when the Markov boundary is not unique, the intersection of all Markov boundaries (or blankets) yields the set of strongly relevant features [Tsamardinos and Aliferis, 2003]. The composition property prevents features to be irrelevant for some B but relevant when considered together for the same B.



(a) Distribution satisfying the intersection property

(b) Distribution **not** satisfying the intersection property

Figure 2.6: Correspondance between relevance and Markov boundaries in case of a distribution (a) satisfying the intersection property with a unique Markov boundary M and (b) not satisfying the intersection property with four Markov boundaries M_1, M_2, M_3, M_4 whose the intersection is the set of strongly relevant features.

¹⁴ Equivalently, for any distributions satisfying the intersection property, there is a unique Markov boundary [Pearl, 1988; Statnikov et al., 2013]. Strictly positive distributions always verifies the intersection property [Nilsson et al., 2007]. This also holds for faithful distributions (to some Bayesian network) satisfying the intersection property and being strictly positive [Tsamardinos and Aliferis, 2003; Tsamardinos et al., 2003a; Aliferis et al., 2010; Statnikov et al., 2013].

DISTRIBUTION PROPERTIES

Let X,Y,Z and W be any four subsets of features from V and $T_i \in V$ be a single variable. Any distribution verifies the following properties [Pearl, 1988; Nilsson et al., 2007; Statnikov et al., 2013]:

- Symmetry: $X \perp\!\!\!\perp Y|Z \Leftrightarrow Y \perp\!\!\!\perp X|Z$,
- Decomposition: $X \perp \!\!\! \perp (Y \cup W)|Z \Rightarrow X \perp \!\!\! \perp Y|Z$ and $X \perp \!\!\! \perp W|Z$,
- Weak union: $X \perp \!\!\! \perp (Y \cup W)|Z \Rightarrow X \perp \!\!\! \perp Y|(Z \cup W),$
- Contraction: $X \perp\!\!\!\perp Y \mid Z$ and $X \perp\!\!\!\perp W \mid (Z \cup Y) \Rightarrow X \perp\!\!\!\perp (Y \cup W) \mid Z$,
- Self-conditioning: X ⊥⊥ Z|Z.

Strictly positive distributions (P) also satisfy [Pearl, 1988; Nilsson et al., 2007; Statnikov et al., 2013]:

• Intersection: $X \perp\!\!\!\perp Y | (Z \cup W)$ and $X \perp\!\!\!\perp W | (Z \cup Y) \Rightarrow X \perp\!\!\!\perp (Y \cup W) | Z$.

Nilsson et al. [2007] also consider two additional classes of distributions: strictly positive distributions that satisfy the composition property (PC):

• Composition: $X \perp\!\!\!\perp Y|Z$ and $X \perp\!\!\!\perp W|Z \Rightarrow X \perp\!\!\!\perp (Y \cup W)|Z$,

and strictly positive distributions that satisfy both composition and weak transitivity (PCWT):

• Weak transitivity: $X \perp\!\!\!\perp Y|Z$ and $X \perp\!\!\!\perp Y|R \cup \{T_i\} \Rightarrow X \perp\!\!\!\perp \{T_i\}|Z$ and $\{T_i\} \perp\!\!\!\perp$ Y|Z.

A more restricted class of distributions is strictly positive distributions that are DAG-faithful (PD) (i.e., faithful to some Bayesian network [Tsamardinos and Aliferis, 2003; Statnikov et al., 2013]). PD is included in PCWT [Nilsson et al., 2007] and verifies all its properties. While PD distributions offer some information about the causal structure (i.e., the Markov boundary of feature Y is the set of direct causes, direct effects, and direct causes of direct effects (i.e., spouses) of Y), PCWT (including in particular jointly Gaussian distributions [Studeny, 2006]) is claimed to be more realistic [Nilsson et al., 2007].

2.4.3 Redundancy

In many applications, and in particular in high-dimensional settings, the information about the output Y to predict is shared and sometimes replicated among several input variables. In neuroimaging for instance, one often observes a strong spatial correlation between voxels (i.e., pixels in 3D image) implying that neighbouring voxels are likely to be exchangeable when it comes to predict the output class [Wehenkel et al., 2018]. The fact that the same information about the output is held by several features is called redundancy. It can be total, i.e., several features carry exactly the same information about the output and are exchangeable, or partial, i.e., several features carry some of the same information about the target.

From a feature selection point of view, features that share similar information about the target, such as neighbouring voxels in neuroimaging, are relevant but not necessarily useful together for a learning algorithm. Taking into account redundancy in feature selection may thus help to reduce the number of selected features.

In this section, we first review and refine formal definitions of feature redundancy , propose a quantitative measure of redundancy, and discuss the relation between redundancy and relevance and between redundancy and correlation.

2.4.3.1 Yu and Liu [2004]'s redundancy

Using the concept of Markov blankets, Yu and Liu [2004] define the following notion of redundancy:

Definition 2.12. Let us consider a subset $B \subset V$ of features and a variable $X_i \in V$ $V \setminus B$. We say that X_i is **redundant** to the set B with respect to the target Y iff (i) X_i is weakly relevant with respect to Y and (ii) there exists a subset $M \subseteq B$ such that $X_i \perp \!\!\! \perp (\{Y\} \cup (V \setminus (M \cup \{X_i\}))|M.$

Condition (i) excludes irrelevant features from consideration, since they are anyhow not useful to predict Y. Condition (ii) implies that B contains a Markov blanket M of variable X_i relative to all other features including the target Y. This subset of variables can thus replace X_i without loss of information, both about Y and about any variables from V not in the set B. According to this definition, and as expected, a strongly relevant feature can thus never be redundant to any subset because it conveys information about Y that can not be found in other features and thus condition (ii) can not be satisfied.

This definition was proposed by Yu and Liu [2004] to identify features that can be safely ignored when B is an intermediate approximate solution in the search for a Markov boundary of the target Y. A relaxed definition could have been adopted by changing condition (ii) simply into $Y \perp X_i \mid B$, but this would have excluded features X_i that might bring complementary information about Y with respect to B when combined with some other features from $V \setminus B$.

Total redundancy 2.4.3.2

Louppe [2014, Definition 7.1] defines totally redundant features as pairs of features X_i and X_i such that

$$H(X_i|X_j) = H(X_j|X_i) = 0.$$
 (2.16)

Note that an asymmetrical version 15 of Equation 2.16 has also been proposed to define redundancy (e.g., [Meyer et al., 2008]). One limitation of these definitions is that they do not involve the output variable Y. Therefore, based on [Louppe, 2014, Lemma 7.1], let us define *total redundancy with respect to* Y as follows:

Definition 2.13. X_i and X_j are totally redundant variables with respect to the target Y if for any conditioning set $B \subseteq V^{-i,j} (= V \setminus \{X_i, X_i\})$, we have:

$$Y \perp X_i \mid B \cup \{X_i\}$$
 and $Y \perp X_i \mid B \cup \{X_i\}$ (2.17)

Equation 2.17 states that X_i provides no additional information about the output once X_i is given, whatever the context B, and vice versa. A direct consequence of this definition is that for all $B \subseteq V^{-i,j}$, we have ¹⁶

$$I(X_i; Y|B) = I(X_i; Y|B),$$
 (2.18)

ie., X_i and X_i are equally informative about Y in all circumstances. Total redundancy defines the ability of one feature to replace entirely the other in any context without loss of information about the output. Two totally redundant features are such that one is irrelevant iff the other is irrelevant. Obviously, none of them can be strongly relevant since Equation 2.17 for $B = V^{-i,j}$ gives $Y \perp X_i | V^{-i,j} \cup \{X_i\} \Leftrightarrow Y \perp X_i | V^{-i}$ and also $Y \perp X_i | V^{-j}$.

Note that Equation 2.16, which implies that X_i and X_j are copies of each other, implies Definition 2.13 (see [Louppe, 2014, Lemma 7.1] for a proof and [Meyer et al., 2008, Equations 3.7-3.9] for a proof in the asymmetrical case) but the converse is not true. Two features might be totally redundant with respect to the target, while not explaining perfectly each other. As defined, total redundancy and Yu and Liu [2004]'s redundancy (Definition 2.12) are also different concepts. Given two totally redundant features X_i and X_j , we do not have necessarily that X_i is redundant with respect to the subset $B = \{X_i\}$ according to Definition 2.12. There might indeed exist a distinct feature $X_k \in V$ such that $X_i \not\perp X_k | X_j$ and thus condition (ii) in Definition 2.12 might not be satisfied. It would be always satisfied however if using Louppe [2014]'s definition of total redundancy (Equation 2.16).

Asymmetric and partial redundancies 2.4.3.3

In this section, we propose and discuss two relaxations of the definitions of redundancy given in the two previous sections.

First, while total redundancy as defined in Definition 2.13 is symmetric, one can also define total redundancy in an asymmetric way:

Definition 2.14. X_i is totally redundant to X_i with respect to Y if $\forall B \subseteq V^{-i,j}$, $X_i \perp \!\!\!\perp Y \mid B \cup X_j$.

In other words, X_i is totally redundant to X_i if it never brings any additional information about Y when X_i is known. X_i and X_i are thus totally redundant if they are totally redundant to each other.

Total redundancy means that X_i is always useless for predicting the output when X_i is known. A notion of partial redundancy could also be defined that relaxes this constraint.

 $^{^{15}}X_i$ is defined as redundant with respect to X_j if $H(X_i|X_j)=0$, which does not imply $H(X_j|X_i)=0$ and the redundancy of X_i with respect to X_i .

¹⁶This is an immediate consequence of Equations 2.23-2.24 and the fact that $Y \perp X_i \mid B \cup \{X_i\} \Rightarrow$ $I(X_i;Y|B \cup \{X_i\}) = 0 \text{ and } Y \perp \!\!\!\perp X_i|B \cup \{X_i\} \Rightarrow I(X_i;Y|B \cup \{X_i\}) = 0.$

Definition 2.15. X_i is partially redundant to X_j with respect to Y if (i) $\exists B \subseteq V^{-i,j}$ such that $X_i \not\perp \!\!\! \perp Y | B \cup X_j$ and (ii) $\forall B \subseteq V^{-i,j}$ such that $X_i \not\perp \!\!\! \perp Y | B \cup X_j$:

$$I(X_i; Y|B) > I(X_i; Y|B \cup \{X_i\}).$$
 (2.19)

Condition (i) excludes X_i from being totally redundant to X_i . Condition (ii) means that the information that X_i brings about the output is always reduced when X_j is known. Having instead $I(X_i;Y|B) < I(X_i;Y|B \cup \{X_j\})$ would mean that X_i is more complementary than redundant to X_i . Note that the equality is impossible since $X_i \not\perp \!\!\! \perp Y | B \cup X_j$ implies that $I(X_i; Y | B \cup \{X_j\}) > 0$.

Interestingly, Definition 2.15 implies that X_i and X_j are both relevant to Y.

Property 2.3. If X_i is partially redundant to X_i with respect to Y, then X_i and X_i are both relevant with respect to the output Y.

Proof. By definition of partial redundance, there exists at least one B such that $X_i \not\perp \!\!\! \perp Y \mid B \cup X_i$. For one such B, condition (ii) implies that:

$$I(X_i; Y|B) > I(X_i; Y|B \cup \{X_i\}) > 0.$$
 (2.20)

From Equation 2.20, we directly have that

$$I(X_i; Y|B) > 0$$
 (2.21)

implying that X_i is relevant with respect to Y.

Then, the first inequality of Equation 2.20 is equivalent to

$$I(X_i; Y|B) - I(X_i; Y|B \cup \{X_i\}) > 0.$$
 (2.22)

The chain rule $(I(X_1, X_2, ..., X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, ..., X_1))$ applied to the mutual information between both features X_i , X_i and Y yields

$$I(X_i, X_j; Y|B) = I(X_i; Y|B) + I(X_j; Y|B \cup \{X_i\})$$
 (2.23)

$$= I(X_i; Y|B) + I(X_i; Y|B \cup \{X_i\})$$
 (2.24)

where Equations 2.23 and 2.24 depend on the order in which X_i and X_j are used. By rearranging terms in 2.23 and 2.24, we have

$$I(X_i; Y|B) - I(X_i; Y|B \cup \{X_i\}) = I(X_i; Y|B) - I(X_i; Y|B \cup \{X_i\}).$$
 (2.25)

Since the left member is strictly positive given Equation 2.22, we thus have

$$I(X_i; Y|B) - I(X_i; Y|B \cup \{X_i\}) > 0$$
 (2.26)

which implies that $I(X_i; Y|B) > 0$ because $I(X_i; Y|B \cup \{X_i\}) \ge 0$ (positivity of conditional mutual information) and $I(X_i; Y|B) > I(X_i; Y|B \cup \{X_i\})$. Therefore X_i is also relevant with respect to Y.

The proof of the previous theorem shows that Equation 2.19 is equivalent to Equation 2.26. In consequence, if X_i reduces the information brought by X_i about Y, then X_i also reduces the information brought by X_i about Y. Nevertheless, partial redundancy is not symmetric because the sets B such that $X_i\not\perp\!\!\!\perp Y|B\cup X_i$ do not necessarily coincide with the sets B such that $X_i \not\perp Y \mid B \cup X_i$.

2.4.3.4 Quantitative measure of redundancy

A measure of redundancy among p random variables X_1, \ldots, X_p can be defined as follows (see, e.g., [McGill, 1954; Watanabe, 1960; Wienholt and Sendhoff, 1996; Jakulin and Bratko, 2003b; Meyer et al., 2008]):

$$R(X_1; X_2; ...; X_p) = \sum_{i=1}^{p} H(X_i) - H(X_1, X_2, ..., X_p)$$
 (2.27)

where $H(X_i)$ and $H(X_1, X_2, ..., X_p)$ are respectively the entropy of X_i and the joint entropy of X_1, X_2, \dots, X_p (see Appendix B). However, like total redundancy (Definition 2.13), this measure does not involve the output variable

Y [Meyer et al., 2008]. Therefore, following the use of $I(X_m; Y|B)$ to quantify feature relevance (see Section 2.4.1.1), one could use similarly multivariate mutual information [McGill, 1954] to quantify redundancy.

Multivariate mutual information is usually defined as follows:

$$I(X;Y;Z) = I(X;Y,Z) - I(X;Z|Y) - I(X;Y|Z)$$
 (2.28)

It can be shown that I(X;Y;Z) is symmetric with respect to a permutation of the roles of X, Y, Z (e.g., I(X;Y;Z) = I(X;Z;Y)) and, applying the chain rule on I(X;Y,Z), that

$$I(X;Y;Z) = I(X;Y) - I(X;Y|Z).$$
 (2.29)

Unlike standard (conditional) mutual information, I(X;Y;Z) can be negative as I(X;Y)can be increased by conditioning on Z. McGill [1954] sees I(X; Y; Z) (Equation 2.29) as the the gain (or loss) of common information between two variables (i.e., X and Y) due to the additional knowledge of a third one (i.e., Z). A negative value is therefore due to an increase of the dependence between X and Y knowing Z. Noting the symmetry, I(X;Y;Z) (Equation 2.28) can also be seen intuitively as a generalisation of the mutual information common to three random variables [Cover and Thomas, 2012].

The degree of redundancy between two features (in a given context B) could then be defined as follows:

Definition 2.16. For a given conditioning set $B \subseteq V^{-i,j}$, the **degree of redundancy** between X_i and X_j with respect to Y is measured by

$$I(X_i; X_j; Y|B) = I(X_i; Y|B) - I(X_i; Y|B \cup \{X_i\}).$$
(2.30)

 $I(X_i; X_i; Y|B)$ has several desirable properties as a measure of the degree of redundancy:

- It is positive as soon as $I(X_i; Y|B) > I(X_i; Y|B \cup \{X_i\})$ or equivalently $I(X_i; Y|B) > I(X_i; Y|B)$ $I(X_i; Y|B \cup \{X_i\})$, which corresponds precisely to condition (ii) of partial redundancy (Definition 2.15).
- It is equal to zero when $I(X_i; Y|B) = I(X_i; Y|B \cup \{X_i\})$, which corresponds to X_i not impacting the information brought by X_i about the output.
- It is negative when X_i and X_j are complementary. For instance, in the case of a XOR problem, X_i and X_i are marginally irrelevant but together perfectly explain the output Y. Mathematically, we have in this case $I(X_i; Y) = I(X_i; Y) = 0$ and $I(X_i; Y|X_i) = I(X_i; Y|X_i) = H(Y)$, which is strictly greater than 0 unless Y is constant. Therefore, $I(X_i; Y|X_i) > I(X_i; Y)$ and thus $I(X_i; X_i; Y) < 0$.

• It is maximal and equal to $I(X_i; Y|B) = I(X_i; Y|B)$ when X_i and X_i are totally redundant, as in this case $I(X_i; Y|B \cup \{X_i\}) = I(X_i; Y|B \cup \{X_i\}) = 0$.

Note that several authors have proposed to use the opposite of Equation 2.16 to quantify the synergy or the complementarity between two features, which is indeed the opposite of redundancy. This measure can also be generalised to more than two features. See, e.g., [Meyer and Bontempi, 2013] for a review of these measures.

2.4.3.5 Redundancy and relevance

Like relevance, redundancy characterises the interest of (de)selecting features. Depending on how the feature selection problem is formulated (see Section 2.4.4), it is often desirable not to select totally redundant features that convey the exact same information about the output as other features. By definition, strongly relevant features always contain some unique information and thus only weakly relevant features can be considered as (totally) redundant with respect to some other features. Figure 2.7 (adapted from Yu and Liu [2004]) illustrates that input features can be divided into four categories: irrelevant, strongly relevant, non-redundant and redundant weakly relevant features. Non-redundant and redundant features are such that the redundant ones are redundant to both the non-redundant ones and the strongly relevant features with respect to the target (according for example to Definition 2.14 extended to sets of features). Since redundancy is a relative notion that is defined for pairs of features (or sets of features), the division of the weakly relevant features is typically not unique. For instance, if two copies of the same (relevant) feature are present, each one of them could play the role of the redundant one to the other leading to at least two divisions.

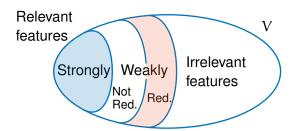


Figure 2.7: Graphical decomposition of the set of input variables V according to the feature relevance. The subset of relevant features can be refined into two degrees of relevance: weak and strong relevance. Weakly relevant features can furthermore be divided into completely redundant (with respect to non-redundant features) and non-redundant features.

2.4.3.6 Redundancy and correlation

Correlation is a statistical measure of the dependence between two numerical random variables. The most common measure of correlation is the Pearson correlation coefficient defined for two random variables A and B as [Pearson, 1896; Lee Rodgers and Nicewander, 1988; Guyon and Elisseeff, 2006]:

$$\rho(A,B) = \frac{cov(A,B)}{\sigma_A \sigma_B} \tag{2.31} \label{eq:2.31}$$

where $cov(A, B) = \mathbb{E}\{(A - \mu_A)(B - \mu_B)\}$ is the covariance between both variables and where μ and σ denote respectively the mean and the standard deviation. When the values of both variables move in the same direction (resp. opposition direction) in a similar fashion (i.e., by keeping a fixed distance), they are perfectly correlated (resp. anti-correlated) and this corresponds to $\rho = 1$ (resp. $\rho = -1$).

Correlation and redundancy are different notions. We saw that duplicated (relevant) features are subsequently totally redundant with respect to the target. Intuitively, one may expect that a high correlation (or anti-correlation) between the values of two features suggests that those features are also redundant. However, correlation does not imply redundancy [Guyon and Elisseeff, 2006]. Figure 2.8 gives examples (inspired from [Guyon and Elisseeff, 2006]) showing that highly correlated features are not necessary redundant. But, if $cov(X_i, X_j) = \pm 1$ then Equation 2.16 holds and thus X_i are totally redundant with respect to any target Y.

2.4.4 Feature selection problems

Besides the objective of size reduction, the problem of feature selection usually can take two flavours [Guyon and Elisseeff, 2003; Nilsson et al., 2007; Genuer et al., 2010; Kursa and Rudnicki, 2011]. Typically, those side-objectives guide the feature selection and determine the subset of features that end up being selected.

Many studies (e.g., with microarray gene-expression data [Ambroise and McLachlan, 2002] or in drug discovery application [Janecek et al., 2008]) showed that dealing with small sets of relevant features usually gives better results and facilitate learning accurate classifiers [Guyon and Elisseeff, 2003; Nilsson et al., 2007; Kursa and Rudnicki, 2011].

In presence of many features, it is common that a large number of features are either irrelevant or redundant (to the target). Such variables are in principle not necessary to predict the output and computational performances of supervised learning algorithms can often be optimised by discarding them [Yu and Liu, 2004]. Discarding some non-redundant features (with respect to those that are kept), may however be detrimental in terms of accuracy. Furthermore, for some specific learning algorithms, it may actually be beneficial in terms of accuracy to keep some redundant features. Moreover, when sample sizes are small compared to the number of features, it may even become beneficial (in terms of accuracy), to discard some non-redundant features (to decrease overfitting). Usually, as the number of selected features grows, it is expected that the performances of a learning algorithm increases and then decreases. The optimal size for the feature subset being the one that maximises the accuracy [Hua et al., 2004], the *minimal-optimal problem* is the first problem of feature selection and consists in finding the smallest optimal subset for a given learning algorithm and a given dataset.

When only accuracy of the learnt predictor is used a criterion to select an optimal subset of features, many weakly relevant features (and sometimes even some strongly relevant one) might be discarded. There is however an interest of identifying all features that are somehow related to the target in order to get a full understanding of the underlying mechanism (e.g., in gene expression analysis [Golub et al., 1999]). The *all-relevant* problem is the second approach of feature selection and consists in finding all relevant features.

Those two approaches are usually complementary for a given application. Let us take the example of a medical diagnosis that consists in predicting a disease. The doctor has to evaluate a given number of factors before making his diagnosis. The number of factors has to be as a small as possible to save time and money. Hence,

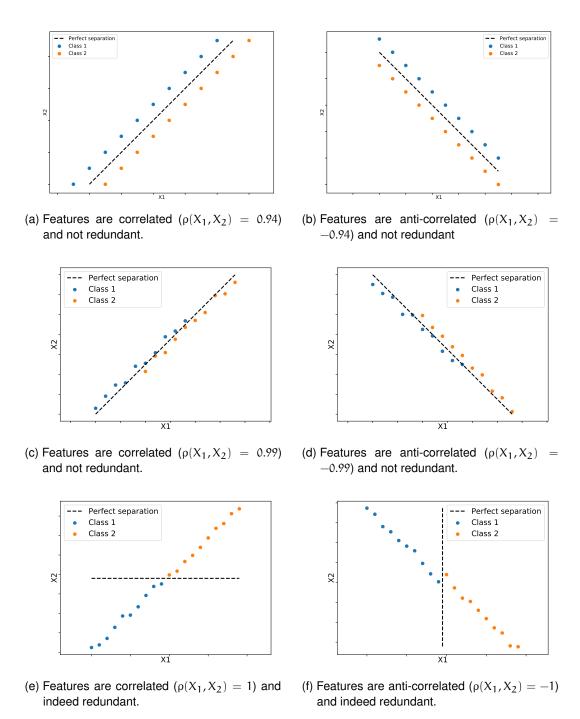


Figure 2.8: Illustrating examples where correlation does not necessary imply redundancy. Figures (a) to (d) show that features can be highly correlated while being not redundant as both features are required to achieve a perfect separation between the two classes. Figures (e) and (f) show that correlated features can indeed be redundant as one feature out of the two is enough to perfectly separate classes. Let us note that in both last examples, both features can individually lead to a perfect separation.

one would want to identify a small set of features that provides the best possible diagnosis. The minimal-optimal approach aims at providing such a feature subset. In different circumstances, for research purposes for instance, the all-relevant approach may be more appropriate. One may want to identify all factors that are related to the output even if some of them are redundant with respect to other.

Both feature selection approaches are further described below.

ALL-RELEVANT PROBLEM The all-relevant problem is defined as follows [Nilsson et al., 2007; Kursa and Rudnicki, 2011]:

Definition 2.17. The all-relevant feature selection problem consists in finding all relevant features. The solution to this problem is the set of all strongly and weakly relevant features.

The solution of this problem is in principle unique, as suggested by Figure 2.9. One further step, in such an analysis, would be to also distinguish between strongly and weakly relevant features.

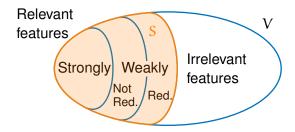


Figure 2.9: The solution S to the all-relevant problem is the union of weakly relevant and strongly relevant features to the target variable. This set includes all relevant features even if there is redundant information about the target.

MINIMAL-OPTIMAL PROBLEM In terms of learning algorithm performances, the minimal-optimal problem is usually defined as follows [Kohavi and John, 1997; Nilsson et al., 2007]:

Definition 2.18. Let $\mathcal A$ be a learning algorithm, V the set of input features and Y be the target feature. The minimal-optimal feature selection problem consists in finding a subset of V of minimal size that minimises the generalisation error of $\mathcal A$.

A solution of this problem is usually a subset of all relevant features, even if for some very specific combinations of problems and algorithms, including irrelevant features may actually be beneficial from the viewpoint of accuracy [Kohavi and John, 1997].

For regression and *calibrated*¹⁷ classification tasks, Tsamardinos and Aliferis [2003] showed that a Markov boundary of minimal size is a solution to the minimal-optimal problem (see [Tsamardinos and Aliferis, 2003, Proposition 3] for more details and see side note on page 47 for a word on Markov blanket discovery algorithms). Therefore, a solution of the minimal-optimal problem is a set made of all

¹⁷A classification problem which requires the exact distribution of predictions of Y and not only the most probable class of Y is said to be calibrated [Tsamardinos and Aliferis, 2003]. Such problems correspond for instance to classification problems where the mean squared loss is used instead of the zero-one loss.

strongly relevant and a maximal subset of non-redundant¹⁸ weakly relevant features [Kursa and Rudnicki, 2011]. Let us however note that when the zero-one loss is used (i.e., only the most probable class of Y is required), Tsamardinos and Aliferis [2003] state that only some features of the Markov boundary are required or features that do not belong to the Markov boundary.



MARKOV BLANKET DISCOVERY ALGORITHMS

Markov blanket and boundary discovery algorithms constitute another broad family of feature selection techniques (see, e.g., [Guyon and Elisseeff, 2003; Tsamardinos et al., 2003a; Aliferis et al., 2010; Statnikov et al., 2013; Tsamardinos et al., 2003b]). They are usually independent of any learning algorithm and are mainly based on graph theory and related to causality. They are however not addressed in this thesis.

Resulting of the multiplicity of Markov boundaries, the minimal-problem problem does not have a unique solution in general. For example, in presence of two totally redundant features, each can be kept (without the other one) giving two valid options. For strictly positive distributions however, the Markov boundary M of Y is unique and corresponds to the set of all strongly relevant variables [Nilsson et al., 2007].

Figure 2.10 shows typical solutions to the minimal-optimal problem with respect to the relevance of features. In the case of a strictly positive distribution, Figure 2.10a gives the unique solution to the minimal-optimal problem which is the set of strongly relevant features. In the case of non-strictly positive distribution, Figure 2.10b illustrates a solution to the minimal-optimal problem which includes in all generality some weakly relevant features and all strongly relevant ones.

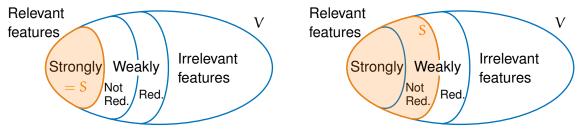
Finding an optimal subset is usually intractable because some distributions may require an exhaustive search of all possible subsets to guarantee optimality [Cover and Van Campenhout, 1977; Kohavi and John, 1997; Blum and Langley, 1997; Yu and Liu, 2004; Nilsson et al., 2007]. With p features, there are 2^p possible subsets which is clearly impractical, especially for high-dimensional datasets. However, letting this search be guided by a heuristic (see Section 2.4.6 and [Guyon and Elisseeff, 2003]) or considering only strictly positive distributions [Nilsson et al., 2007] make this problem more tractable computationally.

2.4.5 Feature selection methods

Feature selection methods are usually classified in three categories depending on how they interact with the learning algorithm: filters, wrappers and embedded methods [Blum and Langley, 1997; Guyon and Elisseeff, 2003; Tsamardinos and Aliferis, 2003; Saeys et al., 2007].

FILTERS A filter approach aims at selecting features independently of the learning algorithm (i.e., without optimising its performance) [Kohavi and John, 1997; Blum and Langley, 1997; Guyon and Elisseeff, 2003, 2006; Saeys et al., 2007; Brown et al., 2012; Chandrashekar and Sahin, 2014]. A filter tries to assess the interest

¹⁸Features providing non-redundant information about the output but that are redundant with some of non-selected features. In other words, relevant features that are included in a Markov boundary but that are not strongly relevant with respect to the output.



- (a) Distribution satisfying the intersection property
- (b) Distribution not satisfying the intersection property

Figure 2.10: Typical solutions S to the minimal-optimal problem for distributions satisfying the intersection property or not. Solutions are Markov boundaries of Y with respect to V.

of keeping features solely from the data in order to then filter out irrelevant features. As a pre-processing step that selects inputs, any learning algorithm can thus be combined with a filtering feature selection.

A common filter method is feature ranking¹⁹ [Stoppiglia et al., 2003b; Blum and Langley, 1997; Guyon and Elisseeff, 2003; Chandrashekar and Sahin, 2014] and consists in ordering features according to a suitable ranking criterion. Any feature relevance measure providing a numerical score can be used (e.g., correlation [Guyon and Elisseeff, 2003] or mutual information [Blum and Langley, 1997; Brown et al., 2012] with the target, decision tree²⁰ [Cardie, 1993], ...). Then the top k features (i.e., with highest value) are selected [Blum and Langley, 1997; Saeys et al., 2007; Chandrashekar and Sahin, 2014]. The number k of selected features is determined using an (arbitrary or not) threshold value or based on a *random probe* (i.e., a random variable is introduced in the process in order to determine which features are statistically better than an artificial irrelevant feature and thus relevant) [Stoppiglia et al., 2003b].

Filter techniques are usually computationally fast and scale very well to high-dimensional datasets [Saeys et al., 2007]. As they are independent of the learning algorithm, they only need to be performed once and for all whatever what follows.

A downside of this independence is that filter techniques totally ignore the performance of the learning algorithm with the selected subset [Kohavi and John, 1997]. In the filtering approach, most proposed techniques (e.g., correlation and mutual information) are univariate: each feature is considered individually and feature dependencies and redundancies are not taken into account [Guyon and Elisseeff, 2003; Saeys et al., 2007; Chandrashekar and Sahin, 2014]. Ignoring such effects may lead to a selected set of features that yields poor performance when compared to other types of (multivariate) feature selection techniques [Saeys et al., 2007]. Besides, a subset of the selected set of features may be sufficient in presence of redundancy [Chandrashekar and Sahin, 2014]. Consequently, multivariate criteria have been proposed to integrate feature dependencies (e.g., based on mutual information [Peng et al., 2005; Meyer et al., 2008; Frénay et al., 2013; Meyer and Bontempi, 2013] or based on Markov blankets [Koller and Sahami, 1996], and see [Brown et al., 2012] for a unifying framework based on conditional likelihood maximisation) but at the cost of scalability and computational speed [Saeys et al., 2007].

¹⁹Feature ranking is sometimes referred to as feature weight based approach as a weight is assigned to each feature [Blum and Langley, 1997; Kira and Rendell, 1992a].

²⁰Feature selection using tree-based models will be discussed in Chapter 4.

In the light of the feature selection problems introduced in Section 2.4.4, two filter approaches are of interest:

THE FOCUS ALGORITHM [Almuallim and Dietterich, 1991a,b, 1994] conducts an exhaustive search among all possible subsets of features for the minimal one providing a perfect discrimination (or the best possible) of the target values [Koller and Sahami, 1996; Kohavi and John, 1997]. It has a preference for a small set of features and suffer from the so-called Min-features bias [Almuallim and Dietterich, 1991a] which may lead to a poor feature selection (see side note on page 49). Nevertheless, it is expected that the selected set of features includes all strongly relevant features and some weakly relevant ones.



MIN-FEATURES BIAS

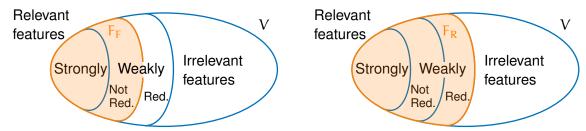
By chance, an irrelevant feature could be sufficient to perfectly determine the target value in the training data (e.g., a unique sample ID such as the social security number in a medical dataset). A learning algorithm receiving such a feature would surely overfit the training data leading to poor performances in generalisation [Kohavi and John, 1997]. A preference towards small set of features - the Min-features bias - would choose that variable as the best subset in comparison with other subsets made of a single variable.

THE RELIEF ALGORITHM [Kira and Rendell, 1992a,b] is an instance of feature ranking and aims at assigning a relevance score to each feature²¹. The selection is then made by considering as relevant (and thus to be kept) features with a relevance score above a given threshold (determined for instance by a statistical method of interval estimation). The selected subset of features is expected to be the set of all relevant features (weak and strong ones) including redundant features.

Figure 2.11 illustrates a typical solution according to the relevance for both algorithms. One can see that Focus algorithm aims to solve the minimal-optimal problem (although ignoring the usefulness of the selected set of features) and that Relief algorithm aims to solve the all-relevant problem [Kohavi and John, 1997].

WRAPPERS A wrapper method aims at selecting a set of features using a learning algorithm as a "black box" [Kohavi and John, 1997; Blum and Langley, 1997; Guyon and Elisseeff, 2003; Saeys et al., 2007]. A set of features is presented to a learning algorithm and the corresponding accuracy performances is used as an estimation of the relative usefulness of the given set of features. The search over all possible subsets is usually guided by a search algorithm (see Section 2.4.6 for more details) "wrapped" around the learning algorithm [Saeys et al., 2007]. At the end, the best set of features is then selected as the one leading to the best performances of the given learning algorithm.

²¹In Kira and Rendell [1992a], the relevance level of the j^{th} feature of the i^{th} sample, denoted x_i^i , is based on two distances: (i) the difference c_i^t between values of x_i^t and x_i^c where x_i^c is the value of the same feature for a sample s which is the closest one with the same class (i.e., $y^j = y^c$) as sample i; (ii) the difference d_i^i between values of x_i^i and x_i^d where x_i^d is the value of the same feature for a sample d which is the closest one with a different class (i.e., $y^j \neq y^d$). The relevance level of a feature X_i is based on an average over all samples of the square of those two distances.



(a) Focus (F_F) aims to solve the minimal-optimal problem (b) Relief (F_R) aims to solve the all-relevant problem

Figure 2.11: Expected set of features selected by Focus (F_F) and Relief (F_R) algorithm according to the feature relevance.

In the wrapper approach, the optimal feature subset search is carried out in interaction with a specific learning algorithm \mathcal{A} . The resulting selected set of features is therefore the most useful for \mathcal{A} but also tailored to it. Wrapper methods benefit from learning algorithm characteristics (e.g., feature dependencies) but depend on its complexity implying a high computational cost. In contrast with filter techniques, the feature selection is coupled with the learning algorithm performances increasing the risk of overfitting. Examples of wrapper methods (e.g., sequential feature selection and sequential backward elimination) are given in Section 2.4.6.

EMBEDDED METHODS An embedded method of feature selection is comprised in the learning algorithm [Blum and Langley, 1997; Guyon and Elisseeff, 2003; Geurts et al., 2006; Saeys et al., 2007; Chandrashekar and Sahin, 2014]. Similarly to wrapper methods, the selected set of features is specific to the learning algorithm. However, the feature subset search and evaluation are incorporated in the training algorithm [Guyon and Elisseeff, 2006] and thus embedded methods are usually less computationally expensive than wrapper methods [Saeys et al., 2007; Chandrashekar and Sahin, 2014]. Two examples of embedded methods are two regularised linear regressions known as Lasso and Ridge regressions. Both methods construct a linear model that minimises its error for a given loss function and uses a subset of the variables while including a penalty term that limits the number of variables used. Similarly with wrappers, the selected set of features may depend on the considered embedded method. Indeed, Lasso and Ridge regressions use different penalisation terms and therefore may select different set of features.

2.4.6 Feature subset search algorithms

Given p features, the number of possible feature subsets (i.e., equal to 2^p) grows exponentially with the number of features making the feature selection space (i.e., the space of all possible subsets of features) very large. Several approaches have been proposed to explore this space. An exhaustive search is optimal but computationally intensive. Heuristic searches have been introduced to explore this space more efficiently. In the rest of this section, we describe well-known search algorithms that will be of interest in this thesis.

EXHAUSTIVE SEARCH [Kira and Rendell, 1992a] consists in exploring the whole feature selection space. All possible subsets are evaluated and the smallest one that maximises a given criterion (which can be a relevance index for a filter approach or the accuracy for a wrapper method for example) is selected.

The optimal subset is thus always found at the expense of computational efficiency. For example, the Focus algorithm [Almuallim and Dietterich, 1991a,b, 1994] examines subsets by increasing order of size and stops as soon as an optimal subset is found. This approach limits the computational burden while preserving optimality [Kira and Rendell, 1992a].

HEURISTIC SEARCH explores more efficiently the search space while trying to find the best (possible) subsets of features. Several approaches aim at reducing the number of subsets to evaluate. A first way consists in limiting the maximal size to $d \leq p$. Only subsets with d or less features are considered but this requires an explicit value of d which is in practice unknown [Kira and Rendell, 1992a; Devijver and Kittler, 1982].

The Sequential Feature Selection (SFS, also known as Forward Selection) [Whitney, 1971; Miller, 1990; Kira and Rendell, 1992a; Blum and Langley, 1997; Jain et al., 2000; Guyon and Elisseeff, 2003; Reunanen, 2003; Chandrashekar and Sahin, 2014] starts by selecting the single feature that maximises the given criterion and then sequentially adds one feature at a time. At each step, each remaining feature is evaluated in combination with already selected features and the best one is permanently added to the current subset. The process stops when all features have been added or when the required size of subset is reached. Conversely, the Sequential Backward Elimination (SBE, also known as *Sequential Backward Selection*) [Marill and Green, 1963; Kira and Rendell, 1992a; Pudil et al., 1994; Kohavi and John, 1997; Blum and Langley, 1997; Jain et al., 2000; Guyon and Elisseeff, 2003; Chandrashekar and Sahin, 2014] starts with all features and then evaluates shrinking feature sets. At each step, the less promising feature (i.e., the one whose removal is the less penalising according to the criterion) is removed, one at a time, until the required subset size is reached.

SFS is more computationally advantageous than SBE as first evaluated subsets are made of few features [Kohavi and John, 1997] (see side node on page 52). Feature dependency is not taken into account as some features may not be very useful individually while being highly informative together [Kira and Rendell, 1992a; Chandrashekar and Sahin, 2014]. However, the backward elimination strategy can theoretically capture feature interactions [Kohavi and John, 1997]. Both approaches do not examine all possible subsets and yield nested feature subsets in the sense that a selected (respectively removed) feature can not be removed (respectively re-selected) even if it would lead to a better subset of features [Pudil et al., 1994; Guyon and Elisseeff, 2003; Reunanen, 2003; Chandrashekar and Sahin, 2014]. Therefore, optimality of the selected subset can not be guaranteed [Pudil et al., 1994; Jain and Zongker, 1997]. More complexed algorithms have been proposed in order to overcome nested subsets. The approach "Plus-1-Minus-r" consists in combining the forward selection and backward elimination in selecting at each step the 1 most promising features and removing the r less promising ones [Stearns, 1976; Kittler, 1978]. Parameters 1 and r need however to be fixed. Sequential Floating Forward Selection (SFFS) follows the sequential search procedure but includes a potential feature elimination at each step [Pudil et al., 1994; Somol et al., 1999]. Similarly, Sequential Floating Backward Elimination (SFBE) includes a potential feature selection at each step [Pudil et al., 1994; Somol

et al., 1999]. Adaptive Sequential Forward Floating Selection (ASFFS) generalises above-mentioned approaches with an adaptive determination of l and r at each step [Somol et al., 1999].

A WORD ON THE COMPLEXITY OF SFS AND SBE

Let us consider a dataset D made of a set F of p features and N samples. We want to solve the minimal-optimal problem. Thus, we need to identify the best feature subset of F according to a function J(E) that evaluates a feature subset $E \subseteq F$ on the dataset **D**. The evaluating function J can either be an independent criterion (in a filtering approach) or an induced model (in a wrapper approach). In both cases, J returns a score that assess the quality of the selected subset E and has a computational cost $\mathcal{O}(J)$ (e.g., that may be the computation cost of the model). In the case of an exhaustive search, all 2^p subsets must be examined in order to find the best one. The overall complexity is therefore $\mathcal{O}(2^p)\mathcal{O}(J)$ but guarantees optimality. This can be conceivably performed, if p is not too large [Guyon and Elisseeff, 2003] but otherwise it is computationally intractable. In the case of a heuristic search^a (either SFS or SBE), the complexity^b is $\mathcal{O}(p^2)\mathcal{O}(J)$, which is much less than $\mathcal{O}(2^p)\mathcal{O}(J)$. Both are then much more efficient than exhaustive search but may not yield optimal results as the procedure may miss some feature interactions. It should be stressed that the evaluating function cost may overburden the overall search complexity and therefore it is important to have an efficient and reliable evaluation of each subset. In Chapter 7, we propose to use a computationally inexpensive model (i.e., a randomised tree, see Chapter 3 for a definition) to perform a multivariate sequential feature selection. Let us also mention that Nilsson et al. [2007] showed that if the distribution is restricted to be strictly positive (i.e., all weakly relevant variables are necessarily redundant and thus can be ignored): the minimal-optimal problem can be solved in polynomial time in the number of features and SBE approaches become consistent (but SFS ones do not).

^aLet us notice that complexities of those sequential search algorithms given in [Kira and Rendell, 1992a] rather correspond to approaches that exhaustively consider all subsets of sizes lower than (resp. greater or equal) d

p for the sequential forward (respectively backward) selection.

^bNote that the size of the selected feature set can be fixed $(d \leq p)$ to stop earlier these sequential searches, but this does not change the complexity.

2.4.7 Discussion

This section aims at reviewing some of the main limitations of feature selection and open problems that motivate some of the research questions considered in the rest of this thesis.

FEATURE RELEVANCE IN THE CONTEXT OF OTHERS Multivariate approaches are usually preferred over univariate ones because they take into account feature dependencies even though they are computationally less efficient. It shows that feature dependencies is crucial in many applications. Relevant features (even strongly) can be marginally irrelevant while being (highly) relevant in combination with other

features [Domingos, 1996; Guyon and Elisseeff, 2003]. A well-known example is the exclusive-OR (XOR) structure [Guyon and Elisseeff, 2003; Kohavi and John, 1997]. Redundancy (another form of feature dependency) may tone relevance or usefulness of features down. Consequently a feature may be not selected (or identified as relevant) while being highly marginally relevant.

The problem of finding all relevant features requires thus to carefully take into account feature dependencies and the only way to do so is to perform an exhaustive search [Nilsson et al., 2007], especially to identify weakly relevant features. For example, a sequential forward selection would systematically fail in the identification of relevant features structured such as cliques, i.e. all features are relevant together but are irrelevant in any subset of the clique. Indeed, let us for instance consider a clique made of two features, i.e. an XOR structure. SFS evaluates the relevance of features in the context of already selected ones. In our example, if both features are marginally irrelevant, then none will be selected preventing also the identification of the other feature. SFS is thus unable to identify features that are only relevant in the context of non-selected ones. Nevertheless, features that make other features relevant need to be relevant as well and thus may be end up being selected [Sutera et al., 2018]. However, if such structures are excluded (e.g., by considering only PCWT distributions), the exhaustive search is not required any more and the allrelevant problem can be solved efficiently [Nilsson et al., 2007]. Complex feature structures such as the clique are studied in Chapter 7 in the context of tree-based feature selection.

Last but not least, the confounding effect is an indirect feature interaction. One input feature may seem irrelevant to the target but another feature, an external feature known as a confounding factor, provides the key to understand the relationship between the input feature and the output. This confounding effect can be enlarged to features that appear at first sight to be irrelevant but, taking into account the context, are indeed relevant. Such feature interactions are studied in Chapter 6.

FEATURE RANKING IS LIMITED FOR INTERPRETATION Feature ranking is extremely limited as it only provides a single ordering of features. This ranking can not render the full complexity of feature interactions or the multiplicity of optimal subsets of features. The subset evaluation function is also critical, e.g. a univariate criterion will only rank features according to their marginal relevance missing potential interactions.

In all generality, the most relevant features are not necessary the best ones (or the only ones) to select [Guyon and Elisseeff, 2003]. The top-ranked feature may be a rather good feature to predict the target but some other features with lower ranks may perfectly discriminate the target together. Redundancy may have lowered the rank of redundant but highly relevant features [Guyon and Elisseeff, 2003]. Selecting a top-ranked feature may also be counter-productive, e.g. selecting only one feature of a clique is not interesting without all the rest of the clique (which might typically be much lower in the ranking). Although very useful, feature selection/ranking methods however only provide very limited information about the often very complex input-output relationships that can be modelled by supervised learning methods. There is no information about feature dependencies in a classical feature ranking. In case of a contextual effect, two similar ranked features may have totally different roles. One may be always relevant while the relevance of the other one depends on the

context. The interpretation is totally different but the rank similarity seems to indicate that they are similarly relevant as well.

Feature ranking does not allow to distinguish among features that are directly related to the output and those that influence it only indirectly. Applications focusing on direct links (e.g., network inference [De Smet and Marchal, 2010; Huynh-Thu et al., 2010; Altay et al., 2011; Marbach et al., 2012], see also Chapter 8), must therefore filter out the indirect component from feature selection methods.

There is thus a high interest in designing new techniques to extract more complete information about input-output relationships than a single global feature subset or feature ranking. A first step towards more interpretable results could be to derive more than one (relevance) score to capture the interest of a feature in several settings. Chapter 6 extends classical tree-based feature ranking to incorporate a contextual analysis.

FINITE SAMPLE SIZE MAKES FEATURE SELECTION MORE DIFFICULT High dimensionality together with small sample-size are nowadays typical in many application domains and it poses a great challenge for classical machine learning techniques [Raudys and Jain, 1991; Braga-Neto and Dougherty, 2004; Molinaro et al., 2005; Saeys et al., 2007] and in particular for feature selection [Sima and Dougherty, 2006; Saeys et al., 2008b; Meinshausen and Bühlmann, 2010; Kuncheva, 2007; Bolón-Canedo et al., 2015; Kuncheva and Rodríguez, 2018]. In such conditions, feature selection is however all the more interesting and may help, for example, to counter-balance the disadvantageous features/samples ratio by reducing the number of variables. Nevertheless, studies show that selecting features in such datasets (e.g., micro-arrays of gene-expression) is less reliable [Jain and Zongker, 1997; Sima and Dougherty, 2006]. In this case, feature selection methods may not necessarily provide a close-to-optimal feature set (i.e., whose error is close to the minimal achievable error) [Sima and Dougherty, 2006; Hua et al., 2009]. They also may be unable to find a satisfying feature subset and this does not imply either that an optimal subset does not exist [Sima and Dougherty, 2006; Hua et al., 2009].

Despite an expensive computational cost, the evaluation function must be properly (cross-)validated²² to avoid the risk of overfitting and overestimated accuracy performances (known as the so-called "peeking phenomenom"²³ or as "selection bias" problem [Ambroise and McLachlan, 2002]) [Reunanen, 2003; Smialowski et al., 2009; Pereira et al., 2009; Diciotti et al., 2013; Kuncheva and Rodríguez, 2018]. Hybrid data (i.e., coexistence of categorical and numerical data) are also worthy of attention [Wang and Liang, 2016; Jiang and Wang, 2016].

In small sample-size conditions, small changes (e.g., addition/removal of samples or noise added to features [Saeys et al., 2008b]) may have a strong influence on the selected feature subset²⁴. For the sake of interpretation for instance, one would usually prefer some stability in the outcomes of feature selection algorithm. In a cross-validation feature selection, this would be highly undesirable to have tremendously

²²Meinshausen and Bühlmann [2010] however claim that cross-validation may fail for high-dimensional data and alternatively propose a stability selection based on subsampling in combination with selection algorithms.

²³It occurs when data dedicated for testing the model is already used in a pre-processing stage such as feature selection. This results in an optimistically biased estimation of accuracy performances for the selected model [Diciotti et al., 2013; Kuncheva and Rodríguez, 2018].

²⁴Let us note that the existence of multiple sets that are equally good may also lead to some instability in selected feature sets [He and Yu, 2010].

different selected feature sets from two folds drawn from the same dataset. Stability of feature selection with respect to sampling variation have drawn researchers' attention as another step towards a more robust feature selection [Kuncheva, 2007; Kalousis et al., 2007; Saeys et al., 2008b,a; Abeel et al., 2009; He and Yu, 2010].

In small sample-size conditions, irrelevant variables may seem relevant due to random fluctuations. Indeed, the risk of having spurious associations between irrelevant features and the output increases with a decreasing sample-size, especially if the number of features is large [Kursa and Rudnicki, 2011]. Discerning barely but truly relevant from falsely relevant features is a common issue in feature selection with high-dimensional datasets. Solutions, such as introducing an artificial random contrast variable [Stoppiglia et al., 2003b; Tuv et al., 2006; Rudnicki et al., 2006; Kursa and Rudnicki, 2011; Huynh-Thu et al., 2012] or using dimensionality reduction techniques (by random projections, see random subspace method [Ho, 1998] in Chapters 3 and 7), are required to do so.

☑ Chapter take-away

Supervised machine learning aims at exploiting a learning set to gain understanding about the interactions among input features and a target output and to build models to make as accurate as possible predictions of the target based on a subset of the inputs. When considering the relation between the input features and the target output, several notions of relevance and redundancy have been defined in the literature and are of interest. These notions may be exploited in many different ways in order to propose feature ranking and feature selection algorithms. Feature selection is often paramount in order to optimize the accuracy of machine learning algorithms, specially in the context of small sample-size and/or high-dimensionality. More and more practical applications are concerned.

Overview

In this chapter we explain the essential ideas of tree-based supervised learning methods. We focus on classification problems, i.e. supervised learning problems where the target variable Y takes a finite number of unordered values called classes. Occasionally we however mention how presented ideas would carry over to the case of regression trees. Our goal is to provide the required notions used in subsequent chapters, while also providing an intuitive understanding of the main features tree-based supervised learning. After a brief introduction, Section 3.2 provides the main building blocks, namely single decision trees and their greedy recursive partitioning based learning algorithm. Then, in section 3.3, we consider tree-based ensemble methods, and more particularly those used in the subsequent chapters.

"May the forest be with you."

3.1 INTRODUCTION

A popular and classical approach to solve a complex problem is the divide-and-conquer strategy. It consists in (recursively) dividing the problem into several sub-problems easier to solve. The solution of the original problem is then a combination of the sub-problem solutions. Based on that strategy, the *recursive partitioning method* aims at simplifying a task to carry out on a set of elements (e.g., sorting, labelling, . . .) by recursively dividing the set into smaller and smaller subsets in such a way that doing this task is easier in each subset than in the original set. For example, sorting can be achieved efficiently using this strategy: the *merge-sort algorithm* recursively divides the list of elements into smaller and smaller groups until each sub-group is easy (or trivial) to sort, and then combines sorted sub-lists.

The decision tree algorithm successfully applies this method to provide a supervised learning model that partitions the input space into distinct (smaller) subspaces [Breiman et al., 1984; Quinlan, 1986, 2014]. As a sub-problem, an output value is then assigned to each subspace. From there, the prediction of a new object simply consists in identifying the subspace in which it falls to retrieve its predicted output value.

Single decision trees are simple and consistent supervised models making them easy to use and to understand. They however suffer from variance, and their accuracy performances are consequently affected.

In order to circumvent variance issues and thus improve model performances, Ho [1998]; Dietterich [2000]; Breiman [2001] were among the firsts to propose to grow an ensemble of trees instead of settling for a single one. Making a predic-

tion by letting every tree vote and then aggregating these votes results in significant improvement in accuracy. Many state-of-the-art algorithms stemmed from that idea, including random forests and boosting methods. In particular, a random forest is an ensemble (i.e., a forest) of randomised trees and is at the centre of this thesis. Randomisation is introduced to create some diversity between trees of the same ensemble. The motivating assumption of this approach is that the prediction of an ensemble of weak models is better than the prediction of a single (supposedly stronger) model.

Furthermore, the success of tree-based methods is also explained by their following common characteristics [Geurts, 2002; Louppe, 2014]:

- NON-PARAMETRIC NATURE by not requiring a priori assumptions on the relationships between inputs and output,
- ABILITY TO HANDLE HETEREGENEOUS DATA by handling learning sets made of a mix of continuous, discrete (ordered or not), and categorical variables (but not necessarily fairly, see Section 4.4.5.3 for more details),
- ROBUSTNESS TO OUTLIERS OR ERRORS IN LABELS by usually avoiding to completely modify the model to fit a few spurious values in the data,
- ROBUSTNESS TO IRRELEVANT OR NOISY VARIABLES by automatically selecting the most useful (and relevant) features to build the tree structure (at least to some extent, see Chapters 4 and 5 for more details),
- INTERPRETABILITY by providing a decision path (with decision trees) or an importance degree for used features (with ensemble methods, see Chapters 4 and 5 for more details),

In this chapter, Section 3.2 describes the decision tree algorithm. Then, Section 3.3 presents ensemble methods as a way of circumventing the high variance of decision trees.

- 3.2 SUPERVISED LEARNING WITH DECISION TREES
- 3.2.1 Semantics of tree based prediction models
- 3.2.1.1 From graph theory to decision tree terminology

In all generality, let G = (V, E) be a graph where V is a finite set of nodes t (also denoted as *vertices* in graph theory), and $E \subset V \times V$ is the set of *edges*. The graph is called *undirected*, if $(t_i, t_i) \in E$ implies that also $(t_i, t_i) \in E$. In graph theory, a tree is an undirected graph in which any two vertices are connected by exactly one (undirected) path.

We use the term tree structure to denote a directed graph obtained from a tree by choosing a node as the root (denoted t_0), and by directing all edges 'away' from this root (see Figure 3.1 for an illustrative example). A branch (t_i, t_{i+1}) is an edge going from t_i towards t_{i+1} where t_i is called the parent of t_{i+1} , and t_{i+1} is a child of ti. A node is internal if it has at least one child, and terminal (also known as leaf node in the tree terminology) if it has no children.¹

¹Internal nodes generally have several children, while every node has exactly one parent. The number of branches of a tree structure is always equal to its number of nodes minus 1.

Figure 3.1 gives an example of a tree structure (i.e., a tree-structured graph). It is represented with the (internal) root node t_0 on top and such that nodes at the same depth (i.e., distance with respect to the root node) are horizontally aligned. Nodes t_1 and t_2 are also internal because they respectively have the children t_3 , t_4 and t_5 , t_6 . Here t_3 , t_4 , t_5 , t_6 are the leaves of the tree.

The following section describes a decision tree model: a tree structure with an additional layer of information.

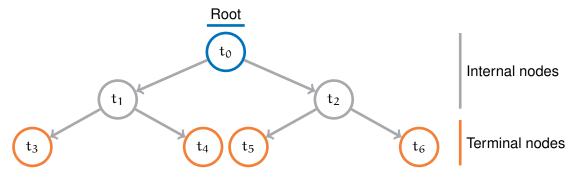


Figure 3.1: Example of a tree-structure.

3.2.1.2 A tree structure shaped by the features

A tree structure G_T recursively partitions the input space $\mathfrak X$ into subsets where each node t is associated to one specific subset $\mathfrak X_t$. The subsets corresponding to the terminal nodes are disjoint and such that their union is the original input space $\mathfrak X$, i.e., $\cup_{t \text{ is terminal}} \mathfrak X_t = \mathfrak X$. The subset corresponding to an internal node is the union of the subsets attached to its children; hence the subset corresponding to the root is always the whole input space. To define all these subsets the tree structure uses features as building blocks. Each internal node typically uses one specific feature, in order to partition its own subset into the subsets corresponding to its children.

In all generality, a *split* s is a partition of a set $\mathcal L$ into a finite number of non-empty and disjoint subsets $\mathcal L^{i}.^{2}$ In other words, every element of $\mathcal L$ belongs to one and only one $\mathcal L^{i}$. A split on a node t, also known as a *test* and denoted by s_{t} , is a split of $\mathcal X_{t}$ using the value of a feature to compute the partition. A *split variable* $v(s_{t})$ is the variable on which the test s_{t} is based and is the one that corresponds to node t in the tree structure.

The cardinality of a split s_t , denoted $|s_t|$, corresponds to the number of created subsets, or equivalently the number of possible test outcomes. Cardinalities may or may not be the same for all t. The cardinality $|s_t|$ also determines the number of children of node t (i.e., the *node cardinality*) and may depend on the number of possible values for the split variable (the *variable cardinality*).

A split is said to be *binary* if exactly two subsets are created. However, a node can be divided in more than two by a so-called *multiway* splits. A multiway split is said to be *exhaustive* if the split cardinality is equal to the number of values of the split variable (i.e., one value per branch).

²i.e. such that $\forall i : \mathcal{L}^i \neq \emptyset$, $\forall i \neq j : \mathcal{L}^i \cap \mathcal{L}^j = \emptyset$, and $\cup_i \mathcal{L}^i = \mathcal{L}$.

Some authors have looked at more exotic splits. An oblique split is made by using a linear combination of several numerical features to create the partition.3 In an even more general framework, multivariate splits also consider complex models (e.g., a decision tree [Botta, 2013]) as separating functions, extending axis-parallel and oblique splits [Gama, 2004]. Fuzzy trees do not longer consider disjoint subsets for children but take advantage of the fuzzy logic to allow some (uncertain) samples to be in several terminal nodes [Janikow, 1998; Olaru and Wehenkel, 2003].

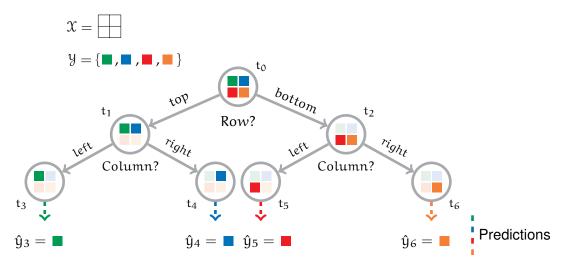


Figure 3.2: Example of a decision tree model: a tree-structure that recursively splits the 2×2 input space with four colours.

3.2.1.3 Decision tree models

A decision tree model $T: \mathcal{X} \to \mathcal{Y}$ recursively partitions the input space \mathcal{X} into subspaces to provide an input-output model in the form of a tree structure (see Figure 3.2 for an illustrative example). The model is such that

- (a) each node t corresponds to one subset $\mathcal{X}_t \subseteq \mathcal{X}$, in particular the one associated to the root node is the input space \mathcal{X} itself,
- (b) each internal node t is labelled with a split s_t,
- (c) each branch going from an internal node t indicates one possible outcome i of the split s_t , and leads to one child c_i of t such that its subset is $\mathcal{X}_{c_i} = \mathcal{X}_t \cap \mathcal{X}^i$ where $X^i \subset X$ is the subset of inputs satisfying outcome i,
- (d) all terminal nodes t have their subsets (called terminal subsets) assigned to a predicted value $\hat{y}_t \in \mathcal{Y}$; \hat{y}_t is also called the label of the leaf t.

Figure 3.2 shows a decision tree model that decomposes an input space of two dimensions (represented by a 2×2 matrix) with four possible output values (i.e., green, blue, red, or orange) using the tree-structure of Figure 3.1. The root node t_0 corresponds to the complete input space x. Its split is made on the vertical axis ("Which row?") and gives two children (t₁ and t₂) corresponding to the two possible outcomes (i.e., top or bottom). Each child has its own subset that is still made of

³Such splits are said to be *oblique* because they produce separating hyperplanes that are not axis-parallel like classical splits made on a single numerical feature. They lead to shorter trees but are more complex to learn [Heath et al., 1993; Murthy and Salzberg, 1995a; Rokach, 2008].

two colours. By splitting them on the horizontal axis ("Which column?"), we obtain four terminal nodes, each with a subset of only one colour. At this point, there is no interest in further partitioning these subsets. The output label associated to each terminal node is immediate and corresponds to the remaining colour. The prediction of the model for a new input value x is the associated value of the terminal node reached by x. Let us observe that, for each internal node t, the input subsets of its children are disjoint and their union is the subset of that node t, i.e., $\mathcal{X}_t = \bigcup_{i=1}^{|s_t|} \mathcal{X}_{c_i}$.

Let us consider a more realistic classification problem, described in Example 3.1, that will be used to illustrate the two following decision tree models.

Example 3.1. Let us consider a classification problem with two input variables X_1 and X_2 with two possible output classes c_1 and c_2 . Figure 3.3 illustrates the learning set where each input variable corresponds to one dimension. At first sight, based on Figure 3.3a, this is not straightforward to give a model that will perfectly separate the two classes. For the sake of illustration, Figure 3.3b gives a decomposition of the input space that provides a perfect separation between objects of different classes.

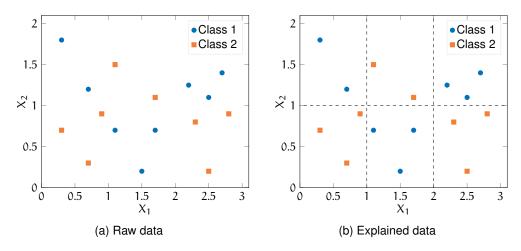


Figure 3.3: Example of a classification problem with two input variables X_1 and X_2 and two possible values of the output y (c_1 and c_2). Blue dots correspond to objects class c_1 while orange squares correspond to objects class c_2 . On the right figure, the underlying decomposition of the input space is explicitly given.

Definition 3.1. A binary decision tree is a decision tree model in which all internal nodes have exactly two children.

This is the case when all splits are binary, that is to say, when there are only two possible outcomes (e.g., *true* or *false*, *yes* or *no*), or when all input features are binary.

A split s divides the input space between the part that satisfies the test \mathcal{X}^s and the rest $\mathcal{X}^{\bar{s}}$. Therefore, the input subspace of the left child c_1 of t (i.e., satisfying the test) is $\mathcal{X}_{c_1} = \mathcal{X}_t \cap \mathcal{X}^s$, and the input subspace of the right child c_r is $\mathcal{X}_{c_r} = \mathcal{X}_t \cap \mathcal{X}^{\bar{s}} = \mathcal{X}_t \cap (\mathcal{X} \setminus \mathcal{X}^s)$.

Figure 3.4 shows a binary classification tree applied on Example 3.1.

Decision trees are typically binary but they can also be built using multiway splits. Figure 3.5 illustrates a multiway decision tree applied on Example 3.1. In comparison with the binary decision tree of Figure 3.4, threeway splits are used at the

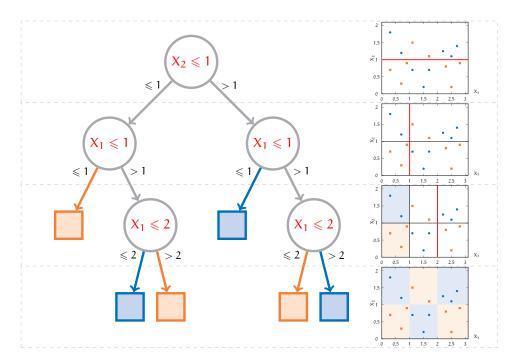


Figure 3.4: Example of a binary classification tree applied on Example 3.1.

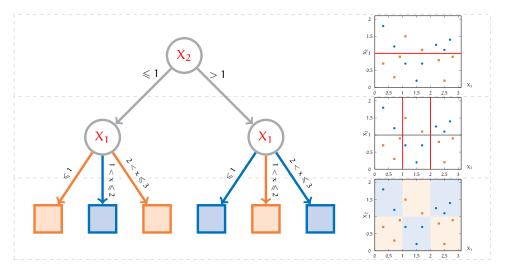


Figure 3.5: Example of a multiway classification tree applied on Example 3.1.

second level to create three children corresponding to three intervals of values of X_1 (in the example, intervals [0,1],[1,2],[2,3]). Notice that while the binary tree uses two more splits, it manages to find the same final partition.

When all input variables are categorical, let us define a decision tree using multiway exhaustive splits:

Definition 3.2. Let all input variables $V = \{X_1, \dots, X_p\}$ be categorical. A **multiway** exhaustive decision tree is a decision tree model in which splits on feature Xi yield exactly $|X_i|$ children, namely one for each possible value of the split variable X_i .

Multiway exhaustive splits⁴ are typically of various cardinalities as they depend on the number of possible values of each split variable. Notice that for such a tree, the maximal depth is limited by the number of features, as each feature can be used at most once along a path.

3.2.2 Learning a decision tree model from data

The tree model aims at fitting at best the partition induced by y over x and thus approximating the Bayes model (i.e., the optimal model yielding the lowest error rate). In practice, the partition induced by \mathcal{Y} over \mathcal{X} is unknown and the input space is only partially observed through a learning set. Given a learning set LS, a decision tree model T^{LS} is learnt on LS and provides a partitioning of LS, denoted φ . While growing the decision tree, the objective is to find the partitioning φ that provides the lowest possible error rate, the optimal induced partitioning φ^* . Assuming that the learning set represents faithfully the input space, φ^* should be close to the partition induced by y over x.

The tree learning algorithms that we consider in this thesis (and which have become a standard in supervised learning) proceed in a top-down fashion, by starting with the root node and progressively developing the tree structure, while at each step choosing a node to split and a way to split the node, until the tree fits the learning sample sufficiently well (see side note on page 64).

This procedure aims at finding a suitable tree structure, and at associating the right class label to each one of its terminal nodes. This thought has been summarised by Breiman et al. [1984] as follows:

"It turns out that the class assignment problem is simple. The whole story [of the construction of a tree] is in finding good splits and in knowing when to stop splitting." [Breiman et al., 1984]

The three next sections are dedicated to a detailed description of these three key steps of a decision tree learning procedure. In Section 3.2.2.1, we describe how to find the variable (and the associated test) that provides a "good" split for a learning subset. In Section 3.2.2.2, we review some stopping criteria that define the end of the building process. In Section 3.2.2.3, how to choose the labels attached to leaves and used for making predictions.

 $^{^4}$ In the rest of this thesis, multiway splits on categorical features will always be exhaustive, i.e., one child for each value and not for only a subset of values. Therefore, the term "exhaustive" is sometimes omitted.



■ GENERIC TOP-DOWN DECISION TREE GROWING ALGORITHM

- Initialization: create the root node of the tree, attach the whole learning set to this node, and set the list of open nodes to contain only this node.
- Recursion: until the list of open nodes is empty, remove a node from the list of open nodes (following a given growing strategy^a), and decide whether this node should be split:
 - If yes, the node becomes a test node, and a good split for it is determined and used to split the learning set of the node into two or more subsets. For each subset a child node is created and inserted in the list of open nodes.
 - If no, the node becomes a leaf and a class label is assigned to it based on its learning subset.

^aWell-known strategies are depth-first, breadth-first or best-first. Each strategy may yield different decision trees if the stop splitting rule is global, i.e. based on the whole tree.

3.2.2.1 Splitting rules

The growing/learning procedure of a decision THE IMPURITY FRAMEWORK tree model recursively divides the learning set in subsets of learning samples LS_t where LS_t is the set of all objects reaching node t (i.e., LS_t = $\{(x,y)|x \in \mathcal{X}_t\}$). For a given node t and its set of learning samples LS_t, let us define $p(c_i|t)$ as the proportion of samples in LS_t such that $y = c_i$, $c_i \in \mathcal{Y}$. The sum of $p(c_i|t)$ for all $c_i \in \mathcal{Y}$ is 1. Based on LS, the learnt model tries to mimic the optimal induced partitioning φ^* .

A good decision tree is one that minimises the generalisation error while minimising some complexity criterion of three, e.g., the size of the tree. Even though several trees can equivalently represent the optimal partitioning ϕ^* , the shorter tree is usually the easiest to interpret and consequently the best one. Naively, one can generate all possible decision trees in order to keep the best one (minimising a criterion depending on the accuracy performances and the complexity of the model). However, even if the number of trees may be finite when the number of (discrete/categorical) features is limited, this number can increase exponentially and becomes intractable from a computational point of view when considering a large number of (continuous) features.

Circumventing the intractability of an exhaustive search for the optimal tree model (giving φ^*), the idea of Breiman et al. [1984]'s heuristic algorithm is to keep splitting nodes until they are (almost⁵) pure. The resulting partitioning is expected to be close to φ*. A node t is pure when all learning samples reaching that node (LS_t) are of the same class label c_i ($p(c_i|t) = 1$, $c_i \in \mathcal{Y}$ and $p(c_i|t) = 0$ for all $c_i \neq c_i, c_i \in \mathcal{Y}$) (see terminal nodes of Figures 3.2, 3.4 and 3.5). Hereafter, we refer to the output distribution of a pure node as a pure distribution. A pure node is always terminal because there is no gain in splitting more its samples. Conversely, the impurity of a node is the largest when all class labels are equally likely $(p(c_i|t) = p(c_i|t)$ for all $c_i, c_i \in \mathcal{Y}$).

 $^{^5}$ The purity of a node is a natural stopping criterion, but some other criteria exist and may stop the growing process before having pure nodes. See Section 3.2.2.2 for more details.

From that, one can logically assume that the purer a node is, the more striking is the majority class making the prediction easier and usually better.

Following the framework of Breiman et al. [1984], let us define an impurity measure i(t) as a non-negative function ϕ that evaluates the purity of a node t from the vector of class proportion samples π (where the jth term of π , $\pi^j = \mathfrak{p}(c_i|t)$) and verifies the following three properties [Breiman et al., 1984; Jolv. 2017]:

- (a) i(t) is minimal (typically equal to 0) when the node t is pure, i.e., $p(c_i|t) = 1$ for some $c_i \in \mathcal{Y}$ and $\forall c_i \neq c_i : p(c_i|t) = 0$,
- (b) i(t) is maximal only when the distribution of output values in LS_t is uniform, i.e. such that $p(c_i|t) = \frac{1}{|y|}$ for any $c_i \in \mathcal{Y}$,
- (c) i(t) is not biased towards some output values (symmetrical with respect to the class proportion samples), e.g., the impurity measures of two nodes t₁ and t₂ are the same if π_2 is a permutation⁶ of π_1 .

THE GOODNESS OF A SPLIT A good split is one that reduces the impurity i(t)of a node t, i.e., such that children of t are purer than t itself. The goodness of a split dividing a node t in two⁷ can be formalised using the impurity measure as follows:

Definition 3.3. Let s be a binary split that divides a node t into a left node t_1 and a right node t_R. The decrease of impurity is

$$\Delta i(s,t) = i(t) - \frac{N_{t_L}}{N_t} i(t_L) - \frac{N_{t_R}}{N_t} i(t_R)$$
 (3.1)

$$= i(t) - p_{t_{I}} i(t_{L}) - p_{t_{R}} i(t_{R})$$
 (3.2)

where N_t is the number of learning samples in node t, N_{t_1} and p_{t_1} (respectively, N_{t_R} and p_{t_R}) are the number of samples and the proportion of samples that fall into t_L (resp., t_R).

We will discuss later on several impurity measures that may be used for growing decision trees. Once the impurity is chosen, the greedy procedure for growing a decision tree consists in searching at each node for the split that yields locally the largest decrease of impurity $\Delta i(s,t)$ among all valid splits.

CANDIDATE SPLITS FOR DIFFERENT TYPES OF FEATURES Let St.m be the set of all candidate splitting functions for node t on feature X_m, consisting of all candidate ways to divide $\mathfrak{X}_{t,\mathfrak{m}}$ in two non-empty subsets, where $\mathfrak{X}_{t,\mathfrak{m}}$ denotes the set of all values of X_m observed in the learning sample of node t.

If X_m is an unordered variable, defining a split amounts to find two non-empty subsets $\mathcal{X}_{t_I,m}$ and $\mathcal{X}_{t_R,m}$ such that every element of $\mathcal{X}_{t,m}$ is in one and only one of them, i.e., $\mathfrak{X}_{t,m} = \mathfrak{X}_{t_l,m} \cup \mathfrak{X}_{t_R,m}$ and $\mathfrak{X}_{t_l,m} \cap \mathfrak{X}_{t_R,m} = \emptyset$. In that case, $S_{t,m}$ can be formally defined as follows:

$$S_{t,m} = \{s(\mathbf{x}) = \mathbb{1}(x_m \in \mathcal{X}_{t_1,m}) | \mathcal{X}_{t_1,m} \subset \mathcal{X}_{t,m}\}$$
 (3.3)

where x is a vector of input values and x_m is the value of X_m . All splits guide samples whose value x_m is in $x_{t_1,m}$ in the left child, while all others go in the right child. Let

⁶The same numerical values but not necessarily in the same order.

 $^{^{7}}$ For the sake of clarity, only binary splits are considered hereafter but one can naturally generalise what follows for multiway splits by considering $|s_t|$ children instead of two.

us note that $\mathfrak{X}_{t_L,\mathfrak{m}}$ must be non-empty, and a proper subset of $\mathfrak{X}_{t,\mathfrak{m}}$ to ensure that $\mathfrak{X}_{t_R,m}$ is also non-empty. A combinatorial analysis gives that the number of possible splits $|S_{t,m}|$ is equal to $2^{|X_{t,m}|-1}-1$ where $|X_{t,m}|$ is the cardinality of $X_{t,m}$.9

If X_m is an *ordered variable*, the logic between values should be preserved by the split. Consequently, the two disjoint non-empty subspaces $\mathfrak{X}_{t_L,m}$ and $\mathfrak{X}_{t_R,m}$ must be such that every element in one subspace has a split variable value strictly lower than the split variable value of any element from the other subspace, i.e., $x_{t_1,m} < x_{t_R,m}$ for all pairs $(x_{t_L}, x_{t_R}) \in \mathcal{X}_{t_L, m} \times \mathcal{X}_{t_R, m}$. An equivalent way to fulfil that condition is to determine a threshold value τ (also called cut-point), and to assign every value below τ to the left child and to the right child otherwise, i.e.:

$$S_m = \{s(x) = \mathbb{1}(x_m \leqslant \tau) | \tau \in \mathfrak{X}_m\}$$
 (3.4)

where τ is a threshold value referred to as the *cut-point* of the split.

In practice, it suffices to consider a single candidate cut-point between each pair of successive values of the concerned feature observed in the learning subset of the node t (in most implementations it is the mid-point). Indeed, different cut-points between a given pair of such successive values yield the same partition of the learning sample of the considered node, and are thus equivalent from the viewpoint of impurity reduction. The number of different splits to consider is thus $|S_{t,m}| = |X_{t,m}| - 1$. Let us however notice that all cut-points between two successive values (as observed in the learning set) are not necessarily equivalent outside the learning set (see Figure 3.6 for an illustrative example).

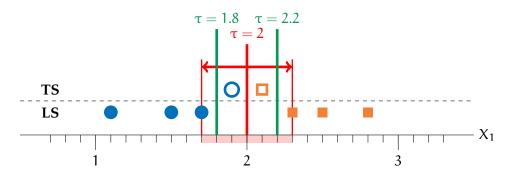


Figure 3.6: Split selection. Projection on the X_1 axis of samples reaching the second node that splits on X_1 (i.e., $X_1 \le 2$) on the left branch (i.e., $X_2 \le 1$) of the decision tree of Figure 3.4. Filled circles and squares are samples from the learning set LS and non-filled ones are samples from the testing set TS (unknown in the learning phase). In practice, all cut-points in the red zone (i.e., between two successive values]1.7, 2.3[) are equivalent on the learning set and $\tau = 2$ was chosen in Figure 3.4. However, other values such as $\tau = 1.8$ or $\tau = 2.2$ also perfectly separate classes in the learning but not on the test set.

Let S_t be the set of splits on all p features and such that $S_t = \bigcup_{i=1}^p S_{t,i}$. The best split s_t^* is therefore

$$s_t^* = \underset{s \in S_t}{\text{arg max}} \Delta i(s, t). \tag{3.5}$$

 $^{^8}$ In practice, it prevents one of the child nodes from having zero learning samples (i.e., $N_{t_{
m L}}=0$ or $N_{t_R}=0$) which corresponds to a split devoid of interest. ⁹Taking into account the fact that exchanging $\mathfrak{X}_{t_L,m}$ with $\mathfrak{X}_{t_R,m}$ leads to an equivalent split.

In practice, Equation 3.5 is solved by exhaustively considering all features and all possible splits on those features (either all cut-points τ or all subsets $\mathcal{X}_{t_L,m}$). This approach however only optimises the split for the current node. Growing a decision tree while foreseeing some future splits is known as (limited) lookahead search and has been shown to provide shorter but not significantly better trees while being computationally more costly [Murthy and Salzberg, 1995b; Louppe, 2014].

SUITABLE IMPURITY MEASURES Any function satisfying the three properties of an impurity measure can be plugged in the decision tree algorithm. Classical impurity measures used for classification problems¹⁰ are the Shannon entropy and the Gini index.

Definition 3.4. The impurity function $i_h(t)$ of a node t derived from Shannon entropy [Shannon and Weaver, 1949] is

$$i_h(t) = -\sum_{j=1}^{C} p(c_j|t) \log_2(p(c_j|t))$$
 (3.6)

where C is the number of possible classes.

Shannon entropy quantifies the uncertainty of a discrete random variable based on its probability density. It is non-negative, maximal for a uniform density, and equal to zero (hence minimal) when only one value has a strictly positive probability. Notice that the entropy-based impurity reduction $\Delta i_h(s,t)$ is actually an estimation, based on the learning subset reaching the node t, of the mutual information between the split outcome and the output t. This impurity reduction is also non-negative, and equal to zero only if the class proportions in the two subsets are identical.

Definition 3.5. The impurity function $i_g(t)$ of a node t derived from Gini index [Gini, 1912] is

$$i_g(t) = \sum_{j=1}^{C} p(c_j|t)(1 - p(c_j|t))$$
 (3.7)

where C is the number of possible classes.

The Gini index quantifies the dispersion of a distribution. The gini-based impurity $\mathfrak{i}_g(t)$ aims at evaluating the error rate of a random labelling of objects from \mathbf{LS}_t following the distribution of labels within node t, p(y|t). That is, the probability of labelling an object with class c_j is given by the probability $p(c_j|t)$ while $1-p(c_j|t)=\sum_{i\neq j}^C p(c_i|t)$ is the probability of error when labelling an object c_j . Similarly to Shannon entropy, $\mathfrak{i}_g(t)$ is non-negative, maximal for a uniform distribution, and equal to zero and hence minimal for a pure distribution. The resulting impurity reduction is also non-negative, and equal to zero only if the class proportions in the two subsets are identical.

EXTENSION TO REGRESSION TREES In order to extend the tree growing algorithm to the case where the output is numerical (i.e. for regression), various alternative goodness of split measures have been defined in the literature. In particular,

¹⁰The above ideas have also been extended to regression problems, where the (empirical) variance is typically used to measure impurity Breiman et al. [1984].

for 'least squares regression', a natural way to do this is to use the same approach as above while using as "impurity" measure the variance of the output Y estimated from a learning subset [Breiman et al., 1984].

Definition 3.6. The "impurity" function $i_{\nu}(t)$ of a node t derived from the variance is

$$i_{\nu}(t) = \frac{1}{N_t} \sum_{y \in \mathcal{Y}_t} (y - \bar{y}_t)^2$$
 (3.8)

where N_t is the number of learning samples in node t and $\bar{y}_t = \frac{1}{N_{\star}} \sum_{y \in \mathcal{Y}_t} y$ is the average of y in LSt.

The variance estimate $i_{\nu}(t)$ is non-negative and equal to zero when all samples have the same target value (equal to the mean value). It also leads to an impurity reduction measure that is non-negative.

3.2.2.2 Stopping rules and pruning

In the previous section, we described how to develop a tree by starting with its root node and splitting its nodes so as to maximise at every step the impurity reduction.

Given the recursive nature of the growing process, there comes a stage when it is no longer possible to further divide a sample set. The splitting process then has no choice but to stop if there is no more valid splits for the node. It occurs in the two following situations, seen as inherent stopping criteria:

- (a) Constant output value: all learning observations reaching the node have the same output value, meaning that the impurity of the learning subset is already equal to zero and hence can not be further reduced,
- (b) Constant input values: all learning observations reaching the node have the same value for every input feature, so that the set of available candidate splits is empty.

Let us note that all learning samples may have the same input values (case (b)) while not having the same output value.

Definition 3.7. A decision tree is said to be **fully developed** if all learning subsets corresponding to its leaves have either a constant output (case (a)) or constant inputs (case (b)) and consequently none of the leaves could have been split in a meaningful way.

Fully developed trees are often overfitting the training data. To limit this phenomenon, additional criteria for stopping to split have been imposed.

- (a) Complexity-based stopping criteria aim at preventing the decision tree from becoming too complex. Typical complexity measures are the total number of nodes or the maximal (or average) depth of the tree.
- (b) Impurity-based stopping criteria stops the growing procedure when the possible impurity reduction is not significant anymore. Indeed, since the growing procedure recursively splits the learning set, the number of learning samples reaching deeper nodes decreases typically rather quickly with the tree depth. Deeper nodes therefore typically yield impurity reductions that are less and less significant from a statistical point of view. Thus it has been proposed to stop splitting if

- i. the size of the learning subset of a node is below a given threshold, or if learning subset sizes of its child nodes would be below a given threshold,
- ii. if the best achievable impurity reduction is too small given the size of the learning subset. Instead of setting explicitly a threshold, some statistical measures (e.g., a χ^2 test or a permutation test) can associate a split impurity reduction to a significance level (e.g., a p-value) for which it is easier to find an interpretable threshold value.

It should be noted that a single criterion may be sufficient to stop the construction of a tree although several can be combined. In practice, all criteria are defined by a hyper-parameter whose value must be carefully chosen. By being too restrictive with their values, these criteria would result in a shallow tree that potentially misses some information about the output in the dataset (i.e., a situation of under-fitting). On the other hand, choosing parameter values that are too permissive would not limit the size of the tree enough, causing over-fitting and sub-optimal performances (in terms of generalisation error). All parameters must therefore be carefully tuned in order to achieve the best trade-off for the size of the tree.

Although those stopping criteria may give in practice good results, they may also lead to sub-optimal trees. A few nodes more or less might indeed sometimes produce a significantly better tree. Another way of finding the best model is to first build a fully developed tree and then choose one of its subtrees a posteriori. Techniques following this approach are known as *post-pruning methods*. In practice, a post-pruning method consists in finding the best subtree $T^* \subseteq T$, obtained by contracting an internal node of the fully developed tree T (i.e., replacing it by a terminal node and dropping all its descendent nodes), say one which minimises a given criterion such as the error rate on a independent test set for example.

Therefore, stopping criteria that preventively control the growing of the tree are usually referred to as *pre-pruning methods*.

3.2.2.3 Labeling the leaves

The prediction $T(x) = \hat{y}(x)$ for an input vector x is obtained by propagating x through the tree (following branches according to its values) and then returning the prediction (or label) \hat{y}_t associated to the terminal node reached by x.

During the learning stage, each terminal node t must thus receive a label $\hat{y}_t \in \mathcal{Y}.$ The choice of \hat{y}_t of course aims at maximizing accuracy and hence essentially depends on the nature of the output variable and on the loss function used to measure accuracy. In practice the output label values found in the learning subset of each leaf are used to choose a label such that in the end the total loss is minimised over the learning set.

FOR CLASSIFICATION TREES AND ZERO-ONE LOSS Let us consider a decision tree model to predict $\mathcal{Y} = \{c_1, \dots, c_J\}$. If the goal is to minimise the probability of mis-classification, the label \hat{y}_t associated to a terminal node t is chosen as the most frequent class (output value) among objects reaching node t. That is

$$\hat{y}_{t} = \underset{c_{j}}{\arg\max} p(c_{j}|t). \tag{3.9}$$

Indeed, in classification tasks, the commonly used loss is the zero-one loss, which for a decision tree and its learning set sums up to

$$L^{0-1} = \sum_{t} \sum_{(x^i, y^i) \in LS_t} 1 (y^i \neq \hat{y}_t),$$

where the outer sum is over all leaves of the tree. And thus, choosing for each leaf its label as the most frequent class in its learning subset LSt therefore minimises the total zero-one loss over the complete learning set.

FOR REGRESSION TREES AND SQUARE LOSS Let us consider a regression tree model ($y \in \mathbb{R}$). If the goal is to minimise the expected square error, the label \hat{y}_t associated to a terminal node t is chosen as the average of all output values of objects reaching this terminal node. That is

$$\hat{y}_{t} = \frac{1}{N_{t}} \sum_{y_{t} \in \mathcal{Y}_{t}} y_{t}. \tag{3.10}$$

Indeed, in regression tasks, the commonly used loss is the square loss, which for a regression tree and its learning set sums up to

$$L^{se} = \sum_{t} \sum_{(x^i, y^i) \in LS_t} (y^i - \hat{y}_t)^2.$$

And thus, choosing for each leaf t its label as the average of all y^i values in LS_t therefore minimises the total square loss over the complete learning set.

Interpretability of decision tree models

One of the main strengths of decision tree models is their interpretability [Hastie et al., 2005]. A decision tree model can be naturally represented in the form of a tree-structured graph or seen as a set of mutually exclusive rules. It recursively partitions the input space into subregions. Each of these regions is described by a sequence of feature-based tests.

A decision tree model also helps to fully understand the reasons for a prediction. By following the path of a sample from the root to the terminal node, one can directly retrieve the explanation for the predicted value. This property is desirable in many domains and in particular in medical applications where a model can provide sensitive results such as a diagnosis or a prognosis. In such cases, understanding the reasons driving the model to some conclusions is crucial as wrong decisions might have severe consequences.

In practice, the tree structure gives all features that are involved in the model. More specifically, the followed branch gives the features used for the prediction in particular and the sequential order in which they are used. In addition to that, one can follow the progress of a prediction by tracking the evolution of output values (i.e., class proportions or output averaged value) within nodes in the path. Figure 3.7 is another graphical representation of the classification tree shown in Figure 3.4 which highlights class proportions within nodes. Note that sometimes left and right nodes are rearranged so that the left child always corresponds to an increase of the same class (even if the splitting function must be reversed). However, it can be laborious to understand each decision/node of a decision tree, especially if it is large or deep

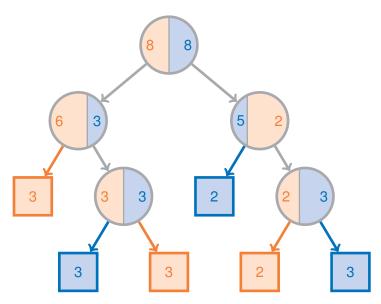


Figure 3.7: Another representation of the binary classification tree in Figure 3.4. In each node, class proportions are represented by the part of the circle filled with the class colour and number of samples of each class are given.

(see [Luštrek et al., 2016] for a study of factors impacting the interpretability of a decision tree).

Furthermore, one may exploit the impurity reductions computed when growing the tree in order to measure the "relevance" of the different input features (see e.g. [Breiman et al., 1984]). Since we will focus on this idea in the subsequent chapters of this thesis, we do not elaborate too much on it here.

On the other hand, an important caveat concerning interpretability stems from the high learning variance of the decision tree growing algorithms [Geurts, 2002] and the so-called "masking effect" [Breiman et al., 1984]. A high learning variance means that small changes to the learning set may lead to large changes in the learnt model. The masking effect denotes situations where several candidate splits on different features yield roughly the same impurity reduction, but one of the features is always slightly better so that none of the other ones has a chance to be selected by the tree-growing algorithm. We highlight both effects on the "XOR" example explained in Figure 3.8.

3.3 TREE-BASED ENSEMBLES

Decision trees are simple and interpretable models but fail to compete with other machine learning algorithms in terms of accuracy. This lack of performances is mostly caused by their very high variance [Geurts, 2002].

This variability stems from the strong sensitivity of the decision tree algorithm to the variability of the learning dataset. Indeed, a small change in the learning set (e.g., due to sampling or noise) may cause significant differences between induced models such as the split choices, the branch depths or the distributions of samples in terminal nodes [Breiman, 1996b; Geurts, 2002]. Any modification has a strong impact on all following decisions because of the recursive nature of the algorithm, resulting in a greatly modified tree structure [Dietterich and Kong, 1995; Schrynemackers, 2015]. In addition, the choice of splits or predictions in deep nodes are made with only few training samples and hence are expected to be of very high

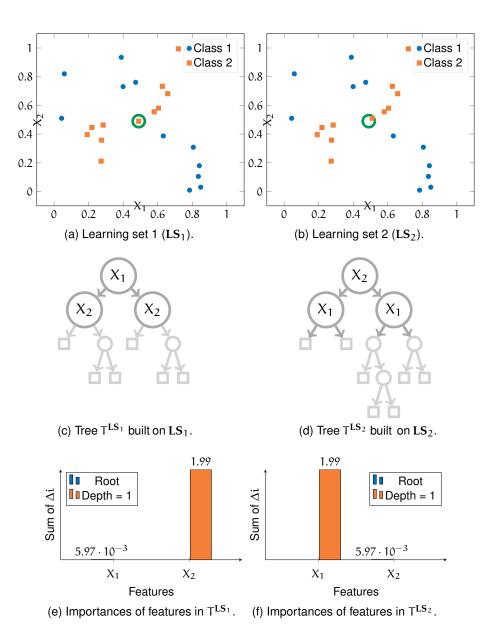


Figure 3.8: Let us consider two highly similar datasets LS₁ and LS₂ made of a set of input features V and a binary output (of two classes). Two features $X_1 \in V$ and $X_2 \in V$ (represented in Figures 3.8a and 3.8b) form a XOR structure that determines the output, i.e. all points with $(X_1 \le 0.5 \text{ and } X_2 \le 0.5)$, or $(X_1 > 0.5 \text{ and } X_2 > 0.5)$ belong to the first class, and to the second class otherwise. Both datasets are identical except one sample (surrounded by a green circle) that has been slightly moved in LS₂. Figures 3.8c and 3.8d show trees built on each learning set respectively. For sake of simplicity, let us assume that X1 and X2 are used on top of the tree and each split has a cut-point at 0.5. In LS₁, X_1 is slightly better than X_2 (masking X_2) and thus selected first, while in LS_2 , the situation is reversed $(X_1 \text{ is now masked by } X_2)$ and X_2 is selected first. The small change only is enough to completely change the (top of the) tree (i.e., the order in which X₁ and X2 are used) and potentially all the rest of the tree, symbolised by shaded different sub-trees (see [Breiman et al., 1984, Figure 5.8] for a complete example). Figures 3.8e and 3.8f show the importances of X_1 and X_2 computed as the (unweighted) sum of Shannon impurity decreases.

variance [Dietterich and Kong, 1995; Geurts, 2002]. Ultimately, the high variance of a decision tree model penalises both its accuracy and its interpretability (at least to some extent).

As a way of increasing the performances, ensemble learning is a technique that is particularly adapted for variance reduction in the context of decision tree models [Louppe, 2014]. Based on the idea of Kwok and Carter [1990]'s 'Multiple decision trees', the principle of this approach consists in combining several different models to achieve better performances than individual ones by aggregating their predictions [Hastie et al., 2005]. Base models of an ensemble are usually built independently of each other and their predictions are either averaged (for a regression task) or aggregated by majority vote (for a classification task). In the same vein, boosting methods do not build independent individual predictors but rather build a sequence of models in which each step builds a predictor trying to refine the predictions of its predecessors.

In what follows, we focus on the first family of methods, usually referred to as averaging methods, where models are built independently and usually differ from each other because of some randomisation introduced in one way or another. We generically denote these methods by "Random forest type of method" to distinguish the family from its particular well-known instance proposed by Leo Breiman and called "Random forests".

3.3.1 Random forest type of methods

Random forest type of methods refers to several tree-based ensemble learning methods based on the idea of randomisation and aggregation. The main common principle is to generate an ensemble of randomised trees (i.e., a forest) in which each individual tree is induced by a randomised version of the classical decision tree growing algorithm, and to combine in a suitable way the predictions of all the elements of this ensemble. Formally, a random forest consists of a collection of N_T tree-structured models $T = \{T_i | i = 1, ..., N_T\}$ used together in the way suggested by Figure 3.9 in order to make predictions.

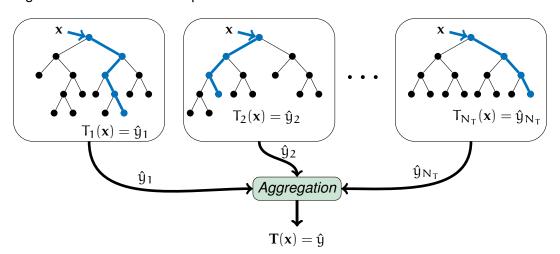


Figure 3.9: Principle of the random forests method. The model T consists of an ensemble of N_T (different) trees. The model prediction $T(x) = \hat{y}$ is the aggregation of the predictions of every individual decision tree model.

The goal of introducing randomisation is to generate diverse tree models, i.e., models whose errors are as much as possible uncorrelated. Indeed, for a given average behavior of the members of the ensemble, the more diverse they are, the smaller is the variance of the ensemble model and the higher is its accuracy (see side note on page 74 and in particular [Hastie et al., 2005; Louppe, 2014; Joly, 2017] for more details).

NUMBER AND DIVERSITY OF TREES IN AN ENSEMBLE

Hastie et al. [2005] motivate the aggregation of several models by giving the variance of the average of:

(a) N_T independent and identically distributed (i.i.d.) random variables, each with a variance of σ^2 , is

$$\frac{1}{N_{T}}\sigma^{2}.$$
 (3.11)

As the number of random variables N_T increases, the variance tends to disappear.

(b) N_T identically distributed (but not independent) (i.d.) random variables, each with a variance of σ^2 and a positive pairwise correlation of ρ , is

$$\rho\sigma^2 + \frac{1-\rho}{N_T}\sigma^2. \tag{3.12}$$

Similarly to the first case, the second term disappears with an increasing N_T . The first term however is independent of N_T but decreases as the variables are de-correlated (i.e., lowering the value of ρ).

Both examples show that trees must as numerous and diverse (i.e., decorrelated) as possible to decrease the variance. It motivates the use of randomisation to generate trees for an ensemble. We refer to Louppe [2014] for a detailed bias-variance decomposition of an ensemble of trees.

In addition to a potential increase of performances, let us note that building a random forest is usually advantageous from a computational point of view. Indeed, the randomisation often cuts the complexity down as it removes heavy computations or reduces the dimensionality of the problem. In addition, the bulk of the learning of a random forest can be parallelised by growing the individual trees independently and exploiting several computers to do so.

Several random forest type of methods have been proposed over the years. They all apply the 'perturb and combine' paradigm and essentially differ from each other only in the way the decision tree procedure is perturbed [Geurts, 2002]. The random perturbation can be introduced in several parts of the algorithm (mainly where the variability is observed), namely at the level of:

(a) the learning set: As discussed in the context of the high variance of decision trees, models are expected to vary if they are built on different learning sets [Breiman, 1996a];

- (b) the split variable selection, i.e., features that are considered at each tree node: not considering all features at each node allows sometimes alternative (e.g. masked) features to be selected;
- (c) the split value selection: the cut-point for numerical features or the binary splitting function for categorical features is chosen at each node at random rather than being optimised in terms of impurity reduction for the learning subset of that node.

Below we explain the involved randomization mechanism of the main random forest type of methods published in the literature 11.

BAGGING - tree-wise learning set randomization

Bagging, standing for bootstrap aggregating [Breiman, 1996a], consists in growing each tree of the ensemble from a bootstrap replicate of the learning set. Given a learning set LS of N samples, a bootstrap sample LS^B is obtained by sampling n samples from LS at random and with replacement [Efron and Tibshirani, 1994]. Let us note that some samples of LS may appear multiple times in LS^B or not at all. On average, around 37% of original samples are not represented in the bootstrap sample [Louppe, 2014], this will be of interest in Section 3.3.2.3. Figure 3.10 sketches the principle of generating bootstrap copies of a learning set, for an ensemble of 5 copies gotten from a learning set of ten samples. Figure 3.11 illustrates the Bagging approach.

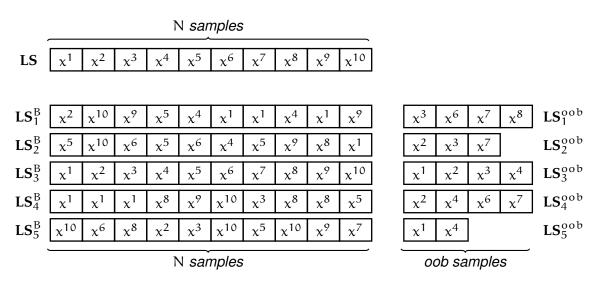


Figure 3.10: Example of five bootstrap replicates of a learning set LS of N = 10 samples. Each x^i represents a sample (x^i, y^i) of the learning set (y^i) is omitted for sake of clarity). On the left, five bootstrap replicates $LS_1^B, LS_2^B, \dots, LS_5^B$ of LS are shown. On the right, sets LS_1^{oob} , LS_2^{oob} , ..., LS_5^{oob} of (out-of-bag) samples that are not used in the corresponding bootstrap samples are highlighted. Sizes of oob sample sets are not necessarily the same. (Figure inspired from Raschka [2016]).

RANDOMIZED TREES - node-wise randomized split selection among best ones With this first randomised version of the decision tree algorithm itself, Dietterich and Kong [1995] extend the idea of Kwok and Carter [1990] and propose to randomise the choice of the split for each node. For a given node t, instead

¹¹See e.g. Louppe [2014] for a more exhaustive list of random forests methods.

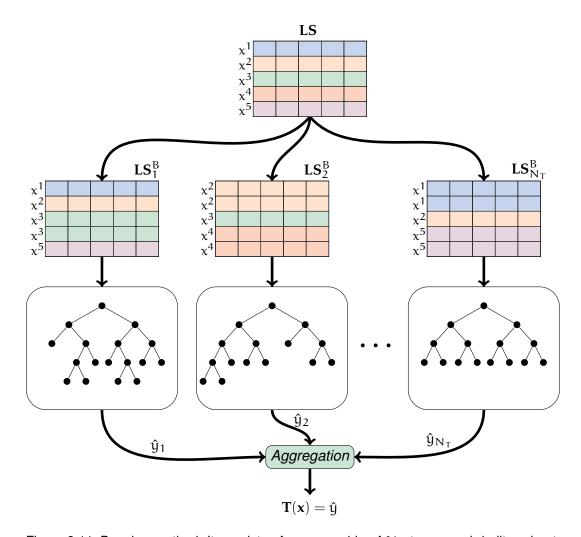


Figure 3.11: Bagging method. It consists of an ensemble of N_{T} trees, each built on bootstrap replicates of LS. Classically, the prediction \hat{y} of the bagging model is the aggregation (majority vote or average) of every individual predictions \hat{y}_i .

of selecting the best split s_t^* , one of the 20 best splits of node t is selected uniformly at random.

RANDOM FEATURE SUBSET - node-wise variable randomization

When the number of variables p is large (e.g., in a handwritten character recognition application), the number of potential splits at each node is typically very large too. In order to avoid a search for the best split among too many possibilities, Amit and Geman [1997] propose to limit the search for the best split among a random subset of only K variables chosen at each node.

RANDOM SUBSPACE - tree-wise variable randomization

Ho [1998] propose to grow each tree of the ensemble on a random subspace, i.e., a learning set in which only K ($\leq p$) features have been randomly chosen. Figure 3.12 illustrates this approach. This method appears as similar to the "Random feature subset" approach, but here one particular tree of the ensemble faces the same subset of features at all its nodes.

RANDOM PATCHES — tree-wise variable and learning set randomization Louppe and Geurts [2012] propose to build an ensemble of trees on random patches where, before building a tree, both a subset of (say K) features and a subset of (say L) learning samples is selected at random. This allows to handle very big datasets and adapt to different types of problems by tuning K and L while keeping $K \times L$ compatible with memory capacity.

RANDOM FORESTS - tree-wise learning set, node-wise variable randomization With Random Forests (RFs), Breiman [2001] combines his idea of bagging with the random feature subset at each node of Amit and Geman [1997] in order to differentiate even more trees by perturbing them in two simultaneous ways. This is undoubtedly the most well known and used version of the random forests methods and more details are given in the following section.

PERFECT RANDOM TREE ENSEMBLES - node-wise split randomization

The novelty of the *Perfect Random Tree Ensembles* (PERT) proposed by [Cutler and Zhao, 2001] is to combine a feature selection totally at random, similar to the random feature subset approach with only one feature considered at each node (i.e., K = 1), and then a random split on that feature. Given an ordered split variable $X_{\mathfrak{m}}$ and a node t, two samples of different output values (classes) in LS_t are selected, say (x^i, y^i) and (x^j, y^j) with $y^i \neq y^j$, and the cut-point τ (the split value) is found as follows $\tau=\alpha x_m^i+(1-\alpha)x_m^j$ where α is drawn uniformly at random between [0,1], $x_{\mathfrak{m}}^{i}$ and $x_{\mathfrak{m}}^{j}$ are respectively the values of variable X_m for samples x^i and x^j .

EXTRA-TREES - node-wise candidate variable and split randomization.

The method of Extremely Randomized Trees or Extra-Trees (ETs) Geurts [2002]; Geurts et al. [2006] draws a random subset of K variables at each node (as the "Random feature subset method") and for each one a single random split, and selects among these K candidate splits the one yielding the largest impurity reduction to split a node. In this method, the cut-point selected for a numerical feature is drawn at each node according to a uniform distribution between the minimum and maximum values of that feature as observed in the local learning subset.

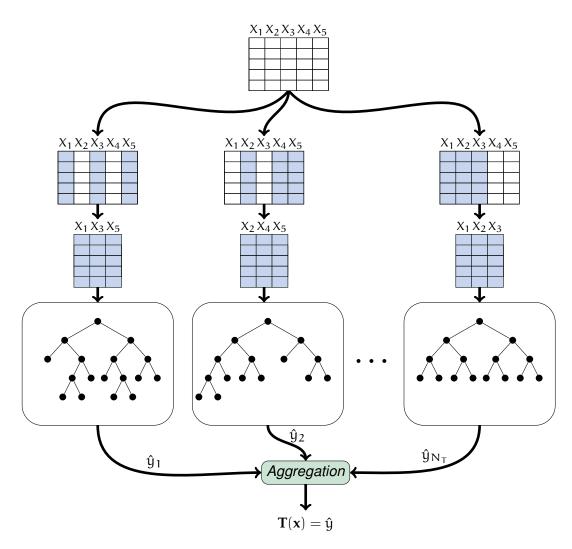


Figure 3.12: Building an ensemble of trees with the random subspace method. Given $\mathfrak{p}=5$ features, each individual tree is learnt on an input subspace made of $\ensuremath{\mathrm{K}}=3$ features that have been randomly sampled.

TOTALLY RANDOMIZED TREES - node-wise split randomization

The method of *Totally Randomized Trees* (TRTs) is a variant of "Extremely randomized trees" maximising the randomization Geurts [2002]; Geurts et al. [2006]. Concretely, it consists in building ETs with K = 1. Node splitting is thus carried independently of the output variable. The method of "Totally randomized trees" is especially of interest in theoretical analyses in the rest of this thesis, in particular in Chapters 4 and 5.

Without further explanation, let us also mention the Rotation Forests method [Rodriguez et al., 2006] which exploits feature extraction principle to build an ensemble of trees on different learning sets.

Random Forests and Extra-Trees: parameters, properties, interpretability 3.3.2

Among all methods, Breiman [2001]'s Random Forests is certainly the most widely known. It was implemented from the very beginning in a freely available and well documented library [Breiman, 2002; Breiman and Cutler, 2003]. Today, it is available within "R" and in the Scikit-learn open-source platform (one of the most used machine learning libraries) which proposes a very efficient and simple to use implementation of both Random Forests and Extra-Trees [Pedregosa et al., 2011]. From a theoretical viewpoint, several authors studied the consistency (i.e., theoretical guarantees that the model converges towards optimality given asymptotic conditions, including a learning set of infinite size) of the method (see, e.g., [Zhao, 2000; Breiman, 2000, 2004; Biau et al., 2008; Biau, 2012; Denil et al., 2014; Scornet et al., 2015]). In conclusion, all the results point in the direction that random forests methods work well in practice (see Louppe [2014] for a review).

In this section, we first go through the different parameters of the Random Forest and Extra-Trees methods and then describe some of their properties that allow us to go beyond a simple predictor, and to some extent interpret the model.

3.3.2.1 Parameters

In this section, we discuss the common parameters of the Random Forest and the Extra-Tree methods. Specific parameters of other random forest type of methods are not mentioned here.

- (a) Randomisation parameter K: It concerns the number of features considered at each node as split variable candidates. Usually given as a function of the number of features, it directly impacts the degree of randomisation of the treebased model. With p features, typical default values for this parameter are K = \sqrt{p} , $K = \log_2 p$ or K = p. Experimentally, it has been shown that \sqrt{p} is usually an appropriate choice for classification tasks, while K = p is often a better choice in case of regression [Hastie et al., 2005; Geurts et al., 2006]. The minimal value, K = 1, implies a maximal randomisation. It may be of interest when all features are a priori known to be more or less equally informative, while large values of K are preferable when a large proportion of irrelevant variables is suspected.
- (b) Number of trees N_T: It defines the number of trees in the ensemble. Intuitively and theoretically, it seems that the number of trees should not be limited as it does not cause over-fitting [Hastie et al., 2005], but performance stabilises

after a certain number of trees depending on the problem considered. However, the number of trees should not be too small either as it has been shown that a certain number of trees is required to achieve the best prediction accuracy or to capture the whole problem structure [Latinne et al., 2001; Genuer et al., 2010; Wehenkel, 2018]. One usually needs to find a good trade-off for the number of trees to achieve good performance while not being too costly in terms of memory or computational resources.

(c) Individual tree complexity: This parameter, unlike the first two, is not only defined by a single value. Several criteria, including of course a simple constraint on the maximal tree depth d, aim at limiting the complexity of the trees. As this corresponds to pre-prune the tree, we retrieve parameters that correspond to the stopping criteria that were discussed in Section 3.2.2.2. In addition to a maximal depth parameter d, n_{min} and n_{leaf} control the growing process of a branch and respectively define the minimal number of samples required to split a node and the minimal number of samples required in child nodes after the split. Δi_{min} and i_{min} respectively prevent the splitting of a node if the impurity reduction is not large enough or if the node has low impurity (i.e., pure enough). N_{nodes} and N_{leaf} control the overall complexity of the tree by defining a maximal number of nodes or leaves.

Let us mention that the choice of the impurity function (typically, Gini or Shannon) for classification tasks is usually left to the discretion of the user.

3.3.2.2 Variable importances

The decision tree model is interpretable. From this model, one can directly read the tree structure giving features that have been used to build the model and how they are split, and the reasons behind a prediction. This was however limited by the high variance of the decision tree model.

When taking an ensemble of trees, the resulting model is indeed more accurate in general but the multiplicity of trees it contains makes it difficult to read and synthesise the information provided by this model. Moreover, because of randomisation, every individual tree structure is also less relevant.

In order to recover some interpretability, the random forest type of algorithms however offer, similarly to single decision trees, the possibility to derive a numerical "importance" value for each feature. This score aims at evaluating the contribution of a feature in the model. Reviewing, studying, and assessing such variable importances derived from tree-based ensemble models is the focus of Chapters 4 and 5. More specifically, Chapter 4 revisits the main variable importance measures, while Chapter 5 is devoted to a detailed analysis of one of these measures in particular, namely the mean decrease of impurity, on which we have focused our research.

3.3.2.3 Out-of-bag samples and estimates

In methods using bootstrapping such as Bagging or Random Forests, for each tree model, there are some samples that have not been used for construction. Given a bootstrap sample set LS_i^B used for tree i, left-out samples $LS_i^{oob} = LS \setminus LS_i^B$ are said to be out-of-bag (OOB) for tree i (see Figure 3.10). These OOB samples can be used to estimate important statistics of the ensemble of trees such as the generalisation error or variable importances (see Section 4.2.2 of Chapter 4).

For each training sample $(x^j, y^j) \in LS$, some trees are built on bootstrap samples that did not include sample j. Let us denote this subset of trees as $T^{-j} = \{T_i^{-j} | i = 1\}$ $1, \ldots, N_T^{-j}$ where N_T^{-j} is the number of such trees. The *out-of-bag error estimate* at $(x^{j}, y^{j}) \in LS$ consists in evaluating the prediction $T^{-j}(x_{i})$ of the ensemble of trees T^{-j} for the input x^{j} . Mathematically, the out-of-bag error estimate over all the learning set is computed as follows

$$\widehat{\operatorname{Err}}^{oob} = \frac{1}{N} \sum_{(\mathbf{x}^{j}, \mathbf{y}^{j}) \in \mathbf{LS}} L(\mathbf{T}^{-j}(\mathbf{x}^{j}), \mathbf{y}^{j})$$
(3.13)

where N is the number of samples in LS. In classification, L and $T^{-j}(x_i)$ are respectively the zero-one loss and the result of a majority vote between all individual predictions $\{T_i^{-j}(x_j)|i=1,\ldots,N_T^{-j}\}$. In regression, L and $T^{-j}(x_j)$ are respectively the MSE loss and the average of all individual prediction, i.e., $\frac{1}{N_{\tau}^{-j}}\sum_{i=1}^{N_{\tau}^{-j}}T_{i}^{-j}(x_{j})$.

The out-of-bag error estimate provides an accurate approximation of the generalisation error (compared to one resulting from a test set of the same size as the training set [Breiman, 1996c] and from a K-fold cross validation 12 [Wolpert and Macready, 1999]). Let us note that the out-of-bag error estimate requires only one ensembles of N_T trees while K-fold cross validation needs to learn K ensemble of N_T trees.

3.3.2.4 Proximity measure

As another by-product, the Random Forests algorithm offers a proximity measure between samples from which a proximity matrix can be derived from the tree-based model [Breiman, 2002; Breiman and Cutler, 2003]. Given a set of N samples, each element (i,j) of the matrix $N \times N$ is the proximity value between samples (x^i,y^i) and (x^{j}, y^{j}) which corresponds to the fraction of trees in which both samples fall in the same terminal node. The intuition is that samples sharing regularly the same terminal node (and thus the same prediction) are close to each other from the point of view of the random forests model. This also provides a comparison of samples that may of high dimensionality and/or made of mixed variables.

This proximity measure can be used to identify structures in the data or for unsupervised learning (see for more details and examples of proximity plot, e.g., [Breiman, 2002; Liaw et al., 2002; Breiman and Cutler, 2003; Hastie et al., 2005; Louppe, 2014; Scornet, 2016]).

¹²K-fold cross validation consists in dividing the learning set into K folds (subsets) of same size and then learning a model on K-1 folds in turn and testing it on the remaining fold.

✓ Chapter take-away

Tree based supervised learning methods have been proposed several decades ago, and studied and used extensively since. With respect to other supervised learning methods, these algorithms are highly scalable, provide interpretable information, but are often suboptimal from an accuracy point of view. In the last twenty years, a significant body of research has been carried out in the machine learning community in order to understand the theoretical features of these methods, and to find out how to improve them. This work has culminated with the idea of building ensembles of randomized trees, rather than one single fully optimized tree. Many different variants of this idea have been proposed over the years, the two most widely used ones being "Random Forests" and "Extremely Randomized Trees". These methods have shown to be very effective in terms of accuracy (among the best general purpose supservised learning algorithms). On the other hand, they lead to less easily interpretable models than the original single decision trees.

Part II CHARACTERISATION OF IMPORTANCE MEASURES

A SURVEY OF THE LITERATURE ABOUT TREE-BASED FEATURE IMPORTANCE MEASURES

Overview

In this chapter we review the literature on tree-based feature importance measures. We start by an intuitive description of this notion and then focus on the two most popular feature importance measures, namely the Mean Decrease of Impurity (MDI) and the Mean Decrease of Accuracy (MDA). We examine theoretical and empirical analyses carried out on those measures and discuss their limitations and biases. Finally, the last parts of this chapter focus on practical applications of those measures, and in particular how to distinguish relevant from irrelevant features based on their importance scores.

Tree-based ensemble methods are known to be powerful methods for modelling complex systems while providing accurate predictions [Auret and Aldrich, 2011]. In many problems, including for example micro-array studies [Archer and Kimes, 2008] or medical prognosis [Wehenkel et al., 2017], a black-box that only provides predictions is however not enough, or even not the main goal. Such applications require indeed to understand how the model is built, to allow some interpretation of results and predictions so as to gain insights on the underlying problem structure [Archer and Kimes, 2008]. However, at first sight, tree based ensemble models are not directly interpretable as the number of trees and the introduction of perturbations in the growing process make their individual interpretation difficult and certainly unreliable [Auret and Aldrich, 2011]. Indeed, two questions are raised among others:

"Is a feature used at the top of only one tree necessarily important?"

"What about features that are only used in a few trees of the ensemble, are they necessarily useless?"

Anticipating this need for interpretability, the Random Forests algorithm (presented in Section 3.3) was proposed together with several built-in measures of feature importance [Breiman, 2001, 2002; Breiman and Cutler, 2003]. Identifying the constitutive elements of the forest model (and their relative importance) is a way to interpret it, and so to gain insight about the underlying problem. Indeed, the variable importance is often presented as a robust statistic to assess the feature contribution in the random forests model of the underlying data generating mechanism [Archer and Kimes, 2008]. Furthermore, these importance measure give an aggregated information, contrasting with the local interpretation of each individual tree.

Concretely, given an ensemble of trees, the principle of *feature importance evaluation* is to derive a numerical score that reflects the "(relative) contribution" of the different candidate features in the learnt model. Based on those scores, one can now evaluate the usefulness of a feature and compare the contributions of two features, whatever the way they are used in the individual trees. A feature having a larger importance score than another one indicates that it is more useful in the learnt model than the other one [Archer and Kimes, 2008]. Conversely, a feature with a very low importance score is not really useful in the learnt model. In addition,

ordering all features according to their importance scores provides a feature ranking [Guyon and Elisseeff, 2006] that may be exploited in different ways.

In this chapter we focus on the subclass of tree-based ensemble methods where all trees are drawn from the same distribution and independently of the others. This choice corresponds, for example, to Tree Bagging, Random Forests, and Totally or Extremely Randomised Trees; but it excludes, for example, Tree Boosting¹ or non-tree-based supervised learning methods. Whenever suitable, we will indicate how the discussed methods could apply to other types of predictors.

Section 4.1 gives an intuitive discussion of the contribution of a feature in a tree-based ensemble model. Section 4.2 provides the definitions of the MDA and MDI measures, the two most used ones, while Sections 4.3 and 4.4 summarise the main theoretical and empirical studies on these measures reported in the literature. Then, the last sections aim at reviewing the main use of those importance measures. In particular, Section 4.5 focuses on techniques to distinguish important features from non-important ones based on their importance scores. Section 4.6 describes several machine learning methods exploiting importance measures or extending them. Section 4.7 is dedicated to other importance measures that have been proposed in the literature. Finally, Section 4.8 aims at describing some practical applications using successfully tree-based feature importance measures.

Remark: in order to make this chapter self-consistent and as complete as possible, we have included in our review results that will be discussed in more details in subsequent chapters of this thesis (and published in [Sutera et al., 2016, 2018]).

4.1 CONTRIBUTION OF A FEATURE TO A TREE-BASED MODEL

In this section, we discuss several possible indicators to evaluate the contribution of a feature in a tree-based predictor. We first look at the role of a feature inside a single decision tree built by the classical CART approach [Breiman et al., 1984] and then consider the case of randomised tree ensembles.

POSITION OF FEATURE SPLITS IN THE TREE Intuitively, the position in the tree structure of the splits using a given feature gives an indication on the importance of that feature: splits close to the root should be more important than those used deeper in the tree. Indeed, in ordre to produce simple trees, the tree growing procedure first considers the most useful splits (corresponding to largest decreases of node impurity) and then refines the model by using less useful ones.

However, this intuitive principle can not be directly transposed to ensemble of randomised trees. In all generality, a feature is used in more than one tree. Instead of a single position, the same feature may be at several (and different) positions in the different trees and one would need to take all of these positions into account to determine which features are the most important ones. For example, a feature might be used deeper in a tree because it has some redundant information with other variables used higher in that tree. Such a feature could be seen as important despite its deep positions in some tree. The randomised nature of the growing procedure (e.g., at the level of split variable selection²) also disrupts the intuitive order in which features are used in the tree. A feature may be used in the top of a tree while

¹Let us note that feature importance can also be derived from ensembles of boosted trees (see, e.g., [Auret and Aldrich, 2011] for a study).

²See Section 3.3 for the other mechanisms.

being barely useful or relevant, e.g., if the split variable selection is randomised, this feature may be considered simultaneously with a lot of noisy irrelevant variables and be the best choice among them.

FEATURE SELECTION FREQUENCY When extended to an ensemble of randomised trees, the position in a tree does not longer reflect the importance of a feature. If we put the node position aside, the decision tree growing procedure still naturally performs a feature selection by selecting the best feature in each node except for the most randomised variant of random forests methods. Intuitively, irrelevant features are not supposed to be selected, or only a very limited number of times by chance, because there is no interest of using them anywhere in the tree. Conversely, relevant features are statistically related to the output and therefore should be regularly used in the model [Konukoglu and Ganz, 2014]. A feature can therefore be seen as important if it is used frequently in many trees. From there, the most straightforward way - although naive - to measure the importance of a feature is to simply count the number of times a feature is used as split variable in all individual trees in the ensemble [Strobl et al., 2007b; Konukoglu and Ganz, 2014; Lundberg and Lee, 2017; Lundberg et al., 2018].

Although it is sometimes not done in the literature, we prefer to normalise the "feature selection importance" by the total number of test nodes of all the trees composing the ensemble, in the following fashion:

Definition 4.1. Let us consider an ensemble $T = \{T_1, ..., T_{N_T}\}$ of N_T trees using a set of input features V to predict an output variable Y. The feature selection **frequency importance measure** Imp^{freq} of $X_m \in V$ in T is the proportion of nodes of the tree ensemble in which X_m has been used as split variable, i.e.,

$$Imp^{freq}(X_m) = \frac{\sum_{i=1}^{N_T} \sum_{t \in T_i} \mathbb{1}(\nu(s_t) = X_m)}{\sum_{i=1}^{N_T} \sum_{t \in T_i} \mathbb{1}}$$
(4.1)

where a node is denoted t and associated to a split s_t with a split variable $v(s_t)$.

Despite its intuitive interest, this importance measure is biased towards features used deeply in trees. Indeed, being selected at the root node only counts for one, while the same feature can be used multiple times deeper in the trees. For example, a barely important feature always selected in each last node of a branch (and providing only marginal impurity reductions) would outscore a feature selected only once at each root node. Moreover, the actual contributions of two features with the same importance (i.e., used the same number of times in the forest model) can be completely different if one yields much larger decreases of impurity than the other. Indeed, some features can be seen many times despite their irrelevance (e.g., because of randomisation) while relevant features are missed because of some undesirable effects (e.g., a masking effect of another feature, see Section 4.4.5 for other examples), impacting directly their importance.

To address those limitations, other criteria of feature importance taking into account the actual contribution of a feature in the learnt predictor should be considered.

TWO WAYS FOR EVALUATING THE ACTUAL CONTRIBUTION OF A FEATURE TO A DECISION TREE PREDICTION As presented in Section 3.2.2.1, the tree growing procedure aims at splitting nodes until all terminal nodes are pure. To that end, each split is optimised by selecting as split variable the feature yielding locally the largest decrease of impurity. The construction of a model is thus completely based on the notion of impurity decrease, and in the eyes of the learning algorithm, a variable is indeed important if it provides a large decrease of impurity. Based on that observation, it makes sense to integrate the amount of impurity decrease obtained thanks to all the splits using a particular feature, in order to evaluate its contribution to making predictions. This rationale leads to the Mean Decrease of Impurity (MDI) importance measure.

Beyond its specific mechanism, the purpose of supervised learning is to enable accurate predictions of the target variable. In this respect, the importance of a feature should be directly related to its contribution to the predictive accuracy of the learnt predictor, or in other ways how this accuracy is affected by not using the concerned feature. This rationale leads to the Mean Decrease of Accuracy (MDA) feature importance measure.

Notice that these two importance measures are not equivalent, since reducing impurity on a learning sample does not necessarily imply increasing accuracy out of the learning sample.

Section 4.2.1 describes the first importance measure based on the contribution of a feature in the building mechanism while Section 4.2.2 presents the second importance measure that associates the contribution of a feature to the impact of its removal on the prediction accuracy.

4.2 MDI AND MDA FEATURE IMPORTANCE MEASURES

In this section, we present the two importance measures, each considering a different aspect of the contribution of features. Section 4.2.1 introduces the *Mean Decrease of Impurity* (MDI) that assesses the importance of a feature based on its average contribution in the impurity reduction in the tree-ensemble growing procedure. Section 4.2.2 defines the *Mean Decrease of Accuracy* (MDA) that evaluates the contribution a feature in terms of its impact on predictive accuracy. Anticipating on the rest of this chapter, let us notice the parallel that can be made with the two feature selection problems (described in Section 2.4.4). The minimal-optimal approach focuses on selecting features that provide the highest accuracy. The all-relevant approach aims at identifying all features that are relevant to the target variable.

4.2.1 MDI importance measure

Used as splitting criterion in decision tree growing [Breiman et al., 1984] and then in tree-based ensemble methods [Breiman, 2001], the computation of impurity and impurity reductions is at the heart of these supervised learning algorithms. Taking advantage of these computations of impurity reductions, Breiman [2002] proposed to evaluate the importance of an input feature $X_{\rm m}$ for predicting the output Y by its *Mean Decrease of Impurity* (MDI), also presented as the empirical improvement in the splitting criterion [Strobl et al., 2007b; Friedman, 2001]³. Concretely, it consists in summing all impurity decreases due to $X_{\rm m}$, weighted by the size of the node (in terms of the relative number of observations reaching that node) and divided by the

³Let us note that the sum of all impurity decreases provided by a feature was already proposed by Breiman et al. [1984] as an importance measure for that feature in a single decision tree.

number of trees composing the ensemble model. For a forest made out of N_T trees, the MDI importance measure is computed as follows:

Definition 4.2. The **Mean Decrease of Impurity importance** Imp^{mdi} of a feature $X_m \in V$ about the output Y is

$$\text{Imp}^{\text{mdi}}(X_m) = \frac{1}{N_T} \sum_{T} \sum_{t \in T: \nu(s_*^*) = X_m} p(t) \Delta i(s_t^*, t) \tag{4.2} \label{eq:4.2}$$

where p(t) is the ratio N_t/N between samples reaching node t (N_t) and the total number of samples (N), and $v(s_t)$ is the split variable of s_t .

This definition of the MDI importance can be applied with any impurity measure, including Gini impurity and Shanon entropy used for decision tree growing, and variance used for regression tree growing (see Section 3.2.2.1).

DISCUSSION The underlying assumption of MDI is that all relevant features, i.e., related to the output and thus important, will show up to be useful to discriminate Y at some point of the ensemble learning, and thus yield a high enough decrease of impurity to lead to their selection as split variable, while, on the contrary, irrelevant features are expected to provide no (too small) impurity decrease in any context, and so will be selected only with very low probability as split variable when growing a tree. It may occur that some noisy features yield (e.g., at nodes with a small number of samples) are still selected, but their (low) impurity decrease should be toned down by the weighting mechanism. Let us however note that the MDI importance can not be negative as a split never increases the impurity of a node, i.e., $\Delta i(s_t^*,t) \ge 0$. Konukoglu and Ganz [2014] see the MDI importance as an extension of the selection frequency importance where the split count is weighted by the actual contribution of the feature, i.e., $\Delta i(s_t^*, t)$. The size of the node p(t) is moreover taken into account to balance deep and shallow nodes. There are more deep nodes than shallow ones but usually with less samples.

One of the main advantages of this measure is its computational efficiency. MDI computation is indeed a direct byproduct of the ensemble learning: all impurity decreases are already computed in order to build the tree ensemble [Breiman and Cutler, 2003]. However, it does not explicitly take into account the quality of the generated model, while being important according to MDI in a poor model does not imply much.

4.2.2 MDA importance measure

In tree ensemble learning methods using bootstrapping (Bagging, Random Forests), a tree of the ensemble does not use all samples for its construction. Using these outof-bag samples, Breiman [2001] proposed to evaluate the importance of an input feature X_m by its Mean Decrease of Accuracy (MDA) based the out-of-bag (OOB) error estimate. To this end, the contribution of a feature in a particular tree is evaluated by the impact of its removal on the OOB error-rate for that tree (which is expected to increase for an important feature). The removal of the feature is simulated by permuting in a random fashion its values in the OOB sample, and by evaluating the impact of this on the prediction accuracy of the tree estimated over its OOB sam-

ple. 4 The contribution of a feature for the whole forest is then obtained by averaging this measure over all trees.5

To formalize this idea, let us first consider a given predictor $f(\cdot) \in \mathcal{Y}^{\mathcal{X}}$ and a given sample \mathcal{D} of input-output pairs (x,y) and some loss function L. Let us denote by $\tilde{\mathbb{D}}_{\mathfrak{m}}$ a modified sample obtained from \mathbb{D} by permuting the values of the variable $X_{\mathfrak{m}}$ randomly (and thus independently of the values of Y and all other input features), and define the MDA-estimate of X_m (in f) over \mathcal{D} by

$$\text{Imp}_f^{\mathfrak{mda}}(X_{\mathfrak{m}},f,\mathfrak{D},\tilde{\mathfrak{D}}_{\mathfrak{m}}) = \frac{1}{|\mathfrak{D}|} \left(\sum_{(x,y) \in \tilde{\mathfrak{D}}_{\mathfrak{m}}} L(f(x),y) - \sum_{(x,y) \in \mathfrak{D}} L(f(x),y) \right). \tag{4.3}$$

This quantity is an empirical estimate, based on the sample \mathfrak{D} , of how much the "removal" of variable $X_{\mathfrak{m}}$ influences the accuracy of f as a predictor of y. Its value depends on the particular permutation $\tilde{\mathbb{D}}_m$ used. This dependence can be factored out by averaging over a uniform distribution of permutations, yielding

$$Imp_{f}^{mda}(X_{m}, f, D) = \mathbb{E}_{\tilde{D}_{m}}\{Imp_{f}^{mda}(X_{m}, f, D, \tilde{D}_{m})\}. \tag{4.4}$$

Now, consider a learning set LS of input-output pairs and a tree growing algorithm Algo. Denote by $T = \{T_1, \dots, T_{N_T}\}$ an ensemble of trees where each tree T_i is grown by Algo on a bootstrap replicate LSi of LS, and evaluated on the corresponding OOB sample ($LS_i^{oob} = LS \setminus LS_i$). The MDA importance of a feature X_m derived from Algo is defined as follows:

Definition 4.3. The Mean Decrease of Accuracy Importance Imp^{mda}_{Algo} of a feature X_m about the output Y derived from a bagged version of Algo applied on the learning sample LS is

$$Imp_{Algo}^{mda}(X_m, Algo, \mathbf{LS}) = \frac{1}{N_T} \sum_{i=1}^{N_T} Imp_f^{mda}(X_m, T_i, \mathbf{LS}_i^{oob}, \widetilde{\mathbf{LS}}_{i,m}^{oob}). \tag{4.5}$$

The underlying assumption of MDA is that all important features are related to the output Y, and thus contribute to the ability of the model to predict Y. The permutation of the values of a feature X_m breaks the statistical link between $X_{\rm m}$ and Y, and thus mimics predictions made without using feature $X_{\rm m}$, which are expected to be worse if X_m is an important feature.

A high (and positive) importance value indicates that the variable is important and its removal strongly reduces the accuracy of the tree ensemble-based predictor.

Contrary to MDI, MDA importances can take negative values [Genuer et al., 2010].

Discussion of MDI versus MDA 4.2.3

Both methods can be used for classification and regression problems. MDA depends explicitly on the loss function used, whereas MDI depends explicitly on the impurity

⁴Therefore, the MDA importance is also known in the literature as the *permutation importance*.

⁵Notice that the original definition of MDA importance derived from a Random Forest, as introduced in Breiman [2001], is quite different from the current one adopted later on by several authors (e.g., [Hastie et al., 2009; Genuer et al., 2010; Biau and Scornet, 2016; Gregorutti et al., 2017]); in the original definition, the impact of removing a feature on the accuracy of the whole ensemble model was evaluated, instead of the now used average impact on the accuracy of the individual terms of the ensemble model. It is the more recent interpretation to which we refer in our work.

measure used. Both Imp^{mdi} and Imp^{mda} are random quantities depending on the random learning sample and on the tree ensemble randomisation; Imp^{mda} further depends on the random permutations of the values of X_m. While MDI is defined only for tree-based models, MDA can be used with any bagged supervised learning algorithm, and with slight modification in the loss-estimation method with any supervised learning algorithm.

4.3 THEORETICAL ANALYSES

Supported by the broad success of tree-based methods in applied research (see, e.g., [Svetnik et al., 2003; Díaz-Uriarte and De Andres, 2006; Cutler et al., 2007; Statnikov et al., 2008; Ghimire et al., 2010; Zaklouta et al., 2011; Nayak et al., 2016; Belgiu and Drăgut, 2016]), many authors studied tree-based variable importances to increase their understanding of the methods. Some theoretical analyses about the consistency of the Random Forests algorithm were already mentioned in Section 3.3.2. But only a few works focused on tree-based variable importances from a theoretical point of view and this section aims at summarising these results and at providing the reader with a better understanding of their theoretical properties.

Mechanisms for building a tree-based ensemble, and consequently to derive importance measures, are highly complex because of their randomisation and their data-dependent nature. For that reason, theoretical studies on MDI and MDA usually deal with that complexity by considering either a simplified version of the tree-based algorithm [Ishwaran, 2007], an asymptotic setting [Louppe et al., 2013; Louppe, 2014; Sutera et al., 2018], or even a specific class of supervised learning problems [Gregorutti et al., 2017].

In the present section, we first review the main known theoretical properties of the importance measures focusing on so-called asymptotic conditions, i.e., when the ensemble of trees and the training sample are both assumed to be of infinite sizes. We then discuss theoretical analyses studying the impact of feature correlation or redundancy on importance measures. Empirical analyses of these measures in real settings are discussed in the next section.

Notational conventions

In the present and subsequent sections, MDI and MDA importances derived in asymptotic conditions, i.e. their population versions, are respectively denoted ${
m Im} {
m p}_{\infty}^{{
m mdi}}$ and ${\rm Imp}_{\infty}^{{\mathfrak m}{\rm d}{\mathfrak a}}$. Additional parameters are specified as subscript or superscripts when they have an influence on the importance measure.

4.3.1 Asymptotic properties of MDA

Following Gregorutti et al. [2017], let us introduce the population version of the MDA importance measure (Equation 4.3) in the context of least-squares regression problems. Denote by P(Y,X) the joint distribution of inputs and all outputs, and by $\tilde{P}_m(Y,X)$ the joint distribution obtained by replacing in P the factor $P(X_m|Y,X^{-m})$ by the marginal distribution of $P(X_m)$, i.e. by breaking any link between X_m with the output and all other input features will leaving the marginal distribution of $X_{\rm m}$ unchanged. Denote also by f_B the Bayes model with respect to the original distribution P and the square loss-function (i.e. $L(y, y') = (y - y')^2$):

$$f_{B}(X) = \mathbb{E}_{P} \{Y|X\},$$

where the subscript P indicates the distribution used for computing the conditional expectation. Then the population version of MDA introduced by Gregorutti et al. [2017] is defined as follows

$$\operatorname{Imp}_{\infty}^{\mathfrak{mda}}(X_{\mathfrak{m}}) = \mathbb{E}_{\tilde{P}_{\mathfrak{m}}}\left\{ (Y - f_{B}(X))^{2} \right\} - \mathbb{E}_{P}\left\{ (Y - f_{B}(X))^{2} \right\}. \tag{4.6}$$

Notice that this quantity is non-negative, since $f_{\rm B}$ is the Bayes model with respect to the original distribution⁶.

Obviously⁷.

both ${\rm Imp}_{\rm f}^{{\mathfrak m}{\rm d}{\mathfrak a}}(X_{\mathfrak m},f_B,{\mathfrak D},\tilde{{\mathfrak D}}_{\mathfrak m})$ (Equation 4.3) and ${\rm Imp}_{\rm f}^{{\mathfrak m}{\rm d}{\mathfrak a}}(X_{\mathfrak m},f_B,{\mathfrak D})$ (Equation 4.4) are unbiased and consistent finite sample estimates of $Imp_{\infty}^{mda}(X_m)$.

On the other hand, while Equation 4.6 only depends on the joint distribution between Y and X, the "Bagging" estimate of Equation 4.5 also depends on the base learner Algo used. The consistency of ${\rm Im} \mathfrak{p}_{Algo}^{\mathfrak{m} d \mathfrak{a}}$ with respect to ${\rm Im} \mathfrak{p}_{\infty}^{\mathfrak{m} d \mathfrak{a}}$ thus depends on the properties (and obviously the consistency) of the base learner. In particular, [Gregorutti et al., 2017] note that this consistency was shown by Zhu et al. [2015] under several hypotheses, including the use of purely random forests Biau et al. [2008] and the independence between features⁸.

Additive regression model.

To handle the complexity of the theoretical analysis of the MDA importance measure, [Gregorutti et al., 2017] consider the particular case of a joint distribution P satisfying the following additive regression model

$$Y = \sum_{j=1}^{p} f_{j}(X_{j}) + \epsilon$$
 (4.7)

where ϵ is such that $\mathbb{E}\{\epsilon|X\}=0$ and $\mathbb{E}\{\epsilon^2|X\}$ is finite (and where all functions f_i are measurable) implying that $f_B(x) = \sum_{j=1}^p f_j(x_j)$.

In this setting, Gregorutti et al. [2017] show that the MDA importance of a variable X_{m} is

$$Imp_{\infty}^{mda}(X_m) = 2var\{f_m(X_m)\}. \tag{4.8}$$

$$\mathbb{E}_{\tilde{P}_{\mathfrak{m}}}\left\{(Y-f_{B}(X))^{2}\right\}=\mathbb{E}_{P}\left\{\mathbb{E}_{\tilde{X}_{\mathfrak{m}}\sim P(X_{\mathfrak{m}})}\left\{(Y-f_{B,\tilde{X}_{\mathfrak{m}}}(X))^{2}\right\}\right\},$$

where $f_{B,\tilde{X}_m}(X)$ returns the value of f_B at \tilde{X}^m obtained from X by replacing X_m by \tilde{X}_m and leaving all other features unchanged. Inverting the two expectations, one gets:

$$\mathbb{E}_{\tilde{X}_{\mathfrak{m}} \sim P(\tilde{X}_{\mathfrak{m}})} \left\{ \mathbb{E}_{P} \left\{ (Y - f_{B, \tilde{X}_{\mathfrak{m}}}(X))^{2} \right\} \right\}.$$

By definition of f_B, the inner expectation, and thus also the outer expectation, is greater or equal to

 $\mathbb{E}_{P}\left\{(Y-f_{B}(X))^{2}\right\}\!\text{, which proves that }\mathrm{Imp}_{\infty}^{\mathfrak{mda}}(X_{\mathfrak{m}})\text{ is non-negative.}$ $^{7}\mathrm{The \ two \ terms \ in \ Equation \ 4.3 \ are \ indeed \ unbiased \ and \ consistent \ sample \ estimates \ of \ the \ two}$ population mean square errors in 4.6.

⁶More formally, we can rewrite the first term of 4.6 as

⁸This assumption is quite strong and excludes works on correlated features for instance.

Equation 4.8 states that the MDA importance of a feature is (twice) the variance of the contribution $f_m(X_m)$ of X_m in the additive Bayes model (Equation 4.7). In the classification setting, Gregorutti et al. [2017] show that this result is not valid with zero-one loss in the case of an additive logistic regression model, as they note that $\mathrm{Imp}_{\infty}^{\mathrm{mda}}(X_m) > 0$ only if the contribution of X_m to P(Y|X) is large enough to change the predicted class.

Zhu et al. [2015] use a slightly different notion of population importance, which is a normalised version of

$$\mathbb{E}\{(f_B(X) - f_B(\tilde{X}^m))^2\}$$

where \tilde{X}^{m} denotes the vector of inputs where the mth coordinate was replaced by an independent copy of X_{m} and the expectation is taken with respect to the joint distribution of Y, the original inputs X, and the independent copy of X_{m} . Under the above additive model, this definition actually coincides with the former notion introduced above, as shown by [Gregorutti et al., 2017].

Simplified permutation scheme.

Instead of considering a specific model and still circumventing the complexity of the permutation scheme, Ishwaran [2007] study a variant of MDA importance sharing similar key properties but implementing another permutation scheme. Instead of permuting the values of a feature $X_{\rm m}$ in oob samples, Ishwaran [2007] propose to "noise up" the feature $X_{\rm m}$ by ignoring all nodes coming after one splitting on $X_{\rm m}$. In practice, it comes to a random left-right assignment of samples in all ignored nodes. The beginning of the tree however remains unchanged. For this setting and assuming that the model can provide a good approximation⁹, the asymptotic behaviour of this variant can be derived.

In particular, [Ishwaran, 2007] focus on the *position bias* and show that variables split close to the root node tend to have a stronger effect on the predictive accuracy than other variables. It seems reasonable that the model performances are highly impacted as most of the tree is ignored when evaluating the importance of a feature close to the root. A similar behaviour is expected in the classical MDA importance. Indeed, the relation between features used at the top of the tree structure and their expected usefulness is obvious.

Nevertheless, some irrelevant features may appear as important in this variant because of the feature noising. Since all nodes are ignored after one splitting on the evaluated feature $X_{\rm m}$, the observed decreases in predictive accuracy is not only due to $X_{\rm m}$ but also to all features used in deeper nodes. Therefore, the importance of $X_{\rm m}$ reflects both the actual contribution of $X_{\rm m}$ and the contribution of all split variables of ignored nodes. The importance of $X_{\rm m}$ can thus be strictly positive even if $X_{\rm m}$ is irrelevant. In response to that, Ishwaran [2007] suggest that non-informative features are more likely used down in trees and thus spurious importance scores should be limited. He also claims that noising up only the right node (i.e., the one using $X_{\rm m}$ to split) is too difficult to be theoretically analysed without additional assumptions.

⁹In details, in asymptotic conditions, the tree-based model must be able to provide a good approximation of the true inputs-output function which implies the consistency of the model and the piecewise constance of the regression function [Ishwaran, 2007].

4.3.2 Asymptotic properties of MDI

Regression tree-based models

According to Friedman [2001], the MDI importance measure is an approximated measure of the relative influence of variables. In the context of regression problems, let us consider a given predictor $f(\cdot) \in \mathcal{Y}^{\mathcal{X}}$. Following Friedman [2001], the relative importance of an input variable X_j in the predictor f is its relative influence on the variation of f over the joint input variable distribution and computed as follows

$$\operatorname{Imp}_{f}^{\inf}(X_{j},f) = \sqrt{\mathbb{E}_{X}\left\{\left(\frac{\partial f(X)}{\partial X_{j}}\right)^{2}\right\} \cdot \operatorname{var}_{X}\{X_{j}\}}.$$
(4.9)

Friedman [2001] note that Equation 4.9 does not strictly exist for piecewise constant functions such as produced by regression tree-based models. Friedman [2001] therefore suggest that the MDI importance measure 10 of X_j was proposed as a surrogate measure to approximate Equation 4.9 for piecewise constant functions and shown to be consistent with expected feature influences in the case of linear relationships between inputs and output variable [Friedman, 2001].

Beyond this intuitive motivation, we now turn to classification problems, and analyse the main properties of the MDI importance measure when it is based on the Shannon entropy as an impurity measure.

Totally randomized decision-tree based ensembles with categorical input features and multiway exhaustive splits.

Following Louppe et al. [2013]; Louppe [2014], let us consider a set $V = \{X_1, \ldots, X_p\}$ of categorical input features and a categorical output Y. For the sake of simplicity, only the Shannon impurity is considered below but most results can be go generalised to other impurity measures [Louppe et al., 2013; Louppe, 2014]. Let us also consider totally randomized trees (defined in Section 3.3) with multiway exhaustive splits (see Section 3.2). In case of categorical variables, each node t is split into $|X_i|$ sub-trees, i.e., one for each possible value of X_i . It implies that features can only be used once and thus limits the depth of a branch to p.

In this setting, the MDI importance of feature $X_m \in V$ for Y computed in asymptotic conditions¹¹ is given by [Louppe et al., 2013]:

$$Imp_{\infty}^{mdi}(X_{m}) = \sum_{k=0}^{p-1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-m})} I(X_{m}; Y|B)$$
 (4.10)

where V^{-m} denotes the subset of features $V\setminus\{X_m\}$, $\mathcal{P}_k(V^{-m})$ is the set of subsets of V^{-m} of cardinality k, and $I(X_m;Y|B)$ is the conditional mutual information of X_m and Y given the variables in the conditioning set B. Additionally, Louppe et al. [2013] show that

$$\sum_{m=1}^{p} Imp_{\infty}^{mdi}(X_m) = I(X_1, \dots, X_p; Y)$$
 (4.11)

 $^{^{10}}$ Actually, the MDI importance computed as the sum of empirical improvement in squared error over all nodes splitting on X_i in a given tree and its average over all trees.

¹¹Infinite learning sample size, infinite ensemble of fully developed (ie., unpruned) totally randomised trees.

where $I(X_1, ..., X_p; Y)$ is the joint mutual information between all features in V and the output Y.

Equation 4.10 shows that each importance can be divided along the interaction degree k, i.e., the number of features in the conditioning set B, and along the combinations of B of fixed size of k features.

Equation 4.11 states that all the information $I(X_1, ..., X_p; Y)$ contained in the set of input variables V about the output Y can be decomposed between the importance of all features. The equality of Equation 4.11 induces that the sum of all importances equals a fixed value (of the joint mutual information). It implies that the increase or decrease of one feature importance is made to the detriment of other importances.

Let us mention that any (conditional) mutual information term involving Y (of the form I(X;Y|B) or $I(X_1,\dots,X_q;Y|B)$ with B potentially empty) is upper bounded by H(Y). It gives in particular that $\sum_{m=1}^{p} Imp_{\infty}^{mdi}(X_m) \leqslant H(Y)$ where the equality indicates that Y is perfectly explained by V (i.e., $I(X_1,...,X_p;Y) = H(Y)$).

Louppe et al. [2013] show also that the form of these expressions remains valid for any impurity measure leading to non negative impurity decreases, including obviously all classical impurity measures such as Shannon-, Gini-, and variance-based ones.

Non-totally randomized trees with multiway exhaustive splits and categorical input features.

Beyond its asymptotic behaviour, [Louppe et al., 2013; Louppe, 2014] establish a relationship between relevance and MDI importance. This relationship follows from the definition of relevance in terms of mutual information (see Definitions 2.7 and 2.8 in Section 2.4.1).

In what follows, results can be extended to MDI importances derived from nontotally randomised trees (i.e., with K > 1). Thus, let us denote the MDI importance computed with totally or non-totally randomized trees depending on the value of K as ${\rm Imp}_{\infty}^{{\rm mdi},1}$ and ${\rm Imp}_{\infty}^{{\rm mdi},K}$ respectively.

In this context, a feature X which is irrelevant for Y with respect to V always verifies $Imp_{\infty}^{mdi,K}(X) = 0$ [Louppe et al., 2013; Sutera et al., 2018]. In case of totally randomised trees (K = 1), a null score is only associated to an irrelevant feature and consequently all relevant features (strongly and weakly) have strictly positive MDI importance scores. Additionally, this result implies that irrelevant features do not impact importance scores of other features. Consequently, the relevant feature MDI importances are thus independent of the number of irrelevant features.

On the contrary, with non-totally randomised trees (K > 1), some relevant features can also have a zero importance score due to the effect of K on the tree construction. Sutera et al. [2018] show that only strongly relevant features are guaranteed to have strictly positive MDI importance score as they convey information about the output that no other variable (or combination of variables) in V conveys Depending on the value of K, some weakly relevant features may have a zero importance score. The randomisation parameter K (when > 1) thus affects the number and nature of relevant variables that can be found.

In the same conditions, [Louppe et al., 2013] also show that the MDI importance derived from pruned trees (i.e., built up to a depth q < p) is equivalent to the ones obtained from unpruned trees built on random subspaces of q variables randomly drawn from V.

4.3.3 Correlated and redundant features

By definition, totally redundant features share exactly the same information about the target variable Y, while correlated features often share information without necessarily being totally redundant with respect to Y. Tree-based or model-based importance measures described so far evaluate the contribution of a feature in the tree-based predictor or in the Bayes model. In the presence of redundant or correlated features, the sum of all contributions can no longer be shared unequivocally between all features. For example, the same "piece" of contribution might be attributed to several totally redundant features as they are interchangeable in the eyes of the model. The rest of this section describes works focusing on that aspect of importance measures.

4.3.3.1 MDA

Additive regression model with centred $f_i(X_i)$ functions.

Gregorutti et al. [2017] continue their theoretical study of the additive model, by analysing the population version of the MDA importance in terms of feature correlations, assuming in addition that all $f_j(X_j)$ functions have zero mean. Under these conditions, Equation (4.7) becomes 12

$$Imp_{\infty}^{md\alpha}(X_m) = 2cov\{Y, f_m(X_m)\} - 2\sum_{k\neq m} cov\{f_m(X_m), f_k(X_k)\} \tag{4.12} \label{eq:4.12}$$

where cov denotes the covariance function. In this alternative formulation, interactions between input features are explicitly shown in the second term.

Additive regression model and a normal distribution.

Gregorutti et al. [2017] further consider the case of normal joint distribution $P_{V,Y} \sim \mathcal{N}_{p+1} (0,Z)$ with a group C of c features $\{X_1,\ldots,X_c\}$ equally correlated with each other and with the output. In order to highlight relationships between block of features, the covariance matrix Z can be expressed as follows

$$Z = \begin{pmatrix} Z_{V} & \boldsymbol{\tau}^{T} \\ \boldsymbol{\tau} & \sigma_{y}^{2} \end{pmatrix} = \begin{pmatrix} \rho & 0 & \boldsymbol{\tau}_{\in C}^{T} \\ 0 & 1 & \boldsymbol{\tau}_{\notin C}^{T} \\ \boldsymbol{\tau}_{\in C} & \boldsymbol{\tau}_{\notin C} & \sigma_{y}^{2} \end{pmatrix}$$
(4.13)

where

- Z_V is the covariance between input features;
- ρ is the covariance sub-matrix $c \times c$ of features in the correlated group such that $cov(X_i, X_i) = 1$ and $cov(X_i, X_j) = \rho$ for all $1 \le i, j \le c$, i.e. $\rho = (1 \rho)I_c + \rho \mathbb{1}\mathbb{1}^T$;
- $\tau_{\in C}$ is a (line)vector of c elements $\tau_{\in C} = \{\tau_C, \dots, \tau_C\}$, i.e. $cov\{X_m, Y\} = \tau_C$ with $0 < m \leqslant c$;
- $\tau_{\not\in C}$ is a (line)vector of (p-c) elements $\tau_{\not\in C} = \{\tau_{c+1}, \ldots, \tau_p\}$, i.e. $cov\{X_j, Y\} = \tau_i$ with c < j < p;

¹²See [Gregorutti et al., 2017, Proposition 2] for a proof; the zero-mean assumption is not essential but simplifies the reading of the expression.

• and σ_{u}^{2} is the variance of Y.

In this setting, Gregorutti et al. [2017] specify the MDA importances as follows:

$$Imp_{\infty}^{md\alpha}(X_m) = 2\alpha_m^2 var\{X_m\} = 2\alpha_m cov\{X_m,Y\} - 2\alpha_m \sum_{k\neq m} \alpha_k cov\{X_m,X_k\} \text{(4.14)}$$

where α 's are deterministic coefficient¹³ equal to $\alpha_m = [Z_V^{-1}\tau]_m$.

For a feature $X_i \notin C$, Equation 4.14 becomes

$$Imp_{\infty}^{mda}(X_j) = 2\tau_j^2 \tag{4.15}$$

where τ_i corresponds to $cov\{X_i, Y\}$.

For a feature $X_i \in C$, Equation 4.14 becomes

$$Imp_{\infty}^{mda}(X_i) = 2\left(\frac{\tau_C}{1 - \rho + c\rho}\right)^2$$
(4.16)

where $\tau_C = cov\{X_i, Y\}$ and $\rho = cov\{X_i, X_k\}$ with $0 \leqslant k \leqslant c, k \neq i$. In the particular case of two copies of the same feature, i.e. c = 2 and $\rho = 1$, it is

$$\operatorname{Imp}_{\infty}^{\operatorname{mda}}(X_{i}) = 2\left(\frac{\tau_{C}}{2}\right)^{2} = \frac{\tau_{C}^{2}}{2}.$$
(4.17)

Equation 4.15 states that the importance of a non-correlated feature is not impacted by potential correlation between other features. Equation 4.16 shows that the importance of a feature correlated with others is influenced by ρ and c. A large number of correlated features c or a strong correlation, i.e. c close to 1, decrease the MDA importance of each individual feature. Combining Equations 4.15 and 4.16 suggests that X_i may appear more important, i.e. corresponds to a higher MDA importance, than X_i even if $\tau_j < \tau_C$ if ρ is large enough. Conversely, anti-correlation ρ < 0 tends to increase the MDA importance.

4.3.3.2 MDI

Totally randomized trees with multiway exhaustive splits and categorical input fea-

Let $X_i \in V$ be a relevant variable with respect to Y and V and let $X_i' \not \in V$ be a new variable such that X_j and X_i' are totally redundant with respect to Y (see Definition 2.17). Louppe [2014] extends the analytical formulation of the MDI importances of X_i and any non-redundant variable $X_i \in V^{-j}$ in order to show the impact of the addition of X_i' . For sake of clarity, only one pair of totally redundant features is considered but see [Louppe, 2014] for a generalisation to c such features.

The asymptotic importance of variable X_i as computed from an ensemble built on $V \cup \{X_i'\} \text{ is}^{14}$:

$$\operatorname{Imp}_{\infty}^{\operatorname{mdi},1}(X_{j}) = \sum_{k=0}^{p-1} \frac{p-k}{p+1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-j})} I(X_{j}; Y|B)$$
(4.18)

¹³See [Gregorutti et al., 2017, Proposition 3] for a proof.

¹⁴See [Louppe, 2014, Proposition 7.2] for a proof.

For any other variable X_1 from V^{-j} , the importance becomes 15

$$\begin{split} & Imp_{\infty}^{mdi,1}(X_{l}) = \sum_{k=0}^{p-2} \frac{p-k}{p+1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-l} \setminus \{X_{j}\})} I(X_{l}; Y | B) \\ & + \sum_{k=0}^{p-2} \left[\sum_{k'=1}^{2} \frac{C_{2}^{k'}}{C_{p+1}^{k+k'}} \frac{1}{p+1-(k+k')} \right] \sum_{B \in \mathcal{P}_{k}(V^{-l} \setminus \{X_{j}\})} I(X_{l}; Y | B \cup \{X_{j}\}) \end{split} \tag{4.19}$$

A comparison of Equations 4.18 and 4.10 shows that the introduction of a variable X_j' totally redundant with X_j decreases the importance of X_j . Indeed, with respect to 4.10, all terms of the sum in 4.18 are multiplied by a factor $\frac{p-k}{p+1} < 1$. Intuitively, this is a consequence of the fact that both X_j and X_j' convey the exact same information about the output and they now both compete to explain the output, as the sum of all importances is not affected by the introduction of X_j' . Indeed, X_j' does not bring any new information about the output with respect to X_j (by definition) and therefore the right side of Equation 4.11 is unchanged. Although we obviously have $\mathrm{Imp}_{\infty}^{\mathrm{mdi},1}(X_j) = \mathrm{Imp}_{\infty}^{\mathrm{mdi},1}(X_j')$ by symmetry, notice that the importance of X_j is not simply divided by a factor 2 since the importances of the other variables are also affected by the introduction of X_j' , as shown in Equation 4.19.

Equation 4.19 shows that the impact of the introduction of X_j' on the importances of the variables in V^{-j} is the combination of two effects. The first sum in 4.19 is over all B composed of variables from V^{-j} . With respect to the corresponding terms in 4.10, each term is multiplied by a factor $\frac{p-k}{p+1}$ strictly lower than 1. The second sum in 4.19 is over all conditionings including X_j and the weights of the corresponding terms are now increased with respect to similar terms in 4.10. Whether or not the importance of X_1 will increase will thus depend on the way X_1 interacts with X_j . If the mutual informations $I(X_1;Y|B,X_j)$ are large $(X_1$ and X_j are complementary), then adding X_j will reinforce these terms and the net effect could be an increase of the importance of X_1 . On the other hand, if these mutual informations are small $(X_1$ and X_j are redundant), the net effect could be a decrease of the importance of X_1 .

4.4 EMPIRICAL ANALYSES

In the previous section, we studied theoretically both importance measures in asymptotic conditions. Although those results are helpful to better understand the mechanisms of MDA and MDI importance measures, they do not provide insights on how they actually behave in practice. In the light of their expected behaviours, the goal of this section is to analyse those two measures in a more realistic setting, i.e. with finite sample size and number of trees. To do so, we review many empirical analyses of their practical behaviours in numerous settings. In particular, we aim at highlighting the main biases and practical limitations of MDI and MDA importance measures in several view angles.

4.4.1 Soundness

Variable importance measures derived from tree-based ensemble methods have been suggested for the identification and selection of relevant features in numer-

¹⁵See [Louppe, 2014, Proposition 7.4] for a proof.

ous applications, e.g. gene selection in micro-array data [Huang et al., 2005; Díaz-Uriarte and De Andres, 2006; Pang et al., 2006; Rodenburg et al., 2008], SNPs in large-scale/genome-wide association study data (GWAS) [Lunetta et al., 2004; Bureau et al., 2005; Botta et al., 2014], proteins [Qi et al., 2006], or, more recently, brain regions involved in neuronal disease in neuroimaging data [Wehenkel, 2018]. Along with this wide practical use, some works have tried to assess the quality of this identification.

Archer and Kimes [2008]; Grömping [2009] show that MDI and MDA feature importance measures manage to identify true predictors in different settings, and results are usually in agreement with other machine learning methods.

In presence of feature interactions, it was also noted that these measures provide interesting alternatives to classical statistical tests because they do not require explicit modelling or assumptions on the problem (e.g., gaussianity, (non-)linearity, or independence) and naturally handle feature interactions [Grömping, 2009; Geurts et al., 2009]. Differences between univariate approaches and tree-based importance scores may additionally be indicative of multivariate interactions [Rodenburg et al., 2008; Auret and Aldrich, 2011]. For example, Lunetta et al. [2004] show that selections of relevant genetic markers (SNPs) provided by random forest feature importance measures outperform those obtained from a standard univariate screening method (i.e., Fisher Exact test), especially in presence of many interacting features.

In presence of correlated features, Archer and Kimes [2008] showed, in a setting similar to Gregorutti et al. [2017]'s (i.e., one group of correlated and equally predictive features, see Section 4.3.3.1), that both Gini MDI and MDA importance measures manage to identify most predictive features in many settings. They however noted that in case of strong correlation (ρ close to 1), the highest importance score may be associated to one feature correlated with the most predictive one. When there were more than one group of predictive correlated features or uncorrelated predictive features, some experiments show that both importance measures are sensitive to correlation structures and this may sometimes impact the reliability and stability of importance scores [Strobl et al., 2008; Nicodemus and Malley, 2009; Toloşi and Lengauer, 2011; Auret and Aldrich, 2011]. Depending on tree parameters and correlation structures, empirical observations seems to diverge. Therefore, a more detailed analysis of those experimental results will be the focus of Section 4.4.5.2.

From another point of view, Lundberg and Lee [2017]; Lundberg et al. [2018] claim that MDI importance measure is not "consistent" in the case of a (non-randomised) single tree. In the chosen example of two equally relevant features, increasing the predictive contribution of one does not necessarily correspond to an increase of its MDI importance. Conversely, the MDA importance measure appears to be "consistent" in this example.

4.4.2 Split randomisation parameter K

In random forest methods, K is the number of features considered at each node as split variable candidates. A low value of K (e.g., K = 1) maximises randomisation as one feature is selected totally at random without optimising the node impurity reduction. Consequently, all features can be selected and all relevant features may be identified. In contrast, high values (e.g., K = p) induce more optimised trees and only strongly relevant features are guaranteed to be identifiable.

The interaction between K and feature importance measures is not clear. For several authors [Auret and Aldrich, 2011; Strobl et al., 2008; Nicodemus et al., 2010], importance measures are more accurate when derived from ensemble of trees built with large K values. In these studies, experiments are carried out on simulated data where the output is a linear combination of several features, i.e. $Y = \alpha_1 X_1 + \alpha_2 X_2 + \cdots + \alpha_p X_p$ where non-zero coefficients correspond to predictive features while zero coefficients refer to non-predictive ones. Additionally, some features may be correlated, possibly in a strong fashion. In this setting, a feature importance measure is said to be inaccurate if it provides importance scores that do not comply with the α_i coefficients of the true model. Below, we argue that feature importance measures should not be necessarily considered as less accurate for low values of K because importance scores do not align with these coefficients, especially when correlated features are not equally contributive in the linear combination as it is the case in their analyses. In particular, as explained above, low values of K might be more appropriate to address the all-relevant problem, even if this leads to importances that do not match coefficients α . Results in these papers are also of interest to discuss biases in feature importance measures due to correlation and we analyse them with this different angle in Section 4.4.5.2.

It should be noted that a non predictive feature X_i ($\alpha_i = 0$) that is strongly correlated with a predictive feature X_i ($\alpha_i > 0$) may therefore be weakly relevant to the target as it may provide part of the information of X_i about Y. Nicodemus et al. [2010] characterised such features that appears to be predictive as long as some other features are not included in the model as "spurious correlation". In our terminology, feature X_i is weakly relevant and totally redundant to X_i with respect to the target. Consequently, coefficients α do not reflect the actual contribution of each feature in a tree-based model.

Authors adopting the minimal-optimal point of view for feature selection (like those mentioned above) concentrate their efforts on identifying only a part of relevant features (i.e., strongly relevant features and a maximal subset of non-redundant ones). It therefore makes sense that redundant features are expected not to be identified as important. However, except in trees built without node-wise split randomisation (i.e., K = p), even totally redundant and weakly relevant features can be selected in tree models if they do not compete at some nodes with features that are most useful (and eventually provide the same information). This explains why [Auret and Aldrich, 2011; Strobl et al., 2008; Nicodemus et al., 2010] observe that feature importance measures seem more accurate for high values of K even if low values of K would be more appropriate when interested in solving the all-relevant problem. In such cases, theoretical results (from [Sutera et al., 2018] and summarised in Section 4.3.2) confirm that high values of K imply that redundant features are more frequently masked by strongly relevant features (with positive coefficients) and therefore importance scores are more similar to coefficients α . Strongly ("truly") relevant features are also expected to be used more often and to recover most of the importance in the tree model. Genuer et al. [2010] indeed observed experimentally that higher values of K increase the importance of truly important variables.

In addition, low-sample conditions imply that only few variables can be evaluated before reaching nodes with too few samples for an accurate impurity estimation (see Section 4.4.5.5). Increasing the value of K may actually improve importance scores for relevant features that are more often chosen near the root. They are estimated

more often and with more samples, potentially making them more stable and more accurately estimated.

Similarly, in presence of many irrelevant features, using a small value of K may induce that numerous splits are made on irrelevant features (because all split variable candidates are irrelevant). On one hand, such splits do not provide information about the target. On the other hand, a feature is selected based on its spurious relationship with the output and is unfairly credited of some importance for it. Less randomised trees (i.e., K close to p) are therefore preferable in such situations. In contrast, if all features are assumed to be equally relevant, then more randomised trees (K close to 1) are more suitable because they consider all features and not just some of them.

From all those observations, a trade-off for the value of K needs to be found in order to identify the right set of relevant features while taking into account the nature of the problem.

4.4.3 Feature ranking stability and number of trees

Typically, the number of trees necessary for good performances grows with the number of features [Liaw et al., 2002]. There is no need to grow more trees when the predictions of a subset of the forest are as good as the predictions of the whole forest. This approach however requires to build an unnecessary large number of trees. Therefore, several works propose simple procedure to determine a priori the number of trees for stable and accurate predictions [Latinne et al., 2001; Hernández-Lobato et al., 2013]. However, these only concern the predictive ability of tree-based ensemble and the number of trees may not be optimal with respect to the feature importance measures. In [Huynh-Thu et al., 2012; Paul et al., 2012], experiments show that the numbers of required trees yielding stable feature selection and predictive performances differ from several orders of magnitude.

Theoretically, feature importance measures only attribute zero importance scores for irrelevant or masked features. However, in practice, this property relies on one fundamental principle: the number of trees is large enough. Indeed, as pointed in [Wehenkel et al., 2018], in case of too small trees and/or high-dimensional datasets (p >> N), some features may have a zero importance value because they never have been considered during the tree growing process. Additionally, some feature importance may have been evaluated in too few occasions to fairly represent its true contribution. For example, two features forming a XOR structure need to be used at least two times such that both features can be used once before each other. Ultimately, one expect that their averaged importances over a sufficient number of evaluations is the same for both features. In that context, Wehenkel [2018] uses the idea of the so-called coupon collector's problem and derives a minimal number of trees (for given parameters N, p and K) that should be built to have some minimum guarantee that all features are seen at least once.

Even if all features have been considered and receive an importance score, the interpretation of feature importance measures is only possible if results are stable enough, i.e., do not vary significantly if a few additional trees are taken into account, for another ensemble of same size or if small changes are made to the dataset [Strobl et al., 2008; Saeys et al., 2008b]. Typically, it has been suggested and observed that increasing the number of trees in the forest improves the stability of feature importance measures [Liaw et al., 2002; Archer and Kimes, 2008; Genuer et al., 2010; Paul et al., 2012]. In practice, Liaw et al. [2002] however observed that importance scores may vary from one ensemble to another while ranking of importances is usually more stable for the same number of trees. In a discussion about stability of ranked gene lists (which aims at identifying a short-list of genes of interest for further analyses), Boulesteix and Slawski [2009] state that the rank of a particular feature is usually as important as its value from a practical point of view. Saeys et al. [2008b] note that the analyses of selected features typically require much effort and time and this stresses the need for a stable feature ranking and robust feature selection techniques, especially for model interpretation in biomedical applications [Tolosi and Lengauer, 2011].

Assuming enough trees and a stabilised feature ranking, it appears in several data sets that the most important features have typically the highest importance scores [Auret and Aldrich, 2011]. This also suggests that efficient feature selection can be performed by selecting the best k features, where k can be determined by selecting a judicious importance thresholds so as to minimise the number of selected irrelevant features (false positive). Section 4.5 focuses on approaches proposed in the literature to determine this threshold. However, a stable feature ranking does not imply that importance scores are reliable, i.e. that one feature better ranked than another is not necessarily more important. Feature importance measures may be sensitive to different factors, such as the presence of correlated features, and provides unfair importance scores. In Section 4.4.5, we review the main sources of unfairness (biases) that have been studied in literature.

4.4.4 Importance measures vs prediction performances

Tree-based feature importance measure is usually seen as a side-product of the random forest model. However, a model optimised so as the maximise its performances is typically not adjusted for measuring feature importances [Van der Laan, 2006]. For example, Paul et al. [2012] show that the number of trees yielding stable prediction performances is smaller of several orders of magnitude than what is required for a stable feature selection. The number of trees should then be carefully chosen. In relation with Section 4.4.2, randomisation parameter K is usually considered as crucial to obtain good accuracy performances, by controlling the randomisation of the model (and thus the bias-variance trade-off). In classification (respectively, in regression), empirical studies typically suggest that $K = \sqrt{p}$ (resp., K = p) is an appropriate and often optimal value with respect to prediction accuracy [Geurts et al., 2006; Strobl et al., 2008]. It has however been noticed that model performances is usually not related to the goodness of tree ensemble parameters for variable importance purposes [Auret and Aldrich, 2011; Huynh-Thu et al., 2012]. Theoretical results suggest that low values of K are more suitable for feature importance measures as K = 1 is the only way to guarantee that all relevant features can be identified, but this usually requires a larger number of trees to consider all features. Conversely, higher values of K tend to focus more on strongly relevant features. In terms of prediction accuracy, larger values (e.g., $K = \sqrt{p}$ or K = p) are more suitable, especially in presence of many irrelevant features, to avoid useless but will definitely prevent some weakly relevant features to be identified. As a result of this discussion, one should carefully choose tree-based parameters and find an appropriate trade-off between feature importance measures (selection or ranking) and prediction performances.

4.4.5 Biases

In what follows, we discuss some experimental results that reveal the presence of biases that affect one or both importance measures. In this work, an importance measure is biased if its use in practical conditions differs from its expected and theoretical behaviour. In particular, it is biased if it does not equally treat similar variables, i.e. it does not attribute the same importance score to all features that are equally relevant (or irrelevant) [Dobra and Gehrke, 2001]. For example, let us consider two features that are completely independent of the output and are thus irrelevant. An unbiased measure would attribute the same score for both variables while a biased one may have a systematic preference for one of them resulting in a higher importance score.

We refer to an importance score over-estimation (respectively, under-estimation) as a positive bias (resp. negative bias). For example, an importance measure that gives a positive score to an irrelevant feature, that should receive a zero importance, is positively biased.

As a preamble, let us note that MDA feature importance measure relies on the tree structure that has been induced using an impurity criterion. Therefore, some biases that affect impurity measures and thus MDI importance measures, may sometimes also affect MDA. For example, if a feature is never selected because it produces for some reasons no impurity decrease, its permutation does not change the accuracy performances of the model. Conversely, it is also possible that MDA importance measure reduces the importance of features that have been unfairly selected. For the sake of example, let us imagine a bias favouring the selection of redundant features, each providing strictly positive impurity decreases (partly due to noise). Permuting the value of one variable may be ineffective on the prediction of the model, yielding to a null MDA importance scores while the corresponding MDI value might be slightly higher.

4.4.5.1 Bias due to masking effect

Source of bias: tree-based method randomisation parameter K.

Masking effect was already mentioned in several occasions in this thesis as a consequence of non-totally randomised split variable selections. In Section 3.2.3, we showed that the inversion between masked and masking features by the means of a small change in the learning set can induce totally different decision tree model (and non randomised), illustrating the high variance of the decision tree algorithms. In Sections 4.3.2 and 4.4.2, we highlighted that large values of K increase the range of masking effect, resulting in giving preference to strongly relevant features that can not be masked to the detriment of weakly relevant features. The masking effect is maximal when K = p. In this section, we discuss the impact of the masking effect on importance measures.

The masking effect denotes situations where several candidate splits on different variables yield roughly the same impurity reduction, but one is always slightly better so that none of the other ones has a chance to be selected by the tree-growing algorithm. Concretely, some branches are never explored as splits are never selected. This induces a positive bias for importances of masking features as they are more frequently selected and their contributions is prioritised over features carrying similar information about the target, i.e., in case of two redundant features with one

masking the other, the first one always receives credit for its information because the second one is never selected before. In contrast, importance of masked features are negatively biased and under-estimated. Let us note that this bias impacts both importance measures as it affects the building of the tree models.

A straightforward way to reduce this bias is to reduce the value of K. This bias can be totally removed by using totally randomised trees (K = 1) but this usually requires to increase the number of trees and might jeopardise the predictive performance of the model in presence of many irrelevant features. However, in order to reach global optimality of the ensemble [Strobl et al., 2008], it may also be necessary to unveil some feature interactions (e.g., cliques where features are marginally irrelevant and thus unlikely to be selected at first sight) or feature importances (e.g., the second feature in an imbalanced XOR¹⁶).

4.4.5.2 Bias due to correlation

Source of bias: presence of correlated features in learning samples.

Random forest methods are popular in many scientific fields for their ability to handle high-dimensional datasets, as it is particularly the case in biomedical applications. In addition, it is quite common in biomedical studies that features are strongly correlated with each other and this strong correlation usually has a biological explanation. For example, co-regulated genes in expression data are expected to be similar as they relate to the same molecular pathway [Tolosi and Lengauer, 2011]. Neighbouring pixels/voxels in biomedical images are likely associated to the same biological entities (e.g., neurons) implying a spatial correlation [Wehenkel et al., 2018]. These examples have motivated several empirical studies of feature importance measures in presence of correlated features.

We however need to distinguish two different biases due to correlation that have been identified in the literature: a preference for correlated features with respect to uncorrelated ones and a preference for correlated groups of smaller sizes. In what follows, let us note that the correlation structure is not the same in both parts. All features in a group share the same predictive power to study the effect of the size of correlated feature groups [Tolosi and Lengauer, 2011] while features within the same group can vary in their information about the target in order to highlight preference for correlated features [Strobl et al., 2008; Nicodemus et al., 2010].

PREFERENCE FOR (UN)CORRELATED FEATURES In their experimental studies, Strobl et al. [2008]; Nicodemus and Malley [2009]; Nicodemus et al. [2010] analyse feature importance measures in presence of correlated features that are not equally contributive in the prediction of the output. Several effects are observed in those studies.

Gini MDI importance measure appears to be biased in the presence of correlation [Nicodemus and Malley, 2009]. Strobl et al. [2008] observe that correlated features are positively biased with MDA feature importance measure. Strobl et al. [2008]; Nicodemus et al. [2010] report that correlated features are more frequently selected at the first split of the tree (when K > 1). Nevertheless, across all splits, Nicodemus et al. [2010] observe a slight preference for selection of uncorrelated features. Most

¹⁶An imbalanced XOR is the example used in Section 3.2.3. Two features form a XOR but one is always slightly more marginally relevant and is thus always selected first, obtaining therefore a lower importance score than the other one.

of these results are studied for different values of K, including totally randomised trees with K=1 but excluding non-randomised trees with K=p. A first observation is that a correlated feature with zero coefficient in the generating model (see Section 4.4.2 for the description) ends up with larger importance than uncorrelated features with zero coefficient. For Strobl et al. [2008], this phenomenon is due to a spurious correlation that makes a zero coefficient feature marginally informative but conditionally useless. However, such feature carrying redundant information is actually weakly relevant and thus might be selected and contribute to the model. Because of randomisation, it may occur that those features are evaluated without being in competition with their correlated features and end up being selected at some nodes. Such situations are expected to be less likely when the level of randomisation decreases, as observed in those studies with an increasing K. Moreover, correlation does not necessarily imply redundancy (as shown in Section 2.4.3.6 and in [Guyon and Elisseeff, 2006]) and it may slightly increase the predictive contribution of some correlated features with respect to uncorrelated ones with similar coefficients, making them more frequently selected. Simultaneously, non-predictive features that are weakly relevant because of correlation necessarily provide redundant informations. If those features are selected, it reduces the potential interest of selecting correlated features in subsequent nodes in favour of uncorrelated features.

In conclusion, we believe that some of these observations are not actually directly due to the presence of correlation but consequences of masking effect (and the preference for strongly relevant features with high K values) and weakly relevance of features with zero coefficient that benefits from their correlation with highly informative features. Furthermore, Nicodemus and Malley [2009] noticed that prepruning trees by limiting node-size tends to reduce the effect of bias. Therefore, part of observed effects may actually be due to other reasons, such as empirical impurity misestimations in nodes with too few samples.

PREFERENCE FOR SMALLER GROUPS OF CORRELATED FEATURES In many biomedical applications, all features within a correlated group are roughly equivalent (e.g., neighbouring voxels in neuroimaging) and can typically be used interchangeably yielding equally performing tree-based models. One can thus associate a group of correlated features with a certain contribution in the prediction of the output.

Theoretical results, especially MDI importance of totally redundant features (see Section 4.3.3.2), suggest that if features are equivalent¹⁷, they are expected to be equally informative and the importance corresponding to the group contribution is equally shared between all correlated features. This implies that features belonging to larger groups receive smaller importance scores compared to a equally informative group but with less correlated features. Toloşi and Lengauer [2011] refer to this phenomenon as the correlation bias and noted that if the group is large enough, all features may appear as irrelevant (because of their low importance scores), even if they are highly informative about the output.

Let us however mention that due to the masking effect can counter-balance this bias as only some features of the group may collect the whole group importance, implying that some other features are masked and so of lower importances.

 $^{^{17}}$ They are assumed to be strictly equivalent and not masked, or equivalently, that K=1. Moreover, let us consider that they are also identical on other aspects, such as their cardinalities, to prevent other biases.

Bias due to number of categories and scale of measurement

Source of bias: features of various natures and different cardinalities.

It is known for a long time that the Gini impurity is biased in favour of features of higher cardinalities which thus offer more potential splits Breiman et al. [1984]; Kim and Loh [2001]. This phenomenon is usually referred to as the so-called "bias selection". Since (Gini) MDI importance measure is directly derived from impurity decreases within trees, it suffers from the same bias towards features of higher cardinalities and numerous studies reported this selection bias for the MDI importance measure (see, e.g., [Dobra and Gehrke, 2001; Strobl et al., 2007a; Boulesteix et al., 2012]). [Strobl et al., 2007b] noted that features with high cardinality (i.e., categorical features with a large number of categories or continuous ones) offer more potential cut-points (splits on that feature) and are thus more likely to provide a good split with respect to features of lower cardinalities. Consequently, the number of categories and the scale of measurement affects the feature and some features might be more frequently selected by a (Gini-based) impurity criterion yielding biased MDI importance scores and misleading feature ranking.

In contrast, it has been observed that this bias does not impact MDA importance measure [Strobl et al., 2007b; Boulesteix et al., 2012]. The explanation given is that a feature that is more frequently selected does not necessarily improves the oob accuracy and thus may receive low MDA importance scores despite being often used in the model. This however increase the variance of MDA importances [Boulesteix et al., 2012].

Let us note that comparison between continuous and discrete features (i.e., with different domain size) is not specific to trees and has been studied in other context (see, e.g., [Jiang and Wang, 2016]).

Louppe [2014] however suggests that the observed bias in Strobl et al. [2007b]'s study is mainly due to empirical misestimations (see Section 4.4.5.5). Indeed, this bias was also observed when no feature or split value selections are performed (e.g., for Extra-Trees with K = 1 or for totally randomised trees). This suggests that the bias is not only caused by a preference for features of higher cardinalities.

4.4.5.4 Bias due to the category frequencies

Source of bias: features of various category frequencies.

In genetic epidemiology, single nucleotide polymorphisms (SNPs), i.e., variation of a single nucleotide that occurs at a specific position in the genome, are of interest to study some diseases and personalised medicine [Carlson, 2008] and are known to interact with each others. In the context of genetic association studies, all SNPs have the same number of categories but vary in their category frequencies. Experiments reveal that both importance measures prefer informative SNPs with larger minor allele frequency¹⁸ (MAF) with respect to informative SNPs with lower MAF and (Gini) MDI importance measure is still biased in case of non-informative SNPs [Nicodemus, 2011; Boulesteix et al., 2011].

This phenomenon, known as the minor allele frequency bias, highlight the bias due to an unbalance in the category frequencies, or more generally in the value

¹⁸It refers to the frequency of the second most frequent allele value.

distribution. Let us note that the presence of missing values modifies the actual value distribution and therefore may also impact importance measures.

4.4.5.5 Bias due to empirical impurity estimations

Source of bias: number of learning samples N.

In the beginning of this chapter, the size of the learning set was never actually taken into account. In all theoretical analyses, a learning set of infinite size (in asymptotic conditions) assumes that the joint probability density $P_{V,Y}$ is known. Similarly, most empirical studies consider artificial datasets and thus the generating model was also known. In practice however, the learning set size is finite and this may cause empirical misestimations. For multiway splits, Louppe [2014] observes that misestimation bias in (Shannon) MDI importance relates to the misestimations of the mutual information terms $\Delta i(s,t) \approx I(X_i;Y|t)$. For independent random variables X_i and Y, the mean of the distribution of finite sample size estimates of their mutual information is proportional to the cardinalities $|X_i|$ and |Y| and inversely proportional to N_t, the number of samples in node t. This explains why MDI importance measures tend to positively bias importance of features of higher cardinalities. We refer to [Louppe, 2014] for a detailed analysis. The case of binary splits is discussed in Section 4.4.5.6.

Notice that the estimation of impurity measures and impurity decreases, in particular Shannon entropy and mutual information, has been widely studied in general frameworks that are not directly related to tree-based methods (see, e.g., [Moddemeijer, 1989; Beirlant et al., 1997; Paninski, 2003; Schürmann, 2004).

4.4.5.6 Bias due to binary splits and split value selection

Source of bias: tree-based algorithms.

Unlike multiway splits, binary splits do not fully exploit a variable. A binary split only discretises the information contained in a variable and therefore the same variable (if not binary) can be reused several times in the same branch. Therefore, binary splits $\Delta i(s,t)$ actually estimates the mutual information between the output and the split outcome (as mentioned in Section 3.2.2.1) while multiway splits would provide an estimate of the mutual information between the split variable and the outcome. As a consequence, the estimated mutual information $I(X_i; Y)$ is actually a collection of potentially biased estimates provided by all binary splits [Louppe, 2014]. From a different angle, explored branches are not equivalent in binary and multiway trees. A feature can be used several times in binary trees but only once in multiway tree because branches correspond to single value of the split variable. Feature importance scores are therefore not computed from the same sequence of impurity terms and can therefore be different. Louppe [2014] gives an illustrative example of two features whose importance scores are different depending on the kind of tree used to compute them. Moreover, the discretisation directly depends on the split value selection and thus the chosen strategy may have an impact on the feature importance scores. For example, a random split value selection such as in Extra-Trees may induce more splits on the same variable and thus more impurity terms, each providing part of the information contained in the feature, compared to an optimal split value such as in Random Forest that may yield all the information contained

in the feature in only one split. Feature importance scores obtained with one or the other technique can thus also differ [Louppe, 2014].

4.4.5.7 Bias due to bootstrapping

Source of bias: tree-based algorithms and number of learning samples N.

Strobl et al. [2007b] observe that the bootstrap sampling increases the bias due to the cardinality and therefore suggest not to use bootstrap. Moreover, it has been shown experimentally in [Louppe and Geurts, 2012] that bootstrapping is rarely crucial for random forest to obtain good accuracy.

The second observation is not directly due to the bootstrap mechanism but related to the number of OOB samples. The number of samples N has a direct impact on the resolution of MDA importance measure. On average, around 37% of original samples are not represented in the bootstrap sample. Therefore, when computing the MDA importance score on those samples, only granular values of accuracy change can be obtained when N is small because to resolution is limited to approximately 3/N, yielding over- and under-estimations of true feature importances [Archer and Kimes, 2008].

MEANINGFUL THRESHOLDS ON FEATURE IMPORTANCES

Feature importance measures can be used to rank features in order to facilitate the identification of a useful subset of important features. In this way those features having an importance below some threshold would be considered as unimportant and thus eliminated from further consideration. Unfortunately, there is no natural way to choose a "good" threshold on importances [Janitza et al., 2015]. Therefore, in practice, performing feature selection from such a ranking consists in selecting the k top features (i.e., with highest importance scores). This then reduces to the determination of a "good" value of k. This may be trivial if one observes a huge gap between relevant and irrelevant features, however in practice, such differences are not common and importance scores are usually smoothly decreasing when going down in the ranking. In such cases, distinguishing when features are no longer informative and when their importances are due to random fluctuations or some undesirable effects, is much more complicated. In this section, we give a non-exhaustive list of several approaches that allow to find either a threshold separating importance scores of relevant features from irrelevant, or propose to use or derive some statistical measure scores for which thresholds are usually more interpretable [Konukoglu and Ganz, 2014]. Let us note that methods that are not specific to tree-based methods are asterisked.

RANDOM PROBE* [STOPPIGLIA ET AL., 2003B] In the probe feature method, the key idea is to introduce a random feature in the feature ranking technique. This probe is expected to be ranked similarly as other irrelevant features and all features ranked below the probe should be naturally discarded. However, this probe rank can actually be seen as a random variable and its cumulative distribution function can be computed exactly or estimated (through the generation of several realisations of that random variable). One can then choose an acceptable value of risk and derive the corresponding rank position (and the corresponding threshold importance value) in order to discriminate relevant from irrelevant features.

ARTIFICIAL CONTRAST VARIABLES [TUV ET AL., 2006] Similarly to random probes, Tuv et al. [2006] propose to introduce M contrast features that are known to be truly independent of the output and to generate them by randomly permuting values of M input features. By the means of a t-test and a significance level, this allows to identify relevant features as those with importance scores significantly better than those of contrast features. Additionally, they propose to estimate split weights from oob samples and to introduce the mechanism of contrast features in an procedure building iteratively ensemble of trees on kept features and a residual of the target.

FEATURE IMPORTANCE AS A REAL-VALUED PARAMETER* [VAN DER LAAN, 2006] The principle of their approach is to define the wished feature importance measure (in particular, in prediction tasks) as a real-valued parameter and propose estimators for those feature importance parameters, accompanied with a pvalue and confidence interval.

MDA Z-SCORE [BREIMAN AND CUTLER, 2008] As defined in Section 4.2.2, the MDA importance score for a feature X_m derived from an ensemble T of N_T trees consists of the average impact of removing a feature on the accuracy of every tree in the forest. In contrast to this "raw" MDA importance score of a feature, [Breiman and Cutler, 2008] propose a "scaled" version for which the raw importance score is divided by its standard error. This importance measure is usually referred to as the z-score of a feature. If all individual importance scores have the same standard deviation σ , the standard error of the mean of those individual scores is $\sigma/\sqrt{N_T}$ [Strobl and Zeileis, 2008]. The z-score of X_m is therefore given by

$$\operatorname{Imp}_{z}^{\operatorname{mda}}(X_{\operatorname{m}},\operatorname{rf},\mathbf{LS}) = \frac{\operatorname{Imp}_{\operatorname{rf}}^{\operatorname{mda}}(X_{\operatorname{m}},\operatorname{rf},\mathbf{LS})}{\frac{\sigma}{\sqrt{N_{T}}}},$$
(4.20)

where rf is a random forest algorithm. Assuming that individual importance scores are independent because they are computed from independent bootstrap samples [Strobl and Zeileis, 2008], then Equation 4.20 tends towards a normal distribution by the central limit theorem. Therefore, a statistical test can be conducted to check whether the null hypothesis of zero importance for variable X_m (i.e., corresponding to an irrelevant variable X_m) is true or not for a given significance level. However, [Strobl and Zeileis, 2008] find out that the power of this test based on z-scores decreases with an increasing sample size and increases boundlessly with the number of trees and claim that these are undesirable properties for an importance measure.

FEATURE SET PERMUTATION SCHEME [TANG ET AL., 2009] Instead of permuting a single feature, Tang et al. [2009] propose to permute a set of features. In their application, each gene corresponds to a set of SNPs. Permuting all SNPs corresponding to the same gene allows to make a gene-permutation that directly evaluates the importance of the gene.

LABEL PERMUTATION SCHEME [ALTMANN ET AL., 2010] In their work, they use a permutation test to obtain a threshold for the selection of relevant features. Firstly, un-permuted feature importance scores are computed. Secondly, m permutations are generated by randomly permuting the labels and then, for each permutation, "permuted" feature importance scores are computed. From that, p-values can be determined by the fraction of permuted importances that are larger than the un-permuted importances and then a threshold can be chosen from a given significance level. Rodenburg et al. [2008] also suggest a second approach that consists in keeping all features whose importance scores are larger than the mean value of maximal permuted importances. This approach however appears to be very restrictive. Alternatively, Altmann et al. [2010] propose to fit a parametrised probability distribution on permuted importance scores.

- CONDITIONAL PERMUTATION SCHEME [STROBL ET Al., 2008] is an alternative permutation scheme aiming at measuring the impact of a feature on the output conditionally to other features in comparison with the classical permutation scheme, and so to correct for the bias towards correlated features. See side note on page 111 for details on this permutation scheme.
- SEPARATE FEATURE PERMUTATION SCHEME [HAPFELMEIER AND ULM, 2013] Instead of permuting labels or group of features, Hapfelmeier and Ulm [2013] propose to permute feature individually while keeping the output and all other features unchanged. The proposed new permutation scheme aims at measuring only the impact of a feature on the output.
- APPROXIMATE FALSE POSITIVE RATE CONTROL [KONUKOGLU AND GANZ, 2014] Permutation techniques can be intractable for high-dimensional datasets and therefore Konukoglu and Ganz [2014] propose an approach to determine thresholds and control the false positive rate in random forest method at no additional computational cost. Based on the feature selection frequency importance measure (see Section 4.1), they rely their approach on the estimation of the probability that a feature is selected k times in a tree ensemble if it is assumed to be irrelevant to the output. They propose an approximate model for selection frequency in random forest from which one can determine a desired level of false positive rate and obtain an optimal threshold on the selection frequency importance scores.
- CONDITIONAL ERROR RATE* [HUYNH-THU ET AL., 2008; HUYNH-THU, 2012] In order to overcome limitations of classical permutation-based techniques of false positive rate estimation, Huynh-Thu et al. [2008] propose the conditional error rate (CER) as an alternative measure to be associated with each importance threshold τ_i . It estimates the probability to include an irrelevant feature when selecting all features (assumed to be relevant) with an importance score greater or equal to τ_i .
- Note that [Huynh-Thu et al., 2012; Wehenkel, 2018] review statistical interpretation of (tree-based) feature importance scores, including random probe techniques and conditional error rate.
- RANK-BASED CONDITIONAL ERROR RATE [WEHENKEL ET AL., 2017] While the CER is based on the importance scores, Wehenkel et al. [2017] propose an adaptation of CER based on rank for group of features. Let us assume an original order of feature groups ranked by order of decreasing importance scores.



ABOUT PERMUTATION SCHEMES

Let us consider a set V of input features and an output Y. Following [Strobl et al., 2008; Hapfelmeier and Ulm, 2013], we detail hereafter permutation schemes that have been proposed to evaluate the MDA importance of a feature $X_m \in V$. The classical permutation scheme, as described in Section 4.2.2, consists in permuting X_m against both the output Y and the remaining features V^{-m} . It therefore simulates the independence between X_m and both Y and V^{-m} . Mathematically, the evaluated independence (null hypothesis) is

$$X_m \perp \!\!\! \perp (Y \cup V^{-m}) \Rightarrow X_m \perp \!\!\! \perp Y \text{ and } X_m \perp \!\!\! \perp V^{-m} \text{ (decomposition property)}.$$

Note that the converse (⇐) is also verified if the composition property is satisfied (e.g., for a strictly positive distribution). Consequently, a deviation yielding a positive importance can result of a violation of the independence either between X_m and Y, or X_m and V^{-m} .

The label permutation scheme (proposed by [Altmann et al., 2010], see page 109) consists in permuting the output value. On one hand, this breaks all relationships between X_m and Y, but on the other hand it also breaks any relationships between any input feature in V^{-m} and Y. Therefore, the evaluated independence is

$$(X_m \cup V^{-m}) \perp \!\!\! \perp Y$$

and would wrongly attribute to X_m the importance of all input features. Permuting the output values is therefore equivalent to permuting all input feature values jointly (i.e., v^{-m} of V^{-m}). Instead of permuting all input features, Tang et al. [2009] suggest to only permute a group of features P including $X_{\rm m}$. In this work focusing on identifying relevant SNPs (input features) in GWAS^a, they propose to simultaneously permute all SNPs which belong to the same gene. Within this permutation scheme, the evaluated independence is

$$P \perp \!\!\!\perp (Y \cup V \setminus P)$$

which does not allow to evaluate the importance of the single feature $X_{\rm m}$. On the contrary this gives the importance of the group to every feature within this group. Hapfelmeier and Ulm [2013] argue that each feature needs to be permuted separately in order to correctly estimate the importance of a single variable which is not possible by means of a label permutation.

With the conditional permutation scheme, Strobl et al. [2008] suggest to permute $X_{\mathfrak{m}}$ only within groups of observations with $V^{-\mathfrak{m}}=\nu^{-\mathfrak{m}}$ in order to preserve the relationships between X_m and all features in V^{-m} while destroying the link with Y. It corresponds to the following evaluated independence

$$X_m \perp \!\!\! \perp Y | V^{-m}$$

which highlights the conditioning on V^{-m}. Interestingly, it corresponds to the definition of strongly relevant features. The conditional permutation scheme may therefore miss some weakly relevant features that are independent of Y knowing all other features (e.g., redundant features).

^aGenome wide association studies.

The key principle is that a relevant group should not be as well or better ranked than originally once all statistical links within this group and in all groups ranked below (in the original order) are broken. Wehenkel [2018] note that this variant is less restrictive than the original method.

SUBSAMPLING AND DELETE-D JACKNIFE [ISHWARAN AND LU, 2018] Recently, Ishwaran and Lu [2018] study several sampling approaches for estimating MDA importance measure variance, such as double-bootstrap, subsampling and deleted jacknife algorithm. They additionally propose a subsampling approach that can be used to estimate the standard error of MDA importance measure and for defining confidence intervals.

4.6 EXTENSIONS AND DERIVATIONS

In this section, we briefly review some methodologies that exploit feature importance measures or derive their use to perform new tasks.

- RECURSIVE FEATURE ELIMINATION [DÍAZ-URIARTE AND DE ANDRES, 2006] Their approach is an instance of the Sequential Back Elimination (SBE, see 2.4.6) that recursively removes features with the smallest importance scores computed with a tree-based ensemble method.
- ENRICHED RANDOM FORESTS [AMARATUNGA ET AL., 2008] In presence of many irrelevant features, many splits can be made on irrelevant features because all split variables candidates were irrelevant. In order to circumvent that, Amaratunga et al. [2008] propose a weighted random sampling in each node instead of a uniform one. They suggest to determine weights as the p-value of a t-test.
- GUIDED REGULARIZED RANDOM FOREST [DENG AND RUNGER, 2013] Similarly to [Amaratunga et al., 2008], their approach first builds a classical random forest and then use feature importance scores to guide the feature selection process in a second model (i.e., regularized random forest [Deng and Runger, 2012]).
- VARIABLE IMPORTANCE-WEIGHTED FEATURE SELECTION [LIU AND ZHAO, 2017] Similarly with the previous approach, instead of selecting split variable candidates at random, Liu and Zhao [2017] propose to sample features according to their importance scores in order to focus on informative features.
- RANDOM SUBSPACE FOR FEATURE SELECTION [HO, 1998; LAI ET AL., 2006] Inspired from the Random Subspace method proposed by [Ho, 1998], this approach consists in growing each tree of the ensemble on a random subspace of K ($\leq p$) features randomly chosen. Similarly to a classical forest, feature importance scores are then computed for each tree and then aggregated with the difference that at least p - K features have necessarily a zero importance for each tree. One can however expect that all available features can be considered in the tree, even if it is made of only few nodes. Let us note that this approach is also compatible with the Random patches method [Louppe and Geurts, 2012].

SEQUENTIAL RANDOM SUBSPACE [SUTERA ET AL., 2018] In this sequential variant of the random subspace method, the key ideas are that (i) some relevant features may be difficult to identify because they need to be considered conditionally to some other features (ii) which are necessarily relevant. Therefore, the principle is to reuse more frequently features that have already been identified as relevant in order to make the detection of other relevant features easier. In contrast with approaches such as variable importance-weighted feature selection [Liu and Zhao, 2017], one can force the method to keep a part of exploration to discover masked features for instance.

FEATURE SELECTION WITH A KNOCK-OUT STRATEGY [GANZ ET AL., 2015] This approach is interesting in several respects. Uncommonly, they consider the frequency selection as importance measure on which they apply a false positive rate control (see Section 4.5 and [Konukoglu and Ganz, 2014]). Moreover, at each iteration, the identified set of relevant features are removed ("knocked out") in order to force the algorithm to identify remaining relevant features since already identified are no longer available. This method has the merit of looking for all relevant features without taking care of accuracy performances. However, [Sutera et al., 2018] show that relevant features may be required to reveal some others that are more difficult to detect (e.g., a clique) but this may be circumvented by the use of frequency selection instead of other importance measures.

REPRESENTATIVE FEATURE(S) [TOLOŞI AND LENGAUER, 2011] Proposed as a way to reduce the correlation bias, the idea developed in [Tolosi and Lengauer, 2011] is to group several "similar" features into representative feature(s) that can then be used as input features for the model. At the end, the importance scores of the original features can be retrieved as the importance of the representative feature (or the average in case of several representatives).

Wehenkel [2018] reviewes some approaches to determine the representative features and discusses those based on a priori knowledge (e.g., atlas for brain regions [Wehenkel, 2018], self-organizing maps for genes [Rodenburg et al., 2008]) and on neighbouring positions. Let us also mention that similar features can also be identified with techniques such as hierarchical clustering [Rodenburg et al., 2008].

GROUP IMPORTANCE SMOOTHING [WEHENKEL, 2018] Because of masking effect or some other biases, similar features may receive different importance scores. [Wehenkel, 2018] proposes two ways to post-process importance scores in order to rebalance more fairly importance scores among similar features. The first approach consists in sharing the importance score of a feature with its neighbours. The second approach consists in assigning all features of a group (e.g., based on a priori knowledge) the same "group" importance scores that have been derived from the distribution of all importance scores within the group. A group importance is then derived using either the average, the sum, or the maximum of the individual importance scores within the group.

4.7 OTHER IMPORTANCE MEASURES

Previous sections show in several respects that MDI and MDA feature importance measures are not perfect and can not address all needs. Thus, several other importance measures have been proposed in the literature. Some of them are described in this section. Note that we exclude from the following list local feature importance measures, such as Shapley values [Lundberg and Lee, 2017; Lundberg et al., 2018], that evaluate feature importances for a given input vector x, although global feature importance measures can be obtained from such local measures by aggregating them over a sample of input vectors.

- CROSS-VALIDATED MDA FEATURE IMPORTANCE [JANITZA ET AL., 2015] Firstly, they propose an alternative approach to compute the MDA importance measures of cross-validated subsets instead of oob samples. The principle is similar: the accuracy is estimated on samples that have not been used to learn the model, i.e., the remaining fold. Secondly, they propose a new variable importance test that is computationally more efficient than traditional permutation schemes discussed in Section 4.5.
- CONTEXTUAL IMPORTANCE MEASURES [SUTERA ET AL., 2016] MDI importance measures are extended to identify and characterise features whose relevance is context-dependent (i.e., varying depending on the context) or context-independent.
- AUC-BASED PERMUTATION IMPORTANCE MEASURE [JANITZA ET AL., 2013] To overcome the sub-optimality of random forest methods in presence of strongly unbalanced data, Janitza et al. [2013] propose to use an AUC-based criterion instead of an error-rate-based one for the MDA importance measure.
- CHANGE IN CLASS VOTE DISTRIBUTION [PAUL ET AL., 2013] In this work, a new feature importance index is proposed that uses a statistical test to determine whether permuting a variable significantly influences the class vote distribution of the forest. This new importance measure correlates well with MDA importance and has the advantage of providing directly a p-value.

Without more explanations, let us however mention [Sandri and Zuccolotto, 2008; Nembrini et al., 2018] proposing bias-corrected impurity importance measures by the means of the addition of uninformative features (e.g., permutation of original features) among input ones.

4.8 SOME APPLICATIONS EXPLOITING FEATURE IMPORTANCES

To conclude this chapter, we briefly mention in this section two applications of feature importance measures in the biomedical domain.

GENE NETWORK INFERENCE In genomics, regulatory gene network inference consists in the identification of all gene-to-gene interactions from their expression level and reconstruct a network with these interactions. Concretely, one needs to infer a (un)directed graph where each nodes is a biological entity (e.g., a gene) and edges connecting two nodes represent an interaction between them. In all generality, the GENIE3 method aims at inferring a network of p nodes by decomposing it into p independent supervised learning problems. Each feature is in turn considered as the target to predict from all p-1 other features. When these sub-problems are solved by the means of tree-based methods, feature importance measures can be derived and seen as indications of the degree of association between input features and the target. Concretely, in a model predicting X_i from V^{-j} , the importance score of a feature $X_i \in V^{-j}$ is used a the degree of association between node i and node j. Once all sub-problems are solved, the ranking of all gene-gene pairs can be used to reconstruct the global network (e.g., by selecting the stronger interactions). Chapter 8 focuses on that application.

NEUROIMAGING Random forest methods are able to handle high-dimensional dataset $(p \gg N)$, such as neuroimaging datasets, and therefore constitute interesting alternatives to SVM and deep learning methods in the context of neuroimaging datasets. For example, in the context of fMRI datasets, [Langs et al., 2011] use Gini MDI importance to identify interacting brain regions that are activated under experimental stimuli, and [Richiardi et al., 2010] exploit tree-based feature importance measures to determine relevant brain region connections. In the particular case of Alzheimer's disease, Wehenkel et al. [2018]; Wehenkel [2018] exploit feature importance measures to identify important (group of) voxels from Positron Emission Tomography (PET) images in order to identify brain regions involved in the prognosis of the disease.

✓ Chapter take-away

In the litterature, several measures have been proposed to quantify the importance of features from tree-based ensemble models. Because of their ability to handle feature interactions and non-linearities, these measures are interesting alternatives to classical statistical tests. Driven by many successful applications (notably in the biomedical domain), several studies have been carried out to analyse these measures from different perspectives that have revealed several biases from which these measures suffer. Recent theoretical works also give new insight on previous empirical results. One major drawback of standard importance measures is that they lack a statistical interpretation that would allow to naturally determine a threshold value to distinguish truly important from non-important features. Several techniques, mostly based on random permutations, have however been proposed in the literature to address this issue. In addition, new tree-based importance measures have been designed to go further in the exploitation and interpretation of tree-based ensemble models.

Overview

In this chapter, we characterise the Mean Decrease of Impurity (MDI) feature importance measure as computed by an ensemble of randomised trees. First, in asymptotic conditions, we derive a multi-level decomposition of the information jointly provided by all input features about the output, with a particular attention on the link between importance and relevance of features. We also extend the characterisation to take into account the presence of redundant features. We then analyse importance measure properties in the case of non-totally randomised, non-fully developed and binary trees, respectively. Finally, we discuss how these properties may change in the finite case, in particular in the number of trees and samples.

References: This chapter presents results that were published in the following publications:

- G. Louppe, L. Wehenkel, A. Sutera, and P. Geurts. Understanding variable importances in forests of randomized trees. In *Advances in neural informa*tion processing systems, pages 431–439, 2013;
- A. Sutera, C. Châtel, G. Louppe, L. Wehenkel, and P. Geurts. Random subspace with trees for feature selection under memory constraints. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 929–937, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. URL http://proceedings.mlr.press/v84/sutera18a.html.

Note that proofs from the first publication, which were also part of [Louppe, 2014], are not reproduced below.

Nowadays, most of state-of-the-art supervised learning algorithms typically provide a black-box model able to accurately predict the output. In many applications, a particular attention is paid to an understanding of the modelled system, which is typically not possible with a black-box model. Random forest methods, by the means of importance measures, allow to identify important features which are the key elements of the model. This interpretation provides insights to understand the underlying mechanism. Concretely, given an ensemble of trees, one may derive a numerical score for each feature that assesses its importance in the tree-based model. Breiman [2001]; Breiman and Cutler [2003] proposed two importance measures¹. Firstly, the Mean Decrease of Accuracy aims at evaluating the contribution of a feature for predicting the output as the change in accuracy of the model when this feature is permuted. Secondly, the Mean Decrease of Impurity (MDI) relies on the impurity criterion used to grow trees. In this chapter, we only focus on that partic-

¹Note that both importance measures are described in Chapter 4.

ular importance measure. It adds up the weighted impurity decreases $\Delta i(s,t)$ over all nodes t in a tree T where the variable X_m is used to split and then averages this quantity over all trees in the ensemble, i.e.²:

$$\begin{split} & \text{Imp}(X_m) = \frac{1}{N_T} \sum_{t \in T: \nu(s_t^*) = X_m} p(t) \Delta i(s_t^*, t) \\ & \text{with } \Delta i(s, t) = i(t) - \frac{p(t_L)}{p(t)} i(t_L) - \frac{p(t_R)}{p(t)} i(t_R) \end{split} \tag{5.1}$$

where i is the impurity measure (introduced in Section 3.2.2.1), p(t) is the proportion of samples reaching node t, $v(s_t)$ is the variable used in the split s_t , at node t, and t_L and t_R are the left and right successors of t after the split.

In Chapter 4, we outlined a theoretical analysis of both importance measures and then focused on their practical uses, in particular biases that may provide misleading interpretations of feature ranking and importance scores. We also consider several extensions, derivations and applications in which feature importance are typically used.

Despite these numerous works, only few studied theoretically feature importance measures from a theoretical point of view [Ishwaran, 2007; Louppe et al., 2013; Louppe, 2014; Zhu et al., 2015; Gregorutti et al., 2017; Sutera et al., 2018]. In order to go one step further in the understanding of this measure, this chapter aims at providing an in-depth theoretical analysis of the MDI importance derived from ensembles of randomised trees in an infinite sample setting. We also discuss how it may change in the case of finite sample and tree ensemble size conditions.

As a preambule, Section 5.1 first defines the degree of a relevant variable and provide two propositions that characterize minimal conditionings that make relevant variables dependant of the output. Section 5.2 then provides a theoretical characterisation of MDI importance measures in asymptotic conditions in the case of totally randomized and fully developed and presents an interpretable decomposition of the information jointly provided by all input features about the output at several levels of feature interactions. Section 5.3 shows that MDI importance measures can be used to identify relevant features. Section 5.4 extends the characterisation of MDI importance measures to highlight the impact of the presence of redundant features. Sections 5.5 and 5.6 consider respectively non-totally randomised and non-fully developed trees and analyse to what extent MDI properties are still verified. Section 5.7 examines trees made with binary splits and aims at extending the characterisation of multiway trees to binary ones that are more common in practice. Finally, Section 5.8 discusses the finite case and in particular considers a finite number of trees (Section 5.8.1) and a finite number of samples (Section 5.8.2).

Notational conventions

For sake of clarity, the setting under study is reminded at the beginning of the sections and summarised by some of the following parameters (described in Chapter 3): the split selection randomisation parameter K of the random forest algorithm, the maximal depth D of the tree structure, the split cardinality³ $|s_t|$ used in decision

²From now on, this thesis only focuses on the MDI importance measure and the notation is thus simplified accordingly, i.e., Imp(X) is equivalent to $Imp^{mdi}(X)$.

 $^{^3}$ |s_t| = | $v(s_t)$ | denotes a tree built with multiway exhaustive splits as the split cardinality equals the number of values of split variable $v(s_t)$.

trees, the number of trees N_T in the ensemble, and the number of samples N of the learning set. For the sake of completeness, subscripts and superscripts will be used to specify the parameter values of the tree-based method used to derive importance scores: ${\rm Imp}_{N,N_T}^{K,D}$ corresponds to the importance measure computed with an ensemble of N_{T} trees built with a split randomisation parameter K and a maximal depth D on a dataset of N samples.

5.1 DEGREE OF RELEVANT VARIABLES

In addition to the definitions of relevance provided in Section 2.4.1 (Definitions 2.5, 2.4 and 2.6 in terms of conditional independences and Definitions 2.7 and 2.8 in terms of mutual informations), for some results derived below, we need to qualify relevant variables according to their degree:

Definition 5.1. [Sutera et al., 2018, Definition 3] The degree of a relevant variable X, denoted $\deg(X)$, is defined as the minimal size of a subset $B \subseteq V$ such that $Y \perp \!\!\! \perp X \mid B$.

Relevant variables X of degree 0, i.e. such that $Y \perp \!\!\! \perp X$ unconditionally, will be called marginally relevant.

We will say that a subset B such that $Y \perp \!\!\! \perp X|B$ is **minimal** if there is no proper subset B' \subseteq B such that Y $\not\perp$ X|B'. The following two propositions give a characterisation of these minimal subsets.

Proposition 5.1. [Sutera et al., 2018, Proposition 1] A minimal subset B such that $Y \not\perp \!\!\! \perp X \mid B$ for a relevant variable X contains only relevant variables.

Proof. Let us assume that B contains an irrelevant variable X_i . Let us denote by B^{-i} the subset $B \setminus \{X_i\}$. Since X_i is irrelevant, we have $Y \perp \!\!\! \perp X_i \mid B^{-i} \cup \{X\}$. Given that B is minimal we furthermore have $Y \perp X \mid B^{-i}$ where $B^{-i} = B \setminus \{X_i\}$. By using the contraction property of any probability distribution (see side note on page 38), one can then conclude from these two independences that $Y \perp \{X, X_i\} B^{-i}$ and, by using the weak union property, that $Y \perp \!\!\! \perp X|B$, which proves the theorem by contradiction.

Proposition 5.2. [Sutera et al., 2018, Proposition 2] Let B denote a minimal subset such that $Y \perp \!\!\! \perp X \mid B$ for a relevant variable X. For all $X' \in B$, $deg(X') \leq |B|$.

Proof. If we reduce the set of features V to a new set $V' = B \cup \{X\}$, X will remain relevant, as well as all features in B, given Proposition 5.1. So, for any feature X' in B, there exists a subset $B' = B \cup \{X\} \setminus X'$ such that $Y \not\perp X' \mid B'$ and the degree of X'is therefore $\leq |B|$.

These two propositions show that a minimal conditioning B that makes a variable dependent on the output is composed of only relevant variables whose degrees are all smaller or equal to the size of B. Let us note that we will provide in Section 7.3.2 a more stringent characterisation of variables in minimum conditionings in the case of specific classes of distributions.

TOTALLY RANDOMISED AND TOTALLY DEVELOPED TREES

Setting of this section: $K = 1, D = p, |s_t| = |v(s_t)|, N_T \to \infty, N \to \infty$

Let us assume a set $V = X_1, ..., X_p$ of categorical input variables and a categorical output Y. Let us consider a joint probability density $P_{V,Y}$ of X_1, \ldots, X_p, Y and a learning set LS of N observations of X_1, \ldots, X_p , Y independently drawn from that distribution. From LS, an infinitely large ensemble of totally randomised, multiway and fully developed trees is inferred. As a reminder of Chapter 3, such trees are built such that, for each node t, a split variable X_i is selected totally at random among those not yet picked and used to split the node t into $|X_i|$ branches (i.e., one for each value of X_i), until there is no more remaining unused features. Let us note that all branches have the same depth p, because each feature is used once along each branch. For sake of simplicity, we only consider Shannon impurity to evaluate the importances, but results can be extended to some extent to other impurity measures as shown in [Louppe et al., 2013, Appendix I]. Note that in the totally randomized setting, the tree structure does not depend on the impurity measure, but the MDI importance measure derived from this structure obviously does.

In that context, let us consider the MDI importance as defined by Equation 5.1 computed by this ensemble of trees.

Theorem 5.3. [Louppe et al., 2013, Theorem 1] The MDI importance of $X_m \in V$ for Y as computed with an infinite ensemble of fully developed totally randomized trees and an infinitely large learning set is:

$$Imp_{\infty,\infty}^{1,p}(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} I(X_m; Y|B),$$
 (5.2)

where V^{-m} denotes the subset $V\setminus\{X_m\},\,\mathcal{P}_k(V^{-m})$ is the set of all subsets of cardinality k of V^{-m} , and $I(X_m; Y|B)$ is the conditional mutual information of X_m and Ygiven the variables in B. A setting when both learning set and tree ensemble sizes are assumed to be infinitely large is further referred to as asymptotic conditions.

Theorem 5.4. [Louppe et al., 2013, Theorem 2] For any ensemble of fully developed trees in asymptotic learning sample size conditions (e.g., in the same conditions as those of Theorem 5.3), we have that

$$\sum_{m=1}^{p} Imp_{\infty,\infty}^{1,p}(X_m) = I(X_1, \dots, X_p; Y).$$
 (5.3)

Proof. See [Louppe et al., 2013, Appendix C] for a proof.

In Theorem 5.4, the term $I(X_1,...,X_p;Y)$ denotes the information contained in the set of input variables about the output variable and can be computed for a given joint probability density $P_{V,Y}$. Let us notice that this property actually holds for every single tree, and consequently also for any ensemble of N_T trees, and in particular when N_T goes to infinity. Given that $I(X_1, ..., X_p; Y)$ is fixed for a given problem, Theorem 5.3 shows that an increase of the importance of one feature will always come with a decrease of the importance of another feature.

Combining Theorems 5.3 and 5.4 in the context of ensemble of trees, the information contained in the set of inputs variables can be decomposed into the following three-level nested sums:

$$I(X_{1},...,X_{p}:Y) = \sum_{m=1}^{p} \sum_{k=0}^{p-1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-m})} I(X_{m};Y|B)$$
 (5.4)

The first sum is over the variables, the second sum over the degrees k of the interaction terms, and the third sum over all conditioning subsets B of size k. Equivalently, the first two sums can be swapped to yield the following decomposition of $I(X_1,\ldots,X_p;Y)$:

$$I(X_1,...,X_p:Y) = \sum_{k=0}^{p-1} \sum_{m=1}^{p} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} I(X_m;Y|B)..$$
 (5.6)

While Equations 5.3 and 5.4 divide the total output information between the features, computing each term of the outer sum in Equation 5.4 will give a decomposition of $I(X_1, ..., X_p; Y)$ per interaction degree, which highlights how important feature interactions are for predicting the output.

Table 5.1 illustrates these two ways of decomposing $I(X_1,...,X_p;Y)$ in the context of the digit recognition problem of [Breiman et al., 1984] (see Appendix C for a description of this problem). We can observe that almost all inner sum terms $\sum_{B} I(X_m; Y|B)$ are strictly positive implying that large conditioning sets B (corresponding to deep nodes in the tree) still contribute to the total variable importance. In this example, importances monotonically decrease with the degree of interaction k, but this is not always the case (e.g., with XOR-like structures)

	k = 0	k = 1	k = 2	k = 3	k = 4	k = 5	k = 6	$\sum_{\mathbf{k}}$
X_1	0.103	0.085	0.068	0.053	0.042	0.033	0.029	0.413
X_2	0.139	0.126	0.105	0.082	0.060	0.042	0.029	0.582
X_3	0.103	0.091	0.081	0.073	0.066	0.061	0.057	0.531
X_4	0.126	0.114	0.097	0.077	0.058	0.042	0.029	0.542
X_5	0.139	0.123	0.106	0.090	0.076	0.065	0.057	0.657
X_6	0.067	0.056	0.043	0.031	0.020	0.010	0.000	0.226
X_7	0.126	0.098	0.070	0.045	0.025	0.010	0.000	0.372
\sum_{m}	0.802	0.692	0.568	0.450	0.347	0.262	0.200	3.322

Table 5.1: Feature importances as computed with an ensemble of totally randomised trees. Last row (\sum_{m}) corresponds to importances per interaction degree (i.e., summed over over all features, see Equation 5.6) while last column (\sum_k) corresponds to importances per feature (i.e., summed over all interaction degrees, see Equation 5.4). Let us note that the sum of all importances is equal to $I(X_1,...,X_7;Y) =$ $\log_2(10) = 3.322.$

The last sum in Equations 5.4 or 5.6 includes all interaction terms of a given degree and it is weighted in a way that depends only on the combinatorics of possible

interaction terms. Interestingly, the weight $\frac{1}{C_p^k}\frac{1}{p-k}$ in front of each such sum perfectly

$$\frac{|\mathcal{P}_{k}(V^{-m})|}{C_{p}^{k}(p-k)} = \frac{C_{p-1}^{k}}{C_{p}^{k}(p-k)} = \frac{1}{p},$$

which is independent of k. This result is illustrated numerically for several values of p in Figure 5.1. Given that each mutual information term $I(X_m; Y|B)$ is upper bounded by H(Y), each term of the sum over k in Equation 5.4 is upper bounded by $\frac{1}{p}$ H(Y), which does not depend on k. It shows that importance measures are inherently unbiased with respect to interaction degrees.

5.3 IMPORTANCES OF RELEVANT AND IRRELEVANT VARIABLES

Setting of this section: $K = 1, D = p, |s_t| = |v(s_t)|, N_T \to \infty, N \to \infty$

The following theorems characterise the importances of relevant and irrelevant variables. These results can be derived from the equivalence between condition independance and zero conditional mutual information (see Section 4.3.2).

Theorem 5.5. [Louppe et al., 2013, Theorem 3] $X_i \in V$ is irrelevant to Y with respect to V if and only if its infinite sample size importance as computed with an infinite ensemble of fully developed totally randomized trees built on V for Y is 0.

Proof. See [Louppe et al., 2013, Appendix D] for a proof.

Corollary 5.6. Imp $_{\infty,\infty}^{1,p}(X_m) > 0$ iff $X_m \in V$ is relevant with respect to Y.

Proof. It directly stems from Theorem 5.5.

Lemma 5.7. [Louppe et al., 2013, Lemma 4] Let X_i ∉ V be an irrelevant variable for Y with respect to V. The infinite sample size importance of $X_m \in V$ as computed with an infinite ensemble of fully developed totally randomized trees built on V for Y is the same as the importance derived when using $V \cup \{X_i\}$ to build the ensemble of trees for Y.

Proof. See [Louppe et al., 2013, Appendix E] for a proof.

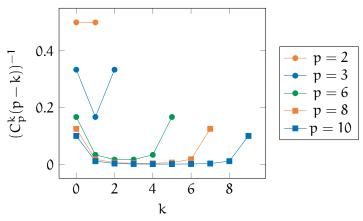
Theorem 5.8. [Louppe et al., 2013, Theorem 5] Let $V_R \subseteq V$ be the subset of all variables in ${\sf V}$ that are relevant with respect to ${\sf Y}$. The infinite sample size importance of any variable $X_m \in V_R$ as computed with an infinite ensemble of fully developed totally randomized trees built on V_R for Y is the same as its importance computed in the same conditions by using all variables in V. That is:

$$Imp(X_{m}) = \sum_{k=0}^{p-1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-m})} I(X_{m}; Y|B)$$

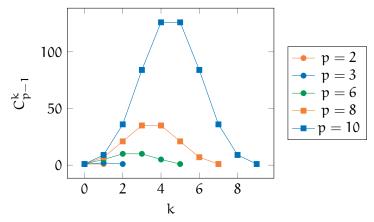
$$= \sum_{l=0}^{r-1} \frac{1}{C_{r}^{l}} \frac{1}{r-l} \sum_{B \in \mathcal{P}_{l}(V_{p}^{-m})} I(X_{m}; Y|B)$$
(5.7)

where r is the number of relevant variables in V_R .

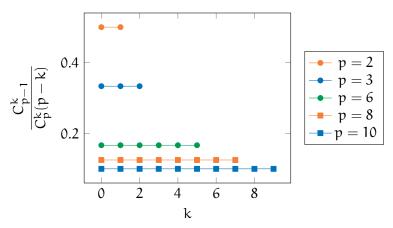
Proof. See [Louppe et al., 2013, Appendix F] for a proof.



(a) Evolution of $\frac{1}{C_p^k(p-k)}$ with respect to k. Note the symmetry.



(b) Evolution of C_{p-1}^k with respect to k. Note the symmetry.



(c) Evolution of $\frac{C_{p-1}^k}{C_p^k(p-k)}$ with respect to k. Note that for a given p, all values are equal.

Figure 5.1: Interpreting the weights in the three-level decomposition of total importance in Equation 5.4. Figure 5.1a shows how the weights of the second level of decomposition evolve with respect to k for several number of features p. Figure 5.1b shows the number of combinations B in the third level of decomposition. Figure 5.1c combines both decompositions and shows that sub-importance terms corresponding to every interaction degree equally contribute to the total importance.

Theorem 5.5 shows that only irrelevant features have a zero importance. They can thus be distinguished from relevant ones based solely on their importance scores. In addition, Lemma 5.7 points out that they do no affect the importance scores of relevant variables and the addition or the removal of irrelevant features have no effect which implies that only relevant features are required to compute importances (Theorem 5.8). Intuitively, splitting on an irrelevant feature X_i instead of a relevant feature X_m at node t only postpones the attribution of the local importance of X_m into the child nodes t_L and t_R , but do not actually change its total importance. Indeed, on one hand, if X_m was used at node t, then the local importance of X_m would be proportional to p(t) (i.e., $p(t)\Delta i(s,t)$). On the other hand, splitting on X_i at node t does not actually change the distribution of samples in t_L and t_R . Therefore, splitting then on X_m at t_L and t_R would provide the sum of local importances $p(t_L)\Delta i(s,t_L) + p(t_R)\Delta i(s,t_R)$. Given that $\Delta i(s,t) = \Delta i(s,t_L) = \Delta (i,t_R)$ because node sample distributions are unchanged by the split on X_i , we have that $(p(t_L) + p(t_R))\Delta i(s,t) = p(t)\Delta i(s,t)$ which shows that splitting on X_i first does not change anything. Similarly, one can recursively apply this reasoning if X_m was used deeper in the tree (i.e., at descendant nodes of t_L or t_R). Let us however note that this result may actually be due to the fact that total importance of a feature X_m is the sum of all local importances in nodes where X_m is used weighted by the number of samples reaching this node p(t). Louppe [2014] suggests that importances computed with another approach consisting in summing local importances over all nodes (e.g., using surrogate splits) would necessarily depend on the total number of nodes in a tree, which depends on the number of features p and not only on the number of relevant features r.

In conclusion, in our opinion, theorems 5.5 and 5.8 exhibit two desirable and sound properties for a feature importance measure.

IMPACT OF REDUNDANT VARIABLES

Setting of this section:
$$K = 1, D = p, |s_t| = |v(s_t)|, N_T \to \infty, N \to \infty$$

Let us consider redundant variables as defined in Section 2.4.3 and in particular totally redundant variables from Definition 2.13. In this section, we analyse how feature importance scores are affected by the presence of (totally) redundant variables.

Proposition 5.9. [Louppe, 2014, Proposition 7.2] Let $X_j \in V$ be a relevant variable with respect to Y and V and let $X'_i \notin V$ be a totally redundant variable with respect to X_i . The infinite sample size importance of X_i as computed with an infinite ensemble of fully developed totally randomized trees built on $V \cup \{X_i'\}$ is

$$Imp_{\infty,\infty}^{1,p}(X_{j}) = \sum_{k=0}^{p-1} \frac{p-k}{p+1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{D}_{\nu}(V^{-j})} I(X_{j}; Y|B)$$
(5.8)

Proof. See [Louppe, 2014, Page 147] for a proof.

As observed in Theorem 5.3, the sum of all importance scores is equal to $I(X_1, \ldots, X_p; Y)$. The addition of X'_i does not actually modify $I(X_1, \ldots, X_j, \ldots, X_p; Y)$ which is equal to

 $I(X_1,...,X_j,X_j',...,X_p;Y)^4$. All importances, including those of non-redundant features, are therefore modified so that the sum of all importances remains the same

Equation 4.18 shows that the importance of a variable decreases if it is totally redundant with other features. Indeed, the addition of a new feature increase the number of feature combinations B and thus the number of terms $(I(X_i; Y|B))$ in the sum. This reflects in the weights of the outer sum of Equation 5.8. Indeed, all weights

$$\frac{1}{C_p^k(p-k)} \text{ are multiplied by a factor } \frac{C_p^k(p-k)}{C_{p+1}^k} = \frac{p-k}{p+1} < 1 \text{ that updates weights}$$

to take into account the new feature, i.e. the ensemble of trees is now built on p+1variables instead of p. Mathematically, the importance of X_i however decreases. By definition of total redundancy, X_j becomes useless if X'_i is given making all those new terms where X'_i is included in B equal to zero. Moreover, X'_i does not either increase the information conveyed by X_i about the target and thus all terms $I(X_i; Y|B)$ where X'_i is not included in B are unchanged. One may notice that the impact of the addition of a totally redundant feature is not simply a division of the original importance score of X_i into X_i and X'_i .

Proposition 5.10. [Louppe, 2014, Proposition 7.4] Let $X_i \in V$ be a relevant variable with respect to Y and V and let $X'_i \notin V$ be a totally redundant variable with respect to X_j . The infinite sample size importance of $X_l \in V^{-j}$ as computed with an infinite ensemble of fully developed totally randomized trees built on $V \cup \{X_i'\}$ is

$$\begin{split} \textit{Imp}_{\infty,\infty}^{1,p}(X_{l}) &= \sum_{k=0}^{p-2} \frac{p-k}{p+1} \frac{1}{C_{p}^{k}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-l} \setminus X_{j})} I(X_{l}; Y | B) + \\ & \hookrightarrow \sum_{k=0}^{p-2} \left[\sum_{k'=1}^{2} \frac{C_{2}^{k'}}{C_{p+1}^{k+k'}} \frac{1}{p+1-(k+k')} \right] \sum_{B \in \mathcal{P}_{k}(V^{-l} \setminus X_{j})} I(X_{l}; Y | B \cup X_{j}). \end{split}$$

Proof. See [Louppe, 2014, Pages 148-149] for a proof.

First, let us note that X_j and X_j' are identical and thus they can be used interchangeably or together without modifying the link between other features and the target. Mathematically, for any conditioning set B and for any variable X_1 , we have that $I(X_l; Y|B, X_i) = I(X_l; Y|B, X_i) = I(X_l; Y|B, X_i, X_i)$. It is the reason why Equation 5.9 is divided in two parts: those terms that do not involve either X_i nor X_i' and those whose B necessarily includes X_i (which is equivalent to include X_i' or both).

Equation 5.9 shows the impact on a non-redundant variable X_i . The first part concerns all B made without V^{-i} . Corresponding conditional mutual information terms are decreased by a factor

$$\frac{C_{\mathfrak{p}}^{k}(\mathfrak{p}-k)}{C_{\mathfrak{p}+1}^{k}(\mathfrak{p}+1-k)} = \frac{\mathfrak{p}-k}{(\mathfrak{p}+1)} \frac{1}{\mathfrak{p}+1-k} < 1.$$

Similarly to Equation 5.8, the first sub-factor updates weights to take into account the additional feature.

The second part concerns all B involving either X_i or X_i' and it shows that the corresponding conditional mutual information weights are accentuated, implying

⁴It can be shown by applying chain rule (I($X_1, X_2, \ldots, X_p; Y$) = $\sum_{i=1}^p I(X_i; Y|X_{i-1}, \ldots, X_1)$) I($X_1, \ldots, X_j, X_j', \ldots, X_p; Y$) while finishing by X_j' . Therefore, the last term is I($X_j'; Y|X_1, \ldots, X_j, \ldots, X_p$) which is, by definition of total redundancy, equal to zero. Then by applying the chain rule backward, we obtain $I(X_1,...,X_j,...,X_p;Y)$.

an increase of importances. Indeed, because of $I(X_l;Y|B,X_i)=I(X_l;Y|B,X_j)=I(X_l;Y|B,X_i,X_j)$, the same B is actually taken into account several times (two more in this case). The net effect of those two parts on the importance of X_j is a trade-off between those two antagonist effects which depends on the interaction of X_j with X_i . Indeed, features that are positively affected by the presence of X_i , e.g. features such that $I(X_j;Y|X_i)>I(X_j;Y)$, may end up with increased importances while importances of features that are either not or negatively impacted by X_i will accordingly decrease (because the fixed value for the sum of all importances).

Without further proof, Louppe [2014] extends Proposition 5.9 and 5.10 to consider the addition of N_c totally redundant features with X_j with respect to Y. Concretely, the effects given above are the same but amplified by the presence of N_c totally redundant features instead of two.

Proposition 5.11. [Louppe, 2014, Proposition 7.5] Let $X_j \in V$ be a relevant variable with respect to Y and V and let $X_j^c \notin V$ (for $c = 1, ..., N_c$) be N_c totally redundant variables with respect to X_j . The infinite sample size importances of X_j and $X_l \in V$ as computed with an infinite ensemble of fully developed totally randomized trees built on $V \cup \{X_1^1, \ldots, X_l^{N_c}\}$ are

$$\begin{split} \textit{Imp}_{\infty,\infty}^{1,p}(X_j) &= \sum_{k=0}^{p-1} \left[\frac{C_p^k(p-k)}{C_{p+N_c}^k(p+N_c-k)} \right] \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-j})} I(X_j;Y|B), \\ \textit{Imp}_{\infty,\infty}^{1,p}(X_l) &= \sum_{k=0}^{p-2} \left[\frac{C_p^k(p-k)}{C_{p+N_c}^k(p+N_c-k)} \right] \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-l} \setminus X_j)} I(X_l;Y|B) + \\ & \hookrightarrow \sum_{k=0}^{p-2} \left[\sum_{k'=1}^{N_c+1} \frac{C_{N_c+1}^{k'}}{C_{p+N_c}^{k+k'}} \frac{1}{p+N_c-(k+k')} \right] \sum_{B \in \mathcal{P}_k(V^{-l} \setminus X_j)} I(X_l;Y|B \cup X_j). \end{split}$$

5.5 NON-TOTALLY RANDOMISED TREES

Setting of this section: K > 1, D = p, $|s_t| = |v(s_t)|$, $N_T \to \infty$, $N \to \infty$

In practice, random forest methods (e.g., Random Forest [Breiman, 2001] or Extra-Trees [Geurts et al., 2006]) are rarely built with K=1 because the growing procedure is then made independently of the data, and may lead to useless tree structures especially if the number of irrelevant features is large. Note that in the case of infinite ensemble size, and assuming that ties are broken deterministically, trees built with K=p (i.e., the maximal value) amount to build classical single trees in a deterministic way.

In contrast with totally randomised trees (with K=1), masking effects may appear when trees are built with K>1. The masking effect denotes situations where several candidate splits on different features yield roughly the same impurity reduction, but one of the features is always slightly better so that none of the other ones has a chance to be selected by the tree-growing algorithm. Note that with multiway splits in particular, each feature is associated to one potential impurity decrease. Some variables may never be selected because some other variables always yield larger impurity decreases, and may thus be "masked". Such effects tend to use first the best variables (in the sense of those yielding the largest impurity decrease at first) while pushing the least promising (i.e., yielding small impurity decreases in comparison to the best ones) towards the leaves. This implies that all feature combinations

are no longer considered: best features are considered alone or conditioned only with the best others used before while the least promising ones are only considered conditioned on most of all other variables. As a result, some branches are never explored and the importance of a variable no longer decomposes into a sum including all $I(X_m; Y|B)$ terms.

To make things clearer, let us consider a simple example. Let X_1 be a variable that perfectly explains Y and let X_2 be a slightly noisy copy of X_1 (i.e., $I(X_1;Y) \approx$ $I(X_2;Y)$, $I(X_1;Y|X_2) = \epsilon$ and $I(X_2;Y|X_1) = 0$). Using totally randomized trees, the importances of X_1 and X_2 are nearly equal – the importance of X_1 being slightly higher than the importance of X_2 :

$$\begin{split} & \operatorname{Imp}_{\infty,\infty}^{1,p}(X_1) &= \frac{1}{2}\operatorname{I}(X_1;Y) + \frac{1}{2}\operatorname{I}(X_1;Y|X_2) = \frac{1}{2}\operatorname{I}(X_1;Y) + \frac{\varepsilon}{2} \\ & \operatorname{Imp}_{\infty,\infty}^{1,p}(X_2) &= \frac{1}{2}\operatorname{I}(X_2;Y) + \frac{1}{2}\operatorname{I}(X_2;Y|X_1) = \frac{1}{2}\operatorname{I}(X_2;Y) + 0 \end{split}$$

In non-totally randomized trees, for K = 2, X_1 is always selected at the root node and X_2 is always used in its children. Also, since X_1 perfectly explains Y, all its children are pure and the reduction of entropy when splitting on X_2 is null. As a result, $\text{Imp}_{\infty,\infty}^{K=2,p}(X_1) = I(X_1;Y)$ and $\text{Imp}_{\infty,\infty}^{K=2,p}(X_2) = I(X_2;Y|X_1) = 0$. Masking effects are here clearly visible: the true importance of X_2 is masked by X_1 as if X_2 were irrelevant, while it is only a bit less informative than X_1 . In the same way, it can also be shown that the importances become dependent on the number of irrelevant variables. Let us indeed consider the following example: let us add in the previous example an irrelevant variable X_i with respect to $\{X_1, X_2\}$ and let us keep K = 2. The probability of selecting X_2 at the root node now becomes positive, which means that $\mathrm{Imp}_{\infty,\infty}^{K=2,p}(X_2)$ now includes $\mathrm{I}(X_2;Y)>0$ and is therefore strictly larger than the importance computed before. For K fixed, adding irrelevant variables dampens masking effects, which thereby makes importances indirectly dependent on the number of irrelevant variables.

Consequently, non-totally randomised trees may be unable to identify all relevant features unlike totally randomised trees (see Corollary 5.6). The following proposition however guarantees that all strongly relevant features will still be identified.

Proposition 5.12. [Sutera et al., 2018]

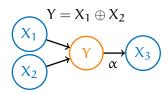
$$\forall K, X_{\mathfrak{m}} \in V: \quad X_{\mathfrak{m}} \text{ strongly relevant } \Rightarrow Imp_{\infty,\infty}^{K,p}(X_{\mathfrak{m}}) > 0.$$

Proof. See proof of Theorem 5.16 with the particular case of q = p.

There is thus no masking effect possible for the strongly relevant features when K > 1. For a given K, the features found will thus include all strongly relevant variables and some (when K > 1) or all (when K = 1) weakly relevant ones. It is easy to show that increasing K can only decrease the number of weakly relevant variables found. Using K = 1 will thus provide a solution for the **all-relevant** problem, while increasing K will provide a better and better approximation of the minimal-optimal problem in the case of strictly positive distributions (see Section 2.4.1).

While strongly relevant variables can not be masked, their importances are not necessarily higher than the importances of weakly relevant variables, i.e., Xi strongly relevant and X_i weakly relevant does not imply that $Imp_{\infty,\infty}^{K,p}(X_i) \ge Imp_{\infty,\infty}^{K,p}(X_i)$. Example 5.1 illustrates this. Unfortunately, strongly relevant variables can thus not be distinguished from weakly relevant ones only using importances.

Example 5.1. Let us consider a problem defined by three binary input variables X_1,X_2 and X_3 , and a binary output Y.



The relationships between input and output variables are the following:

- $Y = X_1 \oplus X_2$ and Y is therefore completely determined by X_1 and X_2 ;
- $X_3 = Y$ with probability α and its value is randomly chosen otherwise (i.e., $X_3 = 0$ with probability $(1 - \alpha)/2$ and $X_3 = 1$ with probability $(1 - \alpha)/2$).

In this case, X₁ and X₂ are strongly relevant with respect to Y while X₃ is only weakly relevant because it is useless when X_1 and X_2 are both known.

For $\alpha = 0.8$, we can compute that $Imp(X_1) = Imp(X_2) = 0.296$ and $Imp(X_3) =$ 0.408. Let us note that for small values of α (e.g., $\alpha = 0.2$), $Imp(X_3) < Imp(X_1) =$ $Imp(X_2)$.

In conclusion, the importances as derived from trees with non-totally randomised split selection do not possess the same properties as those computed with totally randomised trees. The ability to identify all relevant features and the independence with respect to the addition or removal of irrelevant features are both lost. Asymptotically, the use of totally randomised trees seems more appropriate for assessing the importance of features.

But in a finite setting (i.e., a limited number of samples and a limited number of trees), $I(X_m; Y|B)$ terms are not all considered neither for all X_m nor for all B, and/or need to be empirically estimated. Therefore, the use of non-totally randomised trees may help to focus on informative features providing better trees and splits on those features with more samples. Let us note that it could also be of interest in order to avoid useless splits on irrelevant features. Assessing feature importances with K > 1 therefore remains a sound strategy in practice even if some features might be missed and the resulting importances may be biased.

5.6 NON-FULLY DEVELOPED TREES

Setting of this section: K > 1, $D = q(\langle p), |s_t| = |v(s_t)|, N_T \to \infty, N \to \infty$

One key assumption of Theorem 5.3 was that all features are used once in every branch of the tree. However, when trees are no longer fully developed and say limited to a maximal depth q (< p), all combinations are no longer explored and therefore we investigate in this section the ability of identifying relevant features with importance scores derived from pruned trees.

Proposition 5.13. [Louppe et al., 2013, Proposition 6] The importance of $X_m \in V$ for Y as computed with an infinite ensemble of pruned totally randomized trees built up to depth $q \leq p$ and an infinitely large training sample is:

$$Imp_{\infty,\infty}^{1,q}(X_{\mathfrak{m}}) = \sum_{k=0}^{q-1} \frac{1}{C_{\mathfrak{p}}^{k}} \frac{1}{\mathfrak{p} - k} \sum_{B \in \mathcal{P}_{\mathfrak{p}}(V^{-\mathfrak{m}})} I(X_{\mathfrak{m}}; Y|B)$$
 (5.10)

Proof. See [Louppe et al., 2013, Appendix G] for a proof.

Proposition 5.14. [Louppe et al., 2013, Proposition 7] The importance of $X_m \in$ V for Y as computed with an infinite ensemble of pruned totally randomized trees built up to depth q

p and an infinitely large training sample is identical to the importance as computed for Y with an infinite ensemble of fully developed totally randomized trees built on random subspaces of q variables drawn from V.

Proof. See [Louppe et al., 2013, Appendix G] for a proof.

Given Proposition 5.1, the degree of a variable X can not be larger than r-1 and thus as soon as $r \leq q$, we have the guarantee that all relevant variables can be identified with totally randomised trees (K = 1).

Proposition 5.15. If $r \leq q$: $Imp_{\infty,\infty}^{1,q}(X_m) > 0$ iff X is relevant.

Proof. Given Proposition 5.1, for all, and only, the relevant variables, there exists at least one subset B of size |B| < r such that $I(X_m; Y|B) > 0$. The proposition then follows from the fact that Equation 5.10 contains all conditional mutual information terms $I(X_m; Y|B)$ with |B| < r when $r \le q$.

In the case of non-totally randomized trees (K > 1), we lose the guarantee to find all relevant variables even when $r \leq q$. Indeed, there is potentially a masking effect due to K > 1 that might prevent the conditioning needed for a given variable to be relevant to appear in a tree branch. However, we have the following general result:

Theorem 5.16. $\forall K$, if $r \leqslant q$: X_m strongly relevant $\Rightarrow \text{Imp}_{\infty,\infty}^{K,q}(X_m) > 0$

Proof. By definition, $\operatorname{Imp}_{\infty,\infty}^{K,q}(X_m) > 0$ means that there is at least one tree (grown with parameters q and K) in which X_m receives a strictly positive score for its split, i.e. such that Y depends on $X_{\mathfrak{m}}$ conditionally to the variable assignment defined by the path from the root node to the node where X_{m} is used to split. Let us show that one such tree always exists whatever K when X_m is strongly relevant and $r\leqslant q$.

Within the infinite ensemble, let us consider only the trees such that irrelevant variables are tested in each branch only when all relevant variables (including X_m) are exhausted. These trees are always explored whatever the value of K. This derives from the fact that a relevant variable can always be picked with non zero probability at any tree node, except if all relevant variables have been tested above that node. Indeed, except in this latter case, the K tested variables can always include at least one relevant variable. If some relevant variable gets a non zero score, one relevant variable will be automatically used to split since irrelevant variables can only get zero scores. Even when all tested relevant variables get a zero score, one of them can still be selected instead of an irrelevant one given that ties are resolved by randomisation.

Let us denote by τ_R the set of trees as just defined and let us show that X_m gets a non zero score in at least one tree in τ_R .

By definition of relevance and proposition 5.1, X_m strongly relevant implies that there exists at least one assignment of values to all relevant variables but X_m such that conditionally to this assignment, Y is dependent on X_m . In each tree in τ_R , there is a path from the root node to a node where $X_{\rm m}$ is used to split that is compatible with this assignment. Let us assume that X_m always gets a zero score in all these compatible paths and show that this leads to a contradiction.

If all relevant variables are tested above X_m in a compatible path then X_m should receive a non zero score at its node, which would contradict our hypothesis. Thus, X_m can only be tested in a compatible path before all relevant variables have been tested. Given our hypothesis that X_m only gets zero scores, if X_m is used to split in one compatible path, then there exists another tree in τ_R with the same splits above X_m in the compatible path and with the split on X_m replaced by a split on another relevant variables (because of tie randomization or because of the randomisation due to the use of a K < p). In this new tree, X_m is thus used to split at least one level below in the compatible path. Applying this argument recursively, one can thus show that there is at least one tree in τ_R where X_m is the last variable used to split in the compatible path. In this tree, X_m thus gets a non zero score, which contradicts the hypothesis and therefore concludes the theorem.

There is thus no masking effect possible for the strongly relevant features when K > 1 as soon as the number of relevant features is lower than q.

When q < r, we do not have the guarantee any more to explore all minimal conditionings required to find all (strongly or not) relevant variables, whatever the values of K. We nevertheless still have the guarantee to find all (strongly) relevant variables of degree lower than q (proofs are straightforward from proofs of Proposition 5.15 and Theorem 5.16):

Proposition 5.17.

 $X_m \in V$ relevant with respect to Y and $\deg(X_m) < q \Rightarrow \operatorname{Imp}_{\infty,\infty}^{1,q}(X_m) > 0$.

Proposition 5.18.

 $X_{\mathfrak{m}} \in V$ strongly relevant with respect to Y and $deg(X_{\mathfrak{m}}) < q \Rightarrow Imp_{\infty,\infty}^{K,q} > 0$.

5.7 BINARY TREES

Setting of this section: $|s_t| = 2$, $N_T \to \infty$, $N \to \infty$

The last simplifying assumption on tree model is the number of nodes created when splitting a node. So far, we considered multiway trees (i.e., with exhaustive splits) where one branch was created for each value of the split variable. This way of growing trees allows to consider a variable in a branch only once and limits the maximal depth of a tree to the number of features. It also implies that once a variable is used for splitting in a node, all subsequent nodes have access to all the information (about the target) held by this variable.

However, binary trees are most often used in practice. Instead of creating a branch per value, only two branches are created regardless of the cardinality of the split variable. The splitting rule is from now on of the form of a boolean condition (e.g., "less than a given threshold value" or not, "in a subset of values" or not) where samples verifying this condition go in one branch while the others necessarily go in the other branch. As a consequence, a variable can now be used several times in a given tree branch, since a variable potentially only partially delivers its information at each split. There are also now several ways to split a node on the basis of a categorical variable of cardinality greater than two. When growing a tree, a binary split can be determined for such variable either by identifying among a set of predefined candidate binary splits the one that maximizes the impurity reduction (as in the

standard Random Forests method) or by picking one binary split at random among these candidates (as in the Extra-Trees method).

As a consequence of these changes, one can not expect that Theorem 5.3 and formula 5.2 that were derived in the case of multiway trees will remain valid in the case of binary trees (except, trivially, if all variables are binary). And indeed, Example 5.2, taken from Louppe [2014], shows that importances computed from binary trees can be different from importances computed from multiway trees.

Example 5.2. We present here the example as it is given in [Louppe, 2014] and we refer the reader to the original source for more details on the exact computation of importance scores. Let us consider two ordered input variables of different cardinalities: X_1 is a ternary variable (i.e., its cardinality is 3) and X_2 is a binary variable. The output variable Y is defined as $Y = X_1 < 1$ and as a copy of X_2 . The possible combinations of values are given in Table 5.2.

$$\begin{array}{c|ccc}
X_1 & X_2 & Y \\
\hline
0 & 0 & 0 \\
1 & 1 & 1 \\
2 & 1 & 1
\end{array}$$

Table 5.2: Possible combinations of values for X_1, X_2 and Y.

Only two totally randomised trees with multiway splits can be built from this setting as a single node split is sufficient to exhaust a variable of any cardinality (either X1 or X2) and to fully determine the output value. Importances derived from such trees (in asymptotic conditions) are as follows:

$$Imp(X_1) = \frac{1}{2}I(X_1;Y) = \frac{1}{2}H(Y) = 0.459;$$

$$Imp(X_2) = \frac{1}{2}I(X_2;Y) = \frac{1}{2}H(Y) = 0.459.$$

Despite different trees, features are used in exactly half of the trees with the same usefulness and thus their importances are logically identical. Note that since both features perfectly explain the output, their importances do not depend on K.

On the other hand, a binary split can not exhaust X₁ all at once. Using ordered binary splits, four possible decision trees can now be constructed. Assuming that the Extra-Trees split randomization is used and that K is set to 1, the importances of X_1 and X_2 are respectively (in asymptotic conditions):

$$\begin{split} Imp(X_1) &= \frac{1}{4}I(X_1\leqslant 1;Y) + \frac{1}{8}P(X_1\leqslant 1)I(X_1\leqslant 0;Y|X_1\leqslant 1) + \frac{1}{4}I(X_1\leqslant 0;Y) \\ &= 0.375, \\ Imp(X_2) &= \frac{1}{2}I(X_2;Y) + \frac{1}{8}P(X_1\leqslant 1)I(X_2;Y|X_1\leqslant 1) \\ &= 0.541, \end{split}$$

which are strictly different from the importance scores derived from multiway splits.

In this section, our aim is to revisit some of our previous results in the context of binary trees. In Section 5.7.1, we discuss different ways to generate binary splits for unordered and ordered categorical variables, focusing only on sets of candidate binary splits that are totally redundant with the original variable. Example 5.2 shows that importance scores computed with binary trees can be different from those computed with multiway trees. In Section 5.7.2, we show that the links between variable relevances and variable importances that were highlighted in Sections 5.3 and 5.5 are preserved despite this difference. In Section 5.7.3, we illustrate further how binary splits influence variable importance scores on Breiman's digit recognition problem.

Binary splits 5.7.1

As defined in Section 3.2.2.1, binary splits may or may not take into account the value logic, i.e., a potential ordering between the values. An unordered split simply divides all the values into two disjoint sets, while an ordered split creates two partitions consisting of all the values that are respectively either lower or equal, or greater than a given threshold.

In the case of binary variables, both ways of splitting are strictly equivalent. In the case of variables of higher cardinality, they lead to different numbers of candidate splits. For example, there are only two possible ways of splitting a ternary variable of values $\{1, 2, 3\}$ while preserving the order (i.e., $\{\{1\}, \{2, 3\}\}$) and $\{\{1, 2\}, \{3\}\}$). By contrast, there are three possible ways of making two disjoint sets of values if the order is not taken into account (the split $(\{1,3\},\{2\})$) being the additional binary partition that does not preserve the order). In general, a categorical variable of cardinality m will lead to $2^{m-1}-1$ candidate unordered binary splits and to m-1 candidate ordered binary splits.

In addition to these two kinds of binary splits, let us also mention a third one based on the principle of "one value vs. all", where each binary split isolates one value of the variable in one branch and all the others in the other branch. In the case of a ternary variable, it provides the same candidate splits as the unordered binary splits (i.e., ({1}, {2,3}), ({2}, {1,3}), ({3}, {1,2})) but for variables of higher cardinalities, less splits are considered than in the unordered case (see e.g., Figures 5.3a and 5.4a). For a variable of cardinality m, it leads to m candidate binary splits.

All three ways of defining binary splits actually replace a categorical variable X_m by a set of new binary variables, each one corresponding to a candidate binary split defined on X_m . Let us denote by $T_m = \{T_{m,1}, \ldots, T_{m,|T_m|}\}$ the set of binary variables of size $|T_m|$ defined by one of these three families of binary splits. Figures 5.2b, 5.3b and 5.4b illustrate the three sets of binary variables corresponding to the different ways of defining binary splits described above, and Figures 5.2a, 5.3a and 5.4a illustrate all possible splits, in the case of a quaternary variable X_m .

In all three cases, it is easy to show that $T_{\mathfrak{m}}$ and $X_{\mathfrak{m}}$ are totally redundant with respect to the target Y, i.e., mathematically (see Definition 2.13):

$$\forall B\subseteq V^{-\mathfrak{m}}, \quad X_{\mathfrak{m}}\perp \!\!\!\perp Y|B\cup T_{\mathfrak{m}} \quad \text{and} \quad T_{\mathfrak{m}}\perp \!\!\!\perp Y|B\cup \{X_{\mathfrak{m}}\}. \tag{5.11}$$

Thus, collectively, variables in T_m convey the exact same information about the output as the original variable X_m from which they are derived. There is thus no loss in information when replacing multiway splits with binary splits in all three cases. Note that in the case of unordered and one-value-vs-all splits, there are redundancy in T_m in the sense that some variables can be removed from T_m without impacting its total redundancy with X_m .

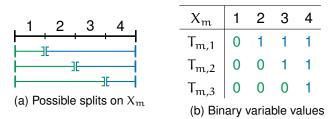


Figure 5.2: Set T_{num} of binary variables corresponding to possible ordered splits, i.e. between two successive values of X_m. Each colour is associated to one of the two branches leaving the node after the spit. For instance, intervals of values in green correspond to the left branch whereas intervals in blue correspond to the right one.

1 0 2 4	$X_{\mathfrak{m}}$	1	2	3	4
1 2 3 4	$T_{m,1}$	0	1	1	1
	$T_{m,2}$	1	0	1	1
	$T_{m,3}$	1	1	0	1
(a) Possible splits on X _m	$T_{m,4}$	1	1	1	0
(a) I coolid opino on \mathcal{M}_{III}	(b) Bin	(b) Binary variable values			

Figure 5.3: Set Toh of binary variables corresponding to possible unordered "one value vs. all" splits, i.e. one-hot encoding of values of X_m. Each colour is associated to one of the two branches leaving the node after the spit. For instance, intervals of values in green correspond to the left branch whereas intervals in blue correspond to the right one.

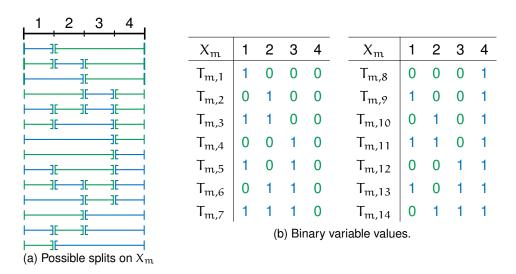


Figure 5.4: Set T_{bp} of binary variables corresponding to possible unordered splits, i.e. all binary partitions of values of X_m. Each colour is associated to one of the two branches leaving the node after the spit. For instance, intervals of values in green correspond to the left branch whereas intervals in blue correspond to the right one.

5.7.2 Relevance in binary trees

In this section, we assume that a binary tree is grown from a set of categorical variables using the Extra-Trees split randomization, i.e., by randomly selecting K variables at each node, picking for each of them a random binary split in its set of candidate binary splits, and finally using the split among K that leads to the most important decrease of impurity (breaking ties at random). The importance of a variable X_m is then obtained by summing total impurity reductions at all nodes where a binary split has been performed on X_m .

In this setting, we would like first to check whether Theorem 5.5, stating that a variable is irrelevant if and only if its infinite sample size importance is 0, remains valid when using (fully developed totally randomized) binary trees instead of multiway ones.

Let us denote by X_m a categorical variables of cardinality greater than 2 and by T_m a set of totally redundant binary variables corresponding to the candidate binary splits used for this variable during tree growing. The following theorem first shows that X_m is relevant if and only if at least one variable in T_m is relevant.

Proposition 5.19. Let $X_m \in V$ be an input variable and let $T_m = \{T_{m,1}, \ldots, T_{m,|T_m|}\}$ be a set of binary variables such that X_m and T_m are totally redundant with respect to Y. There exists a subset $B\subseteq V^{-\mathfrak{m}}$ such that $I(X_{\mathfrak{m}};Y|B)>0$ if and only if there exists a subset $B \subseteq V^{-m}$, a variable $T_{m,i} \in T_m$ and a subset $T' \subseteq T_m \setminus \{T_{m,i}\}$ such that $I(T_i; Y|B \cup T') > 0$.

Proof. Necessary condition: $(\exists B: I(X_m; Y|B) > 0 \Rightarrow \exists B, T_{m,i}, T': I(T_{m,i}; Y|B \cup Proof.)$ T') > 0

As a consequence of the total redundancy between T_m and X_m , we directly have that

$$I(T_m; X_m | B) = I(X_m; Y | B) > 0.$$

Applying the chain rule on $I(T_m; X|B) = I(T_{m,1}, ..., T_{m,|T_m|}; Y|B)$, we have that

$$I(T_m; X_m | B) = \sum_{i=1}^{|T_m|} I(T_{m,i}; Y | B \cup \{T_{m,1}, \dots, T_{m,i-1}\}) > 0$$

which implies that a least one term of the sum should be strictly positive. That is,

$$\exists T_{m,i}: I(T_{m,i}; Y|B \cup T') > 0.$$

where $T' = \{T_{m,1}, ..., T_{m,i-1}\}.$

Sufficient condition: $(\exists B, T_{m,j}, T' : I(T_{m,j}; Y|B \cup T') > 0 \Rightarrow \exists B : I(X_m; Y|B) > 0)$ Given $I(T_{m,j};Y|B \cup T') > 0$, the proof is a direct consequence of the chain rule where variables in T' are used first and then $T_{m,i}$. Indeed,

$$I(T_m;Y|B) = \sum_{j=1}^{|T_m|} I(T_{m,j};Y|B \cup \{T_{m,1},\ldots,T_{m,j-1}\})$$

is therefore necessarily strictly positive and thus

$$I(T_m; Y|B) = I(X_m; Y|B) > 0$$

by total redundancy.

The following proposition further shows that variables whose relevance is conditioned on X_m will remain relevant conditionally to some variables in a totally redundant set T_m .

Proposition 5.20. For any relevant variable $X_i \in V^{-m}$ with respect to Y, there exists a subset B such that $I(X_i; Y|B \cup X_m) > 0$ if and only if there exists a subset $\mathsf{T}'\subseteq\mathsf{T}$ such that $\mathsf{I}(\mathsf{X}_{\mathsf{i}};\mathsf{Y}|\mathsf{B}\cup\mathsf{T}')$ where T is a set of binary variables which is totally redundant with X_m with respect to Y.

Proof. The proof is a direct consequence of the total redundancy between $X_{\mathfrak{m}}$ and T_{m} .

Propositions 5.19 and 5.20 can be combined to show that Theorem 5.5 remains valid in the case of fully developed totally randomized binary trees, when candidate binary splits are totally redundant with the original variables.

Theorem 5.21. Let us assume binary trees constructed by using totally redundant candidate binary splits and the Extra-Trees split randomization. Then, $X_i \in V$ is irrelevant to Y with respect to V if and only if its infinite sample size importance as computed with an infinite ensemble of fully developed totally randomized binary trees built on V for Y is 0.

We do not provide a formal proof of this theorem to not overload the text. Intuitively, the theorem can be proven by noting that when K = 1 and with split randomization, the importance of a variable X_m is a weighted sum of all possible terms $I(T_{m,i};Y|B)$, where $T_{m,i}$ is a binary split based on X_m and B is a subset of binary splits defined on all features (including X_m). Given Propositions 5.19 and 5.20, at least one such term is strictly positive if and only if X_m is relevant.

In the case of multiway trees, Proposition 5.12 shows that strongly relevant variables will be always found whatever the value of K. A similar result can be shown in the case of binary trees.

Let us first characterize the relevance of binary variables in T_m with respect to the relevance of X_m . The following corollary of Proposition 5.19 first shows that if X_m is only weakly relevant, no variable in a totally redundant set T_m can be strongly relevant.

Corollary 5.22. If X_m is weakly relevant with respect to Y, then each $T_{m,j} \in T_m$, with X_m and T_m totally redundant with respect to Y, is either irrelevant or weakly relevant with respect to Y.

Proof. The relevance of some $T_{m,j} \in T_m$ directly results from Proposition 5.19. No $T_{m,j}$ can however be strongly relevant. Indeed, if X_m is weakly relevant with respect to Y, we have that $X_m \perp \!\!\! \perp Y | V^{-m}$ which is equivalent to $T_m \perp \!\!\! \perp Y | V^{-m} \Leftrightarrow$ $T_{m,1},\ldots,T_{m,|T_m|}\perp\!\!\!\perp Y|V^{-m},$ given the total redundancy between T_m and X_m . By weak union, the latter independence implies that:

$$T_i \perp\!\!\!\perp Y | V^{-\mathfrak{m}} \cup T^{-\mathfrak{i}}$$

for all $T_i \in T$.

 $X_{\mathfrak{m}}$ strongly relevant does not ensure that a variable in a totally redundant set $T_{\mathfrak{m}}$ will be strongly relevant (Surely, this can not be the case if T_m contains redundant features), which would have sufficed to show that Proposition 5.12 remains valid for binary trees. However, the following results show that at least one $T_{m,i} \in T_m$ can not be masked by variables in V^{-m} .

Proposition 5.23. Let X_m be a strongly relevant variable with respect to Y and let $T_m = \{T_{m,1}, \dots, T_{m,|T_m|}\}$ be a set of binary variables such that X_m and T_m are totally redundant with respect to Y. There exists at least one variable $T_{m,i} \in T_m$ such that $T_{m,i} \not\perp Y | V^{-m} \cup T'$ for at least one subset $T' \subseteq T_m \setminus \{T_{m,i}\}$.

Proof. Let us assume that one such T_{m,i} does not exist and show that this leads to a contradiction. For all $T_{m,i} \in T_m$ and all $T' \subseteq T_m \setminus \{T_{m,i}\}$ (possibly empty), we thus have $T_{m,i} \perp Y | V^{-m} \cup T'$. Let us consider any ordering of the variables in T_m and let us recursively apply the contraction property. We then have the following sequence of independences: $T_{m,1} \perp Y|V^{-m}$, $T_{m,2} \perp Y|V^{-m} \cup T_{m,1}$ and $T_{m,1} \perp I$ $Y|V^{-m} \text{ gives } \{T_{m,1},T_{m,2}\} \perp \!\!\! \perp Y|V^{-m},\ldots,T_{m,|T_m|} \perp \!\!\! \perp Y|V^{-m} \cup \{T_{m,1},\ldots,T_{m,|T_m|-1}\} \text{ and }$ $\{T_{m,1},\ldots,T_{m,|T_m|-1}\} \perp \!\!\! \perp Y|V^{-m}$ gives $\{T_{m,1},\ldots,T_{m,|T_m|}\} \perp \!\!\! \perp Y|V^{-m}$. The latter independence is impossible because of the strong relevance of X_m that implies that $X_m \not\perp Y|V^{-m}$ and thus $T_m \not\perp Y|V^{-m}$, given that T_m and X_m are totally redundant.

Using this result, one can adapt the proof of Proposition 5.12 in a straightforward way to show the following result (provided without proof):

Proposition 5.24. Let us assume binary trees constructed by using totally redundant candidate binary splits and the Extra-Trees split randomization.

$$\forall K, X_m \in V : X_m \text{ strongly relevant } \Rightarrow \operatorname{Imp}_{\infty,\infty}^{K,p}(X_m) > 0.$$

Theorem 5.21 and Proposition 5.24 thus show that using binary instead of multiway splits fortunately does not affect the ability of variable importances to identify the relevant features and filter out the irrelevant ones.

Importance scores in binary trees 5.7.3

Through Example 5.2, we already know that importance scores are expected to be different in binary trees compared to multiway trees. In this section, we further illustrate this difference in more details by computing variable importance scores in various settings on Breiman et al. [1984]'s digit recognition problem (see Appendix C for a description of this problem).

Figures 5.5a and 5.5b show the importance scores computed respectively from totally randomised (i.e., K = 1) and non-totally randomised (i.e., K = p) Extra-Trees, in which the split is randomly selected (split-wise randomisation). Both figures compare multiway trees and binary trees with either unordered or ordered binary splits. While all variables are binary in the original problem, we artificially increased the cardinality of variable X₁ from 2 to 4 by splitting both values 0 and 1 of this variable each into two new values, respectively {1,2} and {3,4}, with equal probability. This transformation does not change the information brought by X_1 about Y but it allows us to illustrate the effect of the different binary split strategies on importance scores.

When $|X_1| = 2$, all tree growing methods lead to the same importance scores for all variables as expected (the slight differences are due to the use of a finite number of trees). The importance of X_1 is nevertheless decreased when K goes from 1 to p, due to masking effects. When the cardinality of X₁ is increased to 4, we notice that the three splitting strategies lead to different importance scores. With ordered splits (see Figure 5.2 for all candidate splits), all candidate splits are somehow useful because they all provide part (or all for the mid-split) of the information content

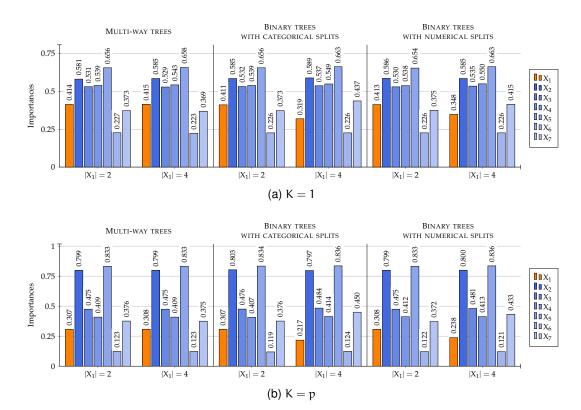


Figure 5.5: Importance scores as computed by an ensemble of 10000 trees with K = 1 (top) and K = p (bottom), with multiway trees (left), binary trees with (unordered) categorical splits (center), and binary trees with ordered splits (right). The considered problem is the digit recognition problem of [Breiman et al., 1984]. $|X_1| = 2$ corresponds to the original problem with only binary variables while $|X_1| = 4$ corresponds to the same problem where the cardinality of X₁ has been artificially increased to 4 (with both values 0 and 1 splitted each into two new values, {0, 1} and $\{2,3\}$ respectively, with equal probability). The six other variables X_2, \ldots, X_7 are left unchanged.

of X₁. By contrast, there are much more candidate unordered splits (see Figure 5.4 for all of them) including several ones that do not provide any information about Y (e.g., the split ($\{1,3\},\{2,4\}$) does not change the distribution of Y). With K=1, split variables are selected totally at random. In the case of ordered splits, any variable except X₁ that is used is granted for all its information while X₁ only receives its full importance in one third of all splits. In the case of unordered splits, the chance of X_1 to be granted of its full importance is even smaller because of the useless splits. In both cases, this gives more opportunity to another variable to capture part of the information contained in X₁ about Y and hence leads to a reduction of the importance of X₁ in the case of binary trees (with respect to multiway trees). A similar effect is observed when K = p. Because of the split randomisation, some splits on X_1 will be uninformative and in such case, X₁ will not be chosen to split the node to the benefit of another variable, leading to an overall decrease of the importance of X₁. Interestingly, the importances of all variables except X_7 are mostly unchanged whatever the splitting strategy. The importance of X_7 is however increased when going from multiway to binary trees with $|X_1| = 4$. This is a consequence of the high redundancy between variables X_1 and X_7 : they are equal for all digits except 7 (see Appendix C). X_7 is thus the variable which benefits the most from the irrelevant splits on X_1 introduced by the binary trees.

Note that importance scores would be different if splits were optimized, instead of randomized, for each variable, as in the standard Random Forests method. In the case of our example, multiway and binary trees would have given the exact same importance scores for all variables even when $|X_1| = 4$, since the optimal split would always be the split $(\{1,2\},\{3,4\})$. It is possible however to design problems where the Random Forests node splitting strategy will make importance scores derived from multiway trees different from importance scores derived from binary trees.

IN NON-ASYMPTOTIC CONDITIONS

Setting of this section: $N_T \nrightarrow \infty$, $N \nrightarrow \infty$

From now on, we do no longer consider asymptotic conditions. This section aims at examining the importance measure in finite settings and investigate results of this chapter in this context. Section 5.8.1 considers a finite number of trees. This suggests that all possible branches (i.e., not masked) are not necessarily explored and/or fairly taken into account. Section 5.8.2 considers a finite number of samples. It implies that $I(X_m; Y|B)$ can not be computed exactly and must be empirically estimated from samples.

5.8.1 With a finite number of trees

As mentioned in Section 4.4.3, in practice, the number of trees in a Random Forest ensemble should be as large as possible in order to achieve the best predictive performances. At some point however, a plateau should be reached and adding more trees will not increase significantly the performance. The impact on feature importance is usually not taken into account however. In this section, we still assume a learning sample of infinite size and study the impact of the number of trees on the properties highlighted so far. We only examine fully developed trees but results in

this section can be easily generalised to non-fully developed trees given the analysis in Section 5.6.

As presented in Equation 5.1, the importance of a feature is computed over all trees and over all nodes of all trees. With an infinite number of trees, we saw in Theorem 5.3 that the relationship between X_m and Y is evaluated for all combinations B in such a way that all terms equally contribute to the total importance. When only a finite number N_T of trees is constructed, some conditionings B (branches) can be missed and thus the importance will only contain a subset of all $I(X_m; Y|B)$ terms. However, we have the following general result:

Proposition 5.25.
$$\forall K, q, \operatorname{Imp}_{\infty, N_{\top}}^{K,q} > 0 \Rightarrow X_{\mathfrak{m}}$$
 is relevant.

Features with strictly positive importance scores are necessarily relevant, since it implies that at least one term $I(X_m; Y|B) > 0$. However, a relevant feature does not necessarily have positive importance score, even a strongly relevant one. In all generality, Proposition 5.12 is thus not valid with a finite number of trees. To give an example, let us consider a XOR scenario with two features X₁ and X₂ such that $I(X_1;Y) = I(X_2;Y) = 0$ and $I(Y;X_1,X_2) > 0$. If a single tree is grown, only the feature tested at the second level will receive a non-zero importance, while both features are (strongly) relevant.

This observation suggests that an undesirable effect of using a finite number of trees is that features that are not examined (or not with the right conditioning set B) have zero importances. Unseen features therefore wrongly appear as irrelevant with respect to Y, like masked features or those with too high degree. Note however that if the composition property is verified, then a single tree (with K = p) can identify all strongly relevant features because strongly relevant features can not be masked.

Theorem 5.26. If K = p and if $P_{V,Y}(V,Y)$ verifies the composition property: $X_m \in V$ *strongly relevant* $\Rightarrow \operatorname{Imp}_{\infty,1}^{p,q}(X_{\mathfrak{m}}) > 0$.

Proof. We want to show that a single tree that is fully developed with K = q is sufficient to give to all strongly relevant features a strictly positive importance score when the distribution over all variables verifies the composition property. Since the tree is fully developed, all features are exhausted in each branch and each leaf corresponds to a possible assignment v to all input features V. Let us assume that a strongly relevant variable X_m does not have a strictly positive importance score and show that this leads to a contradiction.

If X_m does not receive a strictly positive importance score, it means that X_m is never used in a terminal node corresponding to a configuration v^{-m} such that $X_m \not\perp Y|V^{-m} = v^{-m}$ or in an internal node corresponding to a configuration b of $B \subset V^{-m}$ such that $X_m \not\perp \!\!\! \perp Y|B=b$. By definition of strong relevance, we have $X_m \not\perp Y | V^{-m}$ which implies that there exists at least one configuration $V^{-m} = v^{-m,*}$ such that $X_m \not\perp Y|V^{-m} = v^{-m,*}$. Let us consider the path in the tree from the root node to a node where $X_{\mathfrak{m}}$ is tested that matches the values in $\nu^{-\mathfrak{m},*}$. $X_{\mathfrak{m}}$ can not be tested at the end of such path because otherwise it would have got a strictly positive importance score. The node $X_{\mathfrak{m}}$ is thus necessarily used in the path in a node corresponding to a configuration $B = b^*$ that matches for some variables $B \subset V^{-m}$ the configuration $v^{-m,*}$ and such that $X_m \perp \!\!\! \perp Y | B = b^*$. In the same conditioning $B = b^*$, all features in $R = V^{-m} \setminus B$ are also independent of Y conditionally to $B = b^*$, otherwise one of them would have been preferred to X_m to split the node (since K = p means that they were all evaluated when splitting

the node). Given the composition property, we thus have that $(X_m \cup R) \perp \!\!\! \perp Y \mid B = b^*$. The weak union property then implies that $X_m \perp \!\!\! \perp Y | (B = b^*) \cup R$, which means that $X_m \perp \!\!\! \perp Y | (B = b^*) \cup (R = r)$ for all configurations r of the variables in R. This is thus also true for the configuration r^* that matches the configuration $v^{-m,*}$, which shows that $X_m \perp \!\!\! \perp Y | V^{-m} = v^{-m,*}$. This is however impossible by definition of $v^{-m,*}$.

In the same vein, Wehenkel [2018] computed analytically the minimum number of trees such that all features are at least seen once (among the K features selected at a given node) for a given value K. This analysis showed that many trees are needed, in particular when K is small and individual decision trees are small. Note that having seen all features once is obviously not enough to identify all relevant variables, as they need to be tested at least in one of their minimal conditioning sets B and furthermore not to be masked in such case by other variables. The number of trees given in [Wehenkel, 2018] is thus a very minimal bound on the number of trees really needed to find all relevant variables.

Moreover, let us note that even if the number of trees is large enough to consider all possible branches, computed importances with a finite forest are most likely different from theoretical asymptotic importances because all B may not be fairly considered in the forest.

5.8.2 With a finite number of samples

In all analyses carried out so far, assuming a sample set of infinite size actually corresponds to know the data distribution and therefore to compute with exactitude all measures, e.g. node impurity i(t) and node decrease $\Delta i(s,t)$. However, in practice, impurity measurements are estimated from a finite sample set and therefore suffer from an empirical misestimation bias. Concretely, it means that Equation 5.2 of Theorem 5.3 becomes, if we still assume an infinite number of trees,

$$Imp_{N,\infty}^{1,p}(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} \hat{I}(X_m; Y|B)$$
 (5.12)

where $\hat{I}(X_m;Y|B)$ are estimated mutual informations.

Among other authors, Goebel et al. [2005] show that mutual information estimation between two independent variables is positively biased. That is, let us consider two independent discrete random variables X and Y of probability density $P_X(X)$ and $P_Y(Y)$ respectively and such that I(X;Y)=0, their finite sample size estimates $\hat{I}(X;Y)$ are expected to be strictly positive, i.e.,

$$\mathbb{E}\{\hat{\mathbf{I}}(X;Y)\} = \frac{(|X|-1)(|Y|-1)}{2N\ln(2)} > 0$$
 (5.13)

where N is the number of observed samples, |X| and |Y| are respectively the cardinalities of X and Y. In contrast with Theorem 5.5, this however suggests that irrelevant features never have zero importances. Louppe [2014] stresses the linear dependence with variable cardinalities and the inverse dependence of the number of samples, and relates with many empirical studies that observe a bias towards feature of large number of categories and cardinalities [Strobl et al., 2007b].

In details, given three random variables X,Y,Z of probability densities $P_X(X)$, $P_Y(Y)$, P_Z(Z) respectively, Goebel et al. [2005] show that the estimator for conditional mutual information $\hat{I}(X;Y|Z)$ is approximately gamma distributed

$$\hat{I}(X;Y|Z) \sim \Gamma\left(\frac{|Z|}{2}(|X|-1)(|Y|-1), \frac{1}{N\ln(2)}\right)$$
 (5.14)

where $\Gamma(k,\theta)$ is the gamma distribution with a shape parameter k and a scale parameter θ . Let us note that a random variable W such that $W \sim \Gamma(\nu/2, 2c)$ with c > 0, then W also follows a χ^2 (chi-square) distribution⁵ with ν degrees of freedom. In the case of $\hat{I}(X;Y|Z)$, it then follows a χ^2 distributions of |Z|(|X|-1)(|Y|-1) degrees of freedom (and $c = \frac{1}{2N \ln(2)}$). That is, we have that $2N \ln(2) \hat{I}(X;Y|Z)$ converges asymptotically towards a χ^2 distribution with |Z|(|X|-1)(|Y|-1) degrees of freedom, that only depends on feature cardinalities. One can then use a chi-square based statistical test on the mutual information between two features to determine if their are independent. Let us once again note that the number of degrees of freedom increase with features cardinalities.

To avoid false positives, all those results suggest to combine non-totally developed trees, in order not to estimate mutual informations from too few samples at deep nodes, with non-totally randomised trees (K > 1), in order to avoid splitting on irrelevant features at the top nodes, which would unnecessarily reduce the size of the learning sample. Unfortunately, as the previous analyzes show, decreasing tree depth or increasing K will however increase the number of false negatives. There is thus a tradeoff to be found in practice between these two antagonistic effects.

5.9 **RESULT SUMMARY**

The following table summarises the main results exposed in this chapter, with references to the main theorems.

⁵Saporta [2006] define a χ^2 law as follows: Let U_1, U_2, \ldots, U_p be p independent variables, each following $\mathfrak{N}(0,1)$, then the chi-square law with p degrees of freedom, denoted χ_p^2 , is the law of the variable $\sum_{i=1}^{p} U_i^2$.

		K = 1		K > 1		
		D = p	D < p	D = p	D < p	
	Analytical formulation	Thm. 5.3	Prop. 5.13 Prop. 5.14	_	 — 	
Theoretical results	Sum of importances	Thm. 5.4		_		
Importance vs. Relevance in asymptotic conditions	Irrelevant variables	$\Leftrightarrow Imp = 0$ Thm. 5.5	$\lim_{n \to \infty} 1 = 0$	\Rightarrow Imp = 0	$\Rightarrow Imp = 0$	
	Relevant variables		$ \Rightarrow \text{Imp} > 0^{6}$ $ \text{if } r \leqslant q$ $ \text{Prop. 5.15}$	⇒ Imp ≥ 0 Masking effect		
	Strongly relevant variables	⇔ Imp > 0 Cor. 5.6		⇒ Imp > 0 Prop. 5.12	$ \Rightarrow \text{Imp} > 0$ $ \text{if } r \leqslant q^7$ $ \text{Thm. 5.16}$	
	Presence of irrelevant variables	No effect Lem. 5.7, Thm. 5.8		Dampens masking effect		
Binary splits		Same relevance but different importance scores				
Finite settings	Finite number of trees	r $\forall K, D Imp_{\infty,N_T}^{K,D}(X_{\mathfrak{m}}) > 0 \Rightarrow X_{\mathfrak{m}} \text{ is relevant} \\ Prop. \ 5.25$				
	Finite number of samples	misestimation higs $(Imn > 0 \Rightarrow relevance)$				

✓ Chapter take-away

In asymptotic conditions, MDI feature importances can be derived analytically and provide an understandable decomposition of the total information conveyed by input features about the target by feature, by cardinality of the interaction term, and by interaction term. Additionally, the sum of all importances is fixed. The introduction of redundant feature tends to modify all importance scores and not only those of features that conveyed redundant informations. When trees are built totally at random, zero importances are only associated to irrelevant features. When trees are not totally random, the masking effect prevents some weakly relevant feature to be identified and only strongly relevant features are ensured to have positive importance scores. When trees are non-totally developed, interaction terms of larger cardinalities are no longer evaluated and therefore do not enter into account in the importance scores. However, guarantees can be preserved by restricting the number of relevant features or the feature degree. In more realistic and practical settings (i.e., binary trees, finite sample set, finite number of trees), those desirable properties are however usually not preserved.

⁶When r > q, this is only valid for relevant variables X_m such that $deg(X_m) < q$ (see Prop. 5.17).

⁷When r > q, only strongly relevant variables X_m such that $deg(X_m) < q$ can not be masked (see Prop. 5.18).

Part III

EXTENSIONS AND DERIVATIONS OF IMPORTANCE MEASURES

Overview

In this chapter, we extend the random forest feature importances framework to perform a contextual analysis. For many problems, feature selection is often more complicated than identifying a single subset of input features that would together explain the output, as described in Chapter 2 especially. There may be interactions that depend on contextual information, i.e., variables that reveal to be relevant only in some specific circumstances. We briefly discussed in Section 2.4.7 that such feature interactions must be taken into account but a single feature ranking provides only very limited information about such complex relationships. In this setting, our contribution is to extend the MDI feature importance measure (i) to identify variables whose relevance is context-dependent, and (ii) to characterise as precisely as possible the effect of contextual information on the importance of these variables.

References: This chapter is an adapted version of the following publication:

A. Sutera, G. Louppe, V. A. Huynh-Thu, L. Wehenkel, and P. Geurts. Context-dependent feature analysis with random forests. In *Uncertainty In Artificial Intelligence: Proceedings of the Thirty-Second Conference*, 2016.

Terminology and notations have been slightly adjusted for the sake of consistency with the rest of this manuscript. The text has also been processed to minimize overlap with respect to previous chapters.

6.1 MOTIVATION

Supervised learning finds applications in many domains such as medicine, economics, computer vision, or bioinformatics. Given a sample of observations of several inputs and one output variable, the goal of supervised learning is to learn a model for predicting the value of the output variable given any values of the input variables. Another common side objective of supervised learning is to bring as much insight as possible about the relationship between the inputs and the output variable. One of the simplest ways to gain such insight is through the use of feature selection or ranking methods that identify the input variables that are the most decisive or relevant for predicting the output, either alone or in combination with other variables. Among feature selection/ranking methods, one finds variable importance scores derived from random forest models that stand out from the literature mainly because of their multivariate and non parametric nature and their reasonable computational cost. Although very useful, feature selection/ranking methods however only provide very limited information about the often very complex input-output relationships that can be modeled by supervised learning methods. There is thus a high interest in designing new techniques to extract more complete information about input-output relationships than a single global feature subset or feature ranking.

In this chapter, we specifically address the problem of the identification of the input variables whose relevance or irrelevance for predicting the output only holds in specific circumstances, where these circumstances are assumed to be encoded by a specific context variable. This context variable can be for example a standard input variable, in which case, the goal of contextual analyses is to better understand how this variable interacts with the other inputs for predicting the output. The context can also be an external variable that does not belong to the original inputs but that may nevertheless affect their relevance with respect to the output. Practical applications of such contextual analyses are numerous. In some applications, one may be interested in finding variables that are both relevant and independent of the context. For example, in medical studies [see, e.g., Geissler et al., 2000], one is often interested in finding risk factors that are as independent as possible of external factors, such as the sex of the patients, their origins or the data cohort to which they belong. By contrast, in some other cases, one may be interested in finding variables that are relevant but dependent in some way on the context. For example, in systems biology, differential analysis [Ideker and Krogan, 2012] aims at discovering genes or factors that are relevant only in some specific conditions, tissues, species or environments.

Our contribution in this chapter is two-fold. First, starting from common definitions of feature relevance in the literature, we propose a formal definition of contextdependent variables and provide a complete characterization of these variables depending on how their relevance is affected by the context variable. Second, we extend the random forest variable importances framework in order to identify and characterize variables whose relevance is context-dependent or context-independent. Building on existing theoretical results for standard importance scores, we propose asymptotic guarantees for the resulting new measures with respect to the formal definitions.

The chapter is structured as follows. In Section 6.2, we first lay out our formal framework defining context-dependent variables and describing how the context may change their relevance. We describe in Section 6.3 how random forest variable importances can be used for identifying context-dependent variables and how the effect of contextual information on these variables can be highlighted. Our results are then illustrated in Section 6.4 on representative problems. Finally, conclusions and directions of future works are discussed in Section 6.5.

6.2 CONTEXT-DEPENDENT FEATURE SELECTION AND CHARACTERIZATION

Let us consider a set $V = \{X_1, ..., X_p\}$ of p input CONTEXT-DEPENDENCE. variables and an output Y and let us denote by V^{-m} the set $V \setminus \{X_m\}$. All input and output variables are assumed to be categorical, not necessarily binary¹. Let us reconsider the definitions of relevant, irrelevant, and marginally relevant variables based on their mutual information I (as defined in Definitions and 2.8).

Let us now assume the existence of an additional (observed) context variable $X_c \notin V$, also assumed to be categorical.

Inspired by the notion of relevant and irrelevant variables, we propose to define context-dependent and context-independent variables as follows:

¹The case of a non categorical output will be discussed in Section 6.3.5.

Definition 6.1. A variable $X_m \in V$ is context-dependent to Y with respect to X_c iff there exists a subset $B \subseteq V^{-\mathfrak{m}}$ and some values \mathfrak{x}_c and \mathfrak{b} such that \mathfrak{b} :

$$I(Y; X_m | B = b, X_c = x_c) \neq I(Y; X_m | B = b).$$
 (6.1)

Definition 6.2. A variable $X_{\mathfrak{m}} \in V$ is context-independent to Y with respect to X_c iff for all subsets $B\subseteq V^{-\mathfrak{m}}$ and for all values x_c and b, we have:

$$I(Y; X_m | B = b, X_c = x_c) = I(Y; X_m | B = b).$$
(6.2)

Context-dependent variables are thus the variables for which there exists a conditioning set B in which the information they bring about the output is modified by the context variable. Context-independent variables are the variables that, in all conditionings B = b, bring the same amount of information about the output whether the value of the context is known or not. This definition is meant to be as general as possible. Other more specific definitions of context-dependence are as follows:

$$\exists B \subseteq V^{-m}, b, x_c^1, x_c^2 : I(Y; X_m | X_c = x_c^1, B = b) \neq I(Y; X_m | X_c = x_c^2, B = b),$$
 (6.3)

$$\exists B \subseteq V^{-m}, x_c : I(Y; X_m | X_c = x_c, B) \neq I(Y; X_m | B),$$

$$(6.4)$$

$$\exists B \subseteq V^{-m}, b: I(Y; X_m | X_c, B = b) \neq I(Y; X_m | B = b),$$
 (6.5)

$$\exists B \subseteq V^{-m}: I(Y; X_m | X_c, B) \neq I(Y; X_m | B).$$

$$(6.6)$$

These definitions all imply context-dependence as defined in Definition 6.1 but the converse is in general not true. For example, Definition (6.3) misses problems where the context makes some otherwise irrelevant variable relevant but where the information brought by this variable about the output is exactly the same for all values of the context. A variable that satisfies Definition (6.1) but not Definition (6.4) is given in example 6.1. This example can be easily adapted to show that both Definitions (6.5) and (6.6) are more specific than Definition (6.1) (by swapping the roles of X_c and X_2).

Example 6.1. This artificial problem is defined by two input variables X_1 and X_2 , an output Y, and a context X_c . X_1 , X_2 , and X_c are binary variables taking their values in {0, 1}, while Y is a quaternary variable taking its values in {0, 1, 2, 3}. All combinations of values for X₁, X₂, and X_c have the same probability of occurrence 0.125 and the conditional probability $P(Y|X_1, X_2, X_C)$ is defined by the two following rules:

- If $X_2 = X_c$ then $Y = X_1$ with probability 1.
- If $X_2 \neq X_c$ then Y = 2 with probability 0.5 and Y = 3 with probability 0.5.

The corresponding data table is given in Appendix 6.A. For this problem, it is easy to show that $I(Y; X_1 | X_2 = 0, X_c = 0) = 1$ and that $I(Y; X_1 | X_2 = 0) = 0.5$, which means

²In this definition and all definitions that follow, we will assume that the events on which we are conditioning have a non-zero probability and that if such event does not exist then the condition of the definition is not satisfied.

condition (6.1) is satisfied and X_1 is thus context-dependent to Y with respect to $X_{
m c}$ according to our definition. On the other hand, we can show that:

$$I(Y; X_1|X_c = x_c) = I(Y; X_1) = 0.5$$

 $I(Y; X_1|X_2, X_c = x_c) = I(Y; X_1|X_2) = 0.5,$

for any $x_c \in \{0,1\}$, which means that condition (6.4) can not be satisfied for X_1 .

To simplify the notations, the context variable was assumed to be a separate variable not belonging to the set of inputs V. It can however be considered as an input variable, whose own relevance to Y (with respect to $V \cup \{X_c\}$) can be assessed as for any other input. Let us examine the impact of the nature of this variable on context-dependence. First, it is interesting to note that the definition of contextdependence is not symmetric. A variable X_m being context-dependent to Y with respect to X_c does not imply that the variable X_c is context-dependent to Y with respect to X_m . Second, the context variable does not need to be marginally relevant for some variable to be context-dependent, but it needs however to be relevant to Y with respect to V. Indeed, we have the following theorem:

Theorem 6.1. X_c is irrelevant to Y with respect to V iff all variables in V are contextindependent to Y with respect to X_c (and V) and $I(Y; X_c) = 0$.

As a consequence of this theorem, there is no interest in looking for contextdependent variables when the context itself is not relevant⁴.

CHARACTERIZING CONTEXT-DEPENDENT VARIABLES. Contextual analyses need to focus only on context-dependent variables since, by definition, contextindependent variables are unaffected by the context: their relevance status (relevant or irrelevant), as well as the information they contain about the output, remain indeed unchanged whatever the context.

Context-dependent variables may be affected in several directions by the context, depending both on the conditioning subset B and on the value x_c of the context. Given a context-dependent variable X_m , a subset B and some values b and x_c such that $I(Y; X_m|B = b, X_c = x_c) \neq I(Y; X_m|B = b)$, the effect of the context can either be an increase of the information brought by X_m (I(Y; $X_m | B = b, X_c =$ $(x_c) > I(Y; X_m | B = b)$ or a decrease of this information $(I(Y; X_m | B = b, X_c = x_c) < b$ $I(Y; X_m | B = b)$). Furthermore, for a given variable X_m , the direction of the change can differ from one context value x_c to another (at fixed B and b) but also from one conditioning B = b to another (for a fixed context x_c). Example 6.2 below illustrates this latter case. This observation makes a global characterization of the effect of the context on a given context-dependent variable difficult. Let us nevertheless mention two situations where such global characterization is possible:

Definition 6.3. A context-dependent variable $X_m \in V$ is context-complementary (in a context x_c) iff for all $B \subseteq V^{-m}$ and b, we have $I(Y; X_m | B = b, X_c = x_c) \geqslant$ $I(Y; X_m | B = b).$

³This would be the case however if we had adopted the definition (6.6).

⁴This is consistent with Proposition 5.1. All features in a minimal conditioning subset of B are necessarily relevant, including any contextual features.

Definition 6.4. A context-dependent variable $X_{\mathrm{m}} \in V$ is **context-redundant** (in a context x_c) iff for all $B \subseteq V^{-m}$ and b, we have $I(Y; X_m | B = b, X_c = x_c) \leqslant$ $I(Y; X_m | B = b)$.

Context-complementary and redundant variables are variables that always react in the same direction to the context and thus can be characterized globally without loss of information. Context-complementary variables are variables that bring complementary information about the output with respect to the context, while contextredundant variables are variables that are redundant with the context. Note that context-dependent variables that are also irrelevant to Y are always context-complementary, since the context can only increase the information they bring about the output. Context-dependent variables that are relevant to Y however can be either contextcomplementary, context-redundant, or uncharacterized. A context-redundant variable can furthermore become irrelevant to Y (with respect to $V \cup \{X_c\}$) as soon as $I(Y; X_m | B = b, X_c = x_c) = 0$ for all B, b, and x_c .

Example 6.2. As an illustration, in the problem of Example 6.1, X₁ and X₂ are both relevant and context-dependent variables. X1 can not be characterized globally since we have simultaneously:

$$\begin{split} & I(Y; X_1 | X_2 = 0, X_c = x_c) > I(Y; X_1 | X_2 = 0) \\ & I(Y; X_1 | X_2 = 1, X_c = x_c) < I(Y; X_1 | X_2 = 1), \end{split}$$

for both $x_c = 0$ and $x_c = 1$. X_2 is however context-complementary as the knowledge of X_c always increases the information it contains about Y.

RELATED WORKS. Several authors have studied interactions between variables in the context of supervised learning. They have come up with various interaction definitions and measures, e.g., based on multivariate mutual information [McGill, 1954; Jakulin and Bratko, 2003a], conditional mutual information [Jakulin, 2005; Van de Cruys, 2011], or variants thereof [Brown, 2009; Brown et al., 2012]. There are several differences between these definitions and ours. In our case, the context variable has a special status and as a consequence, our definition is inherently asymmetric, while most existing variable interaction measures are symmetric. In addition, we are interested in detecting any information difference occurring in a given context (i.e., for a specific value of X_c) and for any conditioning subset B, while most interaction analyses are interested in average and/or unconditional effects. For example, [Jakulin and Bratko, 2003a] propose as a measure of the interaction between two variables X_1 and X_2 with respect to an output Y the multivariate mutual information, which is defined as $I(Y;X_1;X_2) = I(Y;X_1) - I(Y;X_1|X_2)$. Unlike our definition, this measure can be shown to be symmetric with respect to its arguments. Adopting this measure to define context-dependence would actually amount at using condition (6.6) instead of condition (6.1), which would lead to a more specific definition as discussed earlier in this section.

The closest work to ours in this literature is due to Turney [1996], who proposes a definition of context-sensitivity that is very similar to our definition of contextdependence. Using our notations, Turney [1996] defines a variable X_m as weakly context-sensitive to the variable X_c if there exist some subset $B \subseteq V^{-m}$ and some values y, x_m , b, and x_c such that these two conditions hold:

$$p(Y = y|X_m = x_m, X_c = x_c, B = b) \neq p(Y = y|X_m = x_m, B = b),$$

$$p(Y = y | X_m = x_m, X_c = x_c, B = b) \neq p(Y = y | X_c = x_c, B = b).$$

 X_m is furthermore defined as strongly context-sensitive to X_c if X_m is weakly sensitive to X_c , X_m is marginally relevant, and X_c is not marginally relevant. These two definitions do not exactly coincide with ours and they have two drawbacks in our opinion. First, they do not consider that a perfect copy of the context is contextsensitive, which we think is counter-intuitive. Second, while strong context-sensitivity is asymmetric, the constraints about the marginal relevance of X_m and X_c seems also unnatural.

Our work is also somehow related to several works in the graphical model literature that are concerned with context-specific independences between random variables [see e.g. Boutilier et al., 1996; Zhang and Poole, 1999]. Boutilier et al. [1996] define two variables Y and X_m as contextually independent given some $B \subseteq V^{-m}$ and a context value x_c as soon as $I(Y; X_m | B, X_c = x_c) = 0$. When $B \cup \{X_m, X_c\}$ are the parents of node Y in a Bayesian network, then such context-specific independences can be exploited to simplify the conditional probability tables of node Y and to speed up inferences. Boutilier et al. [1996]'s context-specific independences will be captured by our definition of context-dependence as soon as $I(Y; X_m|B) > 0$. However, our framework is more general as we want to detect any context dependencies, not only those that lead to perfect independences in some context.

6.3 CONTEXT ANALYSIS WITH RANDOM FORESTS

In this section, we show how to use variable importances derived from Random Forests first to identify context-dependent variables (Section 6.3.2) and then to characterize the effect of the context on the relevance of these variables (Section 6.3.3). Derivations in this section are based on the theoretical characterization of variable importances provided in [Louppe et al., 2013], which is briefly reminded in Section 6.3.1. Section 6.3.4 discusses practical considerations and Section 6.3.5 shows how to generalize our results to other impurity measures.

6.3.1 Variable importances 5

Within the random forest framework, Breiman [2001] proposed to evaluate the importance of a variable X_m for predicting Y by adding up the weighted impurity decreases for all nodes t where X_m is used, averaged over all N_T trees in the forest:

$$Imp(X_m) = \frac{1}{N_T} \sum_{t \in T: \nu(s_t) = X_m} p(t) I(Y; X_m | t)$$
(6.7)

where $v(s_t)$ is the variable used in the split s_t at node t, p(t) is the proportion of samples reaching t and I is the mutual information.

According to Louppe et al. [2013], for any ensemble of fully developed trees in asymptotic learning sample size conditions, the Mean Decrease Impurity (MDI) importance (6.7) can be shown to be equivalent to

$$Imp(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} I(Y; X_m | B), \tag{6.8}$$

⁵This section is a reminder of the MDI importance measure and its asymptotic characterisation. See Section 5.2 for more details.

where V^{-m} denotes the subset $V \setminus \{X_m\}$, $\mathcal{P}_k(V^{-m})$ is the set of subsets of V^{-m} of size k. where $\mathcal{P}_k(V^{-m})$ denotes the set of subsets of V^{-m} of size k. Most notably, it can be shown [Louppe et al., 2013] that this measure is zero for a variable X_m iff X_m is irrelevant to Y with respect to V. It is therefore well suited for identifying relevant features.

6.3.2 Identifying context-dependent variables

Theorem 6.1 shows that if the context variable X_c is irrelevant, then it can not interact with the input variables and thus modify their importances. This observation suggests to perform, as a preliminary test, a standard random forest variable importance analysis using all input variables and the context in order to check the relevance of the latter. If the context variable does not reveal to be relevant, then, there is no hope to find context-dependent variables.

Intuitively, identifying context-dependent variables seems similar to identifying the variables whose importance is globally modified when the context is known. Therefore, one first straightforward approach to identify context-dependent variables is to build a forest per value $X_c = x_c$ of the context variable, i.e., using only the data samples for which $X_c = x_c$, and also globally, i.e. using all samples and not including the context among the inputs. Then it consists in deriving from these models an importance score for each value of the context, as well as a global importance score. Context-dependent variables are then the variables whose global importance score differs from the contextual importance scores for at least one value of the context.

More precisely, let us denote by $Imp(X_m)$ the global score of a variable X_m computed using (6.7) from all samples and by $Imp(X_m|X_c = x_c)$ its importance score as computed according to (6.7) using only those samples such that $X_c = x_c$. With this approach, a variable would be declared as context-dependent as soon as there exists a value x_c such that $\text{Imp}(X_m) \neq \text{Imp}(X_m|X_c=x_c)$.

Although straightforward, this approach has several drawbacks. First, in the asymptotic setting of Section 6.3.1, it is not guaranteed to find all context-dependent variables. Indeed, asymptotically, it is easy to show from (6.8) that $Imp(X_m)$ – $Imp(X_m|X_c = x_c)$ can be written as:

$$\begin{split} Imp^{x_{c}}(X_{m}) &\triangleq Imp(X_{m}) - Imp(X_{m}|X_{c} = x_{c}) \\ &= \sum_{k=0}^{p-1} \frac{1}{C_{k}^{p}} \frac{1}{p-k} \sum_{B \in \mathcal{P}_{k}(V^{-m})} (I(Y;X_{m}|B) - I(Y;X_{m}|B,X_{c} = x_{c})). \end{split}$$
 (6.9)

Example 6.1 shows that $I(Y; X_m|B)$ can be equal to $I(Y; X_m|B, X_c = x_c)$ for a context-dependent variable. Therefore we have the property that if there exists an x_c such that $Imp^{x_c}(X_m) \neq 0$, then the variable is context-dependent but the opposite is unfortunately not true. Another drawback of this approach is that in the finite case, we do not have the guarantee that the different forests will have explored the same conditioning sets B and therefore, even assuming that the learning sample is infinite (and therefore that all mutual informations are perfectly estimated), we lose the guarantee that $\operatorname{Imp}^{x_c}(X_m) \neq 0$ for a given x_c implies context-dependence.

To overcome these two problems, we propose the following new importance score to identify context-dependent variables:

$$Imp^{|x_c|}(X_m) \triangleq \frac{1}{N_T} \sum_{T} \sum_{t \in T: v(s_t) = X_m} p(t) |I(Y; X_m|t) - I(Y; X_m|t, X_c = x_c)|$$
 (6.10)

This score is meant to be computed from a forest of totally randomized trees built from all samples, not including the context variable among the inputs. At each node t where the variable X_m is used to split, one needs to compute the absolute value of the difference between the mutual information between Y and X_m estimated from all samples reaching that node and the mutual information between Y and X_m estimated only from the samples for which $X_c = x_c$. The same forest can then be used to compute $Imp^{|x_c|}(X_m)$ for all x_c . A variable X_m is then declared contextdependent as soon as there exists an x_c such that $Imp^{|x_c|}(X_m) > 0$.

Let us show that this measure is sound. In asymptotic conditions, i.e., with an infinite number of trees, one can show from (6.10) that $Imp^{|x_c|}(X_m)$ becomes:

$$Imp^{|x_c|}(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} \sum_{b \in \mathcal{B}} P(B = b)$$

$$\hookrightarrow |I(Y; X_m | B = b) - I(Y; X_m | B = b; X_c = x_c)|.$$
(6.11)

Asymptotically, this measure has now the very desirable property to not miss any context-dependent variable as formalized in the next theorem:

Theorem 6.2. A variable $X_m \in V$ is context-independent to Y with respect to X_c iff $\operatorname{Imp}^{|\mathbf{x}_{\rm c}|}(X_{\rm m})=0$ for all $\mathbf{x}_{\rm c}$.

Given that the absolute differences are computed at each tree node, this measure also continues to imply context-dependence in the case of finite forests and infinite learning sample size. The only difference with the infinite forests is that only some conditionings B and values b will be tested and therefore one might miss the conditionings that are needed to detect some context-dependent variables.

Characterizing context-dependent variables

Besides identifying context-dependent variables, one would want to characterize their dependence with the context as precisely as possible. As discussed in Section 6.3, irrelevant variables (i.e, such that $Imp(X_m) = 0$) that are detected as contextdependent do not need much effort to be characterized since the context can only increase their importance. All these variables are therefore context-complementary.

Identifying the context-complementary and context-redundant variables among the relevant variables that are also context-dependent can in principle be done by simply comparing the absolute value of $Imp^{x_c}(X_m)$ with $Imp^{|x_c|}(X_m)$, as formalized in the following theorem:

Theorem 6.3. If $|\operatorname{Imp}^{x_c}(X_m)| = \operatorname{Imp}^{|x_c|}(X_m)$ for a context-dependent variable X_m , then X_m is context-complementary if $Imp^{x_c}(X_m) < 0$ and context-redundant if $Imp^{x_c}(X_m) > 0.$

Proof. The absolute value of a sum is less than or equal the sum of the absolute value of each terms. The equality is only verified when all terms are of the same sign. Therefore, the sign of $Imp^{x_c}(X_m)$ indicates the sign of all terms and thus verify either the context-complementarity if all terms are negative or the context-redundancy if all terms are positive.

This result allows to identify easily the context-complementary and context-redundant variables. In addition, if, for a context-redundant variable X_m , we have $Imp^{|x_c|}(X_m) =$ $\operatorname{Imp}^{x_c}(X_m) = \operatorname{Imp}(X_m)$, then this variable is irrelevant in the context x_c .

Then it remains to characterize the context-dependent variables that are neither context-complementary nor context-redundant. It would be interesting to be able to also characterize them according to some sort of average effect of the context on these variables. Similarly as the common use of importance $Imp(X_m)$ to rank variables from the most to the less important, we propose here to use the importance $Imp^{x_c}(X_m)$ to characterize the average global effect of context x_c on the variable $X_{\rm m}$. Given the asymptotic formulation of this importance in Equation (6.10), a negative value of $Imp^{x_c}(X_m)$ means that X_m is essentially complementary with the context: in average over all conditionings, it brings more information about Y in context x_c than when ignoring the context. Conversely, a positive value of $Imp^{x_c}(X_m)$ means that the variable is essentially redundant with the context: in average over all conditionings, it brings less information about Y than when ignoring the context. Ranking the context-dependent variables according to $Imp^{x_c}(X_m)$ would then give at the top the variables that are the most complementary with the context and at the bottom the variables that are the most redundant.

Note that, like $Imp^{|x_c|}(X_m)$, it is preferable to estimate $Imp^{x_c}(X_m)$ by using the following formula rather than to estimate it from two forests by subtracting $Imp(X_m)$ and Imp $(X_m|X_c = x_c)$:

$$Imp_{s}^{x_{c}}(X_{m}) = \frac{1}{N_{T}} \sum_{t \in T: \nu(s_{t}) = X_{m}} p(t) (I(Y; X_{m}|t) - I(Y; X_{m}|t, X_{c} = x_{c}))$$
 (6.12)

This estimation method has the same asymptotic form as $Imp(X_m) - Imp(X_m|X_c =$ x_c) given in Equation (6.10) but, in the finite case, it ensures that the same conditionings are used for both mutual information measures. Note that in some applications, it is interesting also to have a global measure of the effect of the context. A natural adaptation of (6.12) to obtain such global measure is as follows:

$$\operatorname{Imp}^{X_{c}}(X_{m}) \triangleq \frac{1}{N_{T}} \sum_{T} \sum_{t \in T: v(s_{t}) = X_{m}} p(t) (I(Y; X_{m} | t) - I(Y; X_{m} | t, X_{c}))$$

which, in asymptotic sample and ensemble of trees size conditions, gives the following formula:

$$Imp^{X_c}(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_k^p} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} (I(Y; X_m | B) - I(Y; X_m | B, X_c)).$$

If $Imp^{X_c}(X_m)$ is negative then the context variable X_c makes variable X_m globally more informative (X_c and X_m are complementary with respect to Y and V). If $\text{Imp}^{X_c}(X_{\mathfrak{m}})$ is positive, then the context variable X_c makes variable $X_{\mathfrak{m}}$ globally less informative (X_c and X_m are redundant with respect to Y and V).

6.3.4 In practice

As a recipe when starting a context analysis, we suggest first to build a single forest using all input variables X_m (but not the context X_c) and then to compute from this forest all importances defined in the previous section: the global importances $Imp(X_m)$ and the different contextual importances, $Imp_s^{x_c}(X_m)$, $Imp^{|x_c|}(X_m)$, and $Imp^{X_c}(X_m)$, for all variables X_m and context values x_c .

Second, variables satisfying the context-dependence criterion, i.e., such that $Imp^{|x_c|}(X_m) > 0$ for at least one x_c , can be identified from the other variables. Among context-dependent variables, an equality between $|\operatorname{Imp}_{s}^{x_c}(X_m)|$ and $\operatorname{Imp}^{|x_c|}(X_m)$ highlights that the context-dependent variable X_m is either context-complementary or context-redundant (in x_c) depending on the sign of $Imp_{x_c}^{x_c}(X_m)$. Finally, the remaining context-dependent variables can be ranked according to $Imp_s^{x_c}(X_m)$ (or $Imp^{X_c}(X_m)$ for a more global analysis).

Note that, because mutual informations will be estimated from finite training sets, they will be generally non zero even for independent variables, leading to false positives in the identification of context-dependent variables. In practice, one could instead identify context-dependent variables by using a test $Imp^{|x_c|}(X_m) > \epsilon$ where ϵ is some cut-off value greater than 0. In practice, the determination of this cutoff can be very difficult. In our experiments, we propose to turn the importances $Imp^{|x_c|}(X_m)$ into p-values by using random permutations. More precisely, 1000 scores $Imp^{|x_c|}(X_m)$ will be estimated by randomly permuting the values of the context variable in the original data (so as to simulate the null hypothesis corresponding to a context variable fully independent of all other variables). A p-value will then be estimated by the proportion of these permutations leading to a score $Imp^{|x_c|}(X_m)$ greater than the score obtained on the original dataset.

X_c	X ₁	X_2	X_3	Υ
$\frac{X_c}{0}$	0	0	0	2
0		0	1	2
0	0 0 0	1	0	2
0	0	1	1	2
0 0 0 0 0	1	0	0	0
0	1	0	1	0
0	1	1	0	1
0	1	1	1	1
1	0	0	0	2 2
1	0 0 0	0	1	2
1	0	1	0	2
1	0	1	1	2
1	1	0	0	0
1	1	0	1	1
1	1	1	0	0
1	1	1	1	1

Table 6.1: Problem 1: Values of X_c , X_1 , X_2 , X_3 , Y.

Generalization to other impurity measures

All our developments so far have assumed a categorical output Y and the use of Shannon's entropy as the impurity measure. Our framework however can be carried

	X ₁	X_2	X_3
$Imp(X_m)$	1.0	0.125	0.125
$Imp(X_m X_c=0)$	1.0	0.5	0.0
$Imp(X_m X_c=1)$	1.0	0.0	0.5
$Imp^{ 0 }(X_m)$	0.0	0.375	0.125
$Imp^{0}(X_{\mathfrak{m}})$	0.0	-0.375	0.125
$Imp^{ 1 }(X_m)$	0.0	0.125	0.375
$Imp^{1}(X_{\mathfrak{m}})$	0.0	0.125	-0.375
$Imp^{X_c}(X_m)$	0.0	-0.125	-0.125

Table 6.2: Problem 1: Variable importances as computed analytically using asymptotic formulas. Note that X_1 is context-independent and X_2 and X_3 are contextdependent.

	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈
$Imp(X_m)$	0.5727	0.7514	0.5528	0.687	0.1746	0.0753	0.1073	0.0
$Imp(X_{\mathfrak{m}} X_{\mathfrak{c}}=0)$	0.4127	0.5815	0.5312	0.5421	0.6566	0.2258	0.372	0.0
$Imp(X_m X_c=1)$	0.6243	0.8057	0.5577	0.7343	0.0	0.0	0.0	0.0
$Imp^{ 0 }(X_m)$	0.2263	0.2431	0.1181	0.2241	0.4139	0.1961	0.2861	0.0
$\operatorname{Imp}^{ 1 }(X_{\mathfrak{m}})$	0.0987	0.0611	0.021	0.0736	0.1746	0.0753	0.1073	0.0
$Imp^0(X_m)$	0.2179	0.2422	0.1111	0.2190	-0.3839	-0.1389	-0.2346	0.0
$Imp^{1}(X_{\mathfrak{m}})$	-0.0516	-0.0543	-0.0049	-0.0473	0.1746	0.0753	0.1073	0.0

Table 6.3: Problem 2: Variable importances as computed analytically using the asymptotic formulas for the different importance measures.

over to other impurity measures and thus in particular also to a numerical output Y. Let us define a generic impurity measure $i(Y|t) \ge 0$ that assesses the impurity of the output Y at a tree node t. The corresponding impurity decrease at a tree node is defined as:

$$G(Y;X_m|t) = i(Y|t) - \sum_{x_m \in \mathcal{X}_m} p(t_{x_m})i(Y|t_{x_m})$$

$$\tag{6.13}$$

with t_{x_m} denoting the successor node of t corresponding to value x_m of X_m . By analogy with conditional entropy and mutual information, let us define the population based measures i(Y|B) and $G(Y; X_m|B)$ for any subset of variables $B \subseteq V$ as follows:

$$i(Y|B) = \sum_{b} P(B=b)i(Y|B=b)$$

$$G(Y;X_m|B) = i(Y|B) - i(Y|B,X_m),$$

where the first sum is over all possible combinations b of values for variables in B. Now, substituting mutual information I for the corresponding impurity decrease measure G, all our results above remain valid, including Theorems 1, 2, and 3 (proofs are omitted for the sake of space). It is important however to note that this substitution changes the notions of both variable relevance and context-dependence. Definition 6.1 indeed becomes:

Definition 6.5. A variable $X_m \in V$ is context-dependent to Y with respect to X_c iff there exists a subset $B \subseteq V^{-m}$ and some values x_c and b such that

$$G(Y; X_m | B = b, X_c = x_c) \neq G(Y; X_m | B = b).$$

When Y is numerical, a common impurity measure is variance, which defines i(Y|t) as the empirical variance var[Y|t] computed at node t. The corresponding $G(X_m; Y|B = b)$ and $G(X_m; Y|B = b, X_c = x_c)$ in Definition (5) are thus defined respectively as:

$$var{Y|B = b} - \mathbb{E}_{X_m|B=b}{var{Y|X_m, B = b}}$$

and

$$var\{Y|B = b, X_c = x_c\} - \mathbb{E}_{X_m|B = b, X_c = x_c}\{var\{Y|X_m, B = b, X_c = x_c\}\}.$$

We will illustrate the use of our framework in a regression setting with this measure in the next section.

6.4 EXPERIMENTS

We first illustrate the different importance measures defined in Section 6.3 on two artificial problems and then exploit them on two real bio-medical datasets.

Problem 1.

The purpose of this first problem is to illustrate the different measures introduced earlier. This artificial problem is defined by three binary input variables X_1, X_2 , and X_3 , a ternary output Y, and a binary context X_c . All samples are enumerated in Table 6.1 and are supposed to be equiprobable. By construction, the output Y is defined as Y = 2 if $X_1 = 0$, $Y = X_2$ if $X_c = 0$ and $X_1 = 1$, and $Y = X_3$ if $X_c = 1$ and $X_1 = 1$.

Table 6.2 reports all importance scores for the three inputs. These scores were computed analytically using the asymptotic formulas, not from actual experiments. Considering the global importances $Imp(X_m)$, it turns out that all variables are relevant, with X_1 clearly the most important variable and X_2 and X_3 of smaller and equal importances. According to $Imp^{|0|}(X_m)$ and $Imp^{|1|}(X_m)$, X_1 is a contextindependent variable, while X2 and X3 are two context-dependent variables. This result is as expected given the way the output is defined. For X₂ and X₃, we have furthermore $Imp^{|x_c|}(X_m) = |Imp^{|x_c|}(X_m)|$ for both values of x_c . X_2 is therefore contextcomplementary when $X_c = 0$ and context-redundant when $X_c = 1$. Conversely, X_3 is context-redundant when $X_c = 0$ and context-complementary when $X_c = 1$. X_2 is furthermore irrelevant when $X_c=1$ (since ${\rm Imp}^1(X_2)={\rm Imp}^{|1|}(X_2)={\rm Imp}(X_2)$) and X_3 is irrelevant when $X_c = 0$ (since $Imp^0(X_3) = Imp^{|0|}(X_3) = Imp(X_3)$). The values of $Imp^{X_c}(X_2)$ and $Imp^{X_c}(X_3)$ suggest that these two variables are in average complementary.

Problem 2.

This second experiment is based on an adaptation of the digit recognition problem initially proposed in Breiman et al. [1984] and reused in Louppe et al. [2013] (see Appendix C for a detailed description). The original problem contains 7 binary variables (X_1,\ldots,X_7) and the output Y takes its values in $\{0,1,\ldots,9\}$. Each input represents the on-off status of one lightning segment of a seven-segment indicator and is determined univocally from Y. To create an artificial (binary) context, we created two copies of this dataset, the first one corresponding to $X_c = 0$ and the second one

to $X_c = 1$. The first dataset was unchanged, while in the second one variables X_5 , X_6 , and X_7 were turned into irrelevant variables. In addition, we included a new variable X_8 , irrelevant by construction in both contexts. The final dataset contains 320 samples, 160 in each context.

Table 6.3 reports possible importance scores for all the inputs. Again, these scores were computed analytically using the asymptotic formulas. As expected, variable X₈ has zero importance in all cases. Also as expected, variables X_5 , X_6 , and X_7 are all context-dependent $(Imp^{|x_c|}(X_m) > 0$ for all of them). They are context-redundant (and even irrelevant) when $X_c = 1$ and complementary when $X_c = 0$. More surprisingly, variables X₁, X₂, X₃, and X₄ are also context-dependent, even if their distribution is independent from the context. This is due to the fact that these variables are complementary with variables X_5 , X_6 , and X_7 for predicting the output. Their context-dependence is thus a consequence of the context-dependence of X_5 , X_6 , X_7 , X_1 , X_2 , X_3 , and X_4 are all almost redundant when $X_c = 0$ and complementary when $X_c = 1$, which expresses the fact that they provide more information about the output when X_5 , X_6 and X_7 are irrelevant ($X_c = 1$) and less when X_5 , X_6 , and X_7 are relevant ($X_c = 0$). Nevertheless, X_8 remains irrelevant in every situation.

Problem 3.

As a third experiment, we consider bio-medical data from the Primary tumor dataset. The objective of the corresponding supervised learning problem is to predict the location of a primary tumor in patients with metastases. It was downloaded from the UCI repository [Lichman, 2013] and was collected by the University Medical Center in Ljubljana, Slovenia. We restrict our analysis to 132 samples without missing values. Patients are described by 17 discrete clinical variables (listed in the first column of Table 6.4) and the output is chosen among 22 possible locations. For this analysis, we use the patient gender as the context variable.

Table 6.4 reports variable importances computed with 1000 totally randomized trees and their corresponding p-values. According to the p-values of $Imp^{|\mathbf{x}_c|}(X_m)$, two variables are clearly emphasized for each context: importances of histologictype and neck both significantly decrease in the first context (female) and importances of peritoneum and abdominal both significantly decrease in the second context (male). While the biological relevance of these finding needs to be verified, such dependences could not have been highlighted from standard random forests importances.

Note that the same importances computed using the asymptotic formulas are provided in Table 6.E.1. Importance values are very similar, highlighting that finite forests provide good enough estimates for this problem.

Problem 4.

As a last experiment, we consider a publicly available brain cancer gene expression dataset [Verhaak et al., 2010]. This dataset collects measurements of mRNA expression levels of 11861 genes in 220 tissue samples from patients suffering from glioblastoma multiforme (GBM), the most common form of malignant brain cancer in adults. Samples are classified into four GBM sub-types: Classical, Mesenchymal, Neural and Proneural. The interest of this dataset is to identify the genes that play a central role in the development and progression of the cancer and thus improve our understanding of this disease. In our experiment, our aim is to exploit importance

		$Imp(X_m)$	Imp(X _m	$ X_{c} = x_{c} $		Imp ^x c	(X _m)			Imp _s	(X _m)	
m		-	$x_c = 0$	$x_c = 1$	$x_c = 0$	pval	$x_c = 1$	pval	$x_c = 0$	pval	$x_c = 1$	pval
0	age	0.2974	0.2942	0.2900	0.1505	0.899	0.1717	0.417	0.0032	0.938	0.0074	0.846
1	histologic-type	0.3513	0.1354	0.4005	0.2265	0.000	0.1183	0.121	0.2159	0.000	-0.0492	0.331
2	degree-of-diffe	0.4415	0.3725	0.4070	0.1827	0.680	0.1724	0.689	0.0690	0.102	0.0345	0.398
3	bone	0.2452	0.2342	0.2220	0.1088	0.396	0.0845	0.904	0.0110	0.717	0.0232	0.410
4	bone-marrow	0.0188	0.0190	0.0131	0.0128	0.892	0.0105	0.980	-0.0001	0.994	0.0057	0.682
5	lung	0.1677	0.1837	0.1420	0.1134	0.448	0.1079	0.397	-0.0160	0.605	0.0257	0.373
6	pleura	0.1474	0.1132	0.1127	0.0613	1.000	0.1026	0.097	0.0342	0.179	0.0348	0.165
7	peritoneum	0.3171	0.2954	0.2084	0.0939	0.968	0.1516	0.000	0.0216	0.710	0.1087	0.000
8	liver	0.2300	0.1844	0.2784	0.0888	0.966	0.1382	0.053	0.0456	0.134	-0.0483	0.100
9	brain	0.0466	0.0334	0.0566	0.0403	0.173	0.0279	0.814	0.0131	0.693	-0.0101	0.751
10	skin	0.0679	0.0310	0.0786	0.0426	0.922	0.0420	0.841	0.0369	0.107	-0.0107	0.663
11	neck	0.2183	0.0774	0.2255	0.1562	0.000	0.0710	0.575	0.1409	0.000	-0.0071	0.764
12	supraclavicular	0.1701	0.1807	0.1344	0.0942	0.379	0.0738	0.884	-0.0106	0.695	0.0357	0.136
13	axillar	0.1339	0.1236	0.0846	0.0748	0.214	0.0663	0.388	0.0103	0.795	0.0493	0.194
14	mediastinum	0.1826	0.1752	0.1613	0.1129	0.266	0.0867	0.853	0.0074	0.767	0.0213	0.404
15	abdominal	0.2558	0.2883	0.1512	0.1419	0.139	0.1526	0.028	-0.0325	0.368	0.1046	0.003

Table 6.4: Problem 3: Importances as computed with a forest of 1000 totally randomized trees. The context is defined by the binary context feature Sex (Sex = 0 denotes female and Sex = 1 denotes male). P-values were estimated using 1000 permutations of the context variable. Grey cells highlight p-values under the 0.05 threshold.

scores to identify interactions between genes that are significantly affected by the cancer sub-type considered as our context variable. This dataset was previously exploited by Mohan et al. [2014], who used it to test a method based on Gaussian graphical models for detecting genes whose global interaction patterns with all the other genes vary significantly between the subtypes. This latter method can be considered as gene-based, while our approach is link-based.

Following [Mohan et al., 2014], we normalized the raw data using Multi-array Average (RMA) normalization. Then, the data was corrected for batch effects using the software ComBat [Johnson et al., 2007] and then log₂ transformed. Following [Mohan et al., 2014], we focused our analysis on only two GBM sub-types, Proneural (57 tissue samples) and Mesenchymal (56 tissue samples), and on a particular set of 32 genes, which are all genes involved in the TCR signaling pathway as defined in the Reactome database [Matthews et al., 2009]. The final dataset used in the experiments below thus contains 113 samples, 57 and 56 for both context values respectively, and 32 variables.

To identify gene-gene interactions affected by the context, we performed a contextual analysis as described in Section 6.3 for each gene in turn, considering each time a particular gene as the target variable Y and all other genes as the set of input variables V. This procedure is similar to the procedure adopted in the Random forests-based gene network inference method called GENIE3 [Huynh-Thu et al., 2010], that was the best performer in the DREAM5 network inference challenge [Marbach et al., 2012]. Since gene expressions are numerical targets, we used variance as the impurity measure (see Section 6.3.5) and we built ensembles of 1000 totally randomized trees in all experiments.

The matrices in Figure 6.1 highlight context-dependent interactions found using different importance measures (detailed below). A cell (i, j) of these matrices corresponds to the importance of gene j when gene i is the output (the diagonal is irrelevant). White cells correspond to non significant context-dependencies as determined by random permutations of the context variable, using a significance level of 0.05. Significant context-dependent interactions in Figures 6.1(a) and (b) were de-

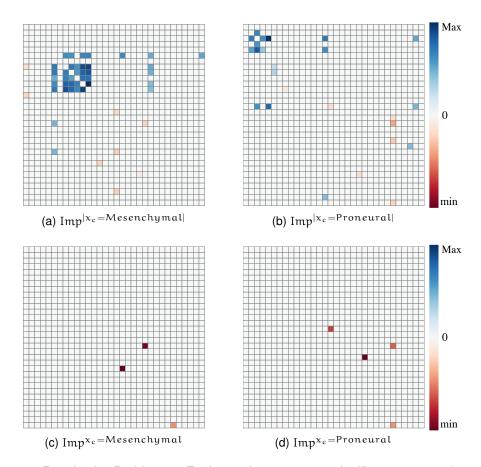


Figure 6.1: Results for Problem 4. Each matrix represents significant context-dependent gene-gene interactions as found using $Imp^{|x_c|}$ in (a)(b) and Imp^{x_c} in (c)(d), in GBM sub-type Mesenschymal in (a)(c) and Proneural in (b)(d). In (a) and (b), cells are colored according to $Imp_s^{x_c}$. In (c) and (d), cells are colored according to Impxc. Positive (resp negative) values are in blue (resp. red) and highlight context-redundant (resp. context-complementary) interactions. Higher absolute values are darker.

termined using the importance $Imp^{|x_c|}$ defined in (6.10), which is the measure we advocate in this paper. As a baseline for comparison, Figures 6.1(c) and (d) show significant interactions as found using the more straightforward score Imp^{x_c} defined in (6.10). In Figures 6.1(a) and (b) (resp. (c) and (d)), significant cells are colored according to the value of $Imp_s^{x_c}$ defined in (6.12). In Figures 6.1(c) and (d), they are colored according to the value of Imp^{x_c} in (6.10) instead. Blue (resp. red) cells correspond to positive (resp. negative) values of Imp^{x_c} or Imp^{x_c} and thus highlight context-redundant (resp. context-complementary) interactions. The darker the color, the higher the absolute value of Imp^{x_c} or Imp^{x_c} .

Respectively 49 and 26 context-dependent interactions are found in Figures 6.1(a) and (b). In comparison, only 3 and 4 interactions are found respectively in Figures 6.1(c) and (d) using the more straightforward score Imp x_c . Only 1 interaction is common between Figures 6.1(a) and (c), while 3 interactions are common between Figures 6.1(b) and (d). The much lower sensitivity of Imp^{x_c} with respect to $Imp^{|x_c|}$ was expected given the discussions in Section 6.3.2. Although more straightforward, the score $Imp^{x_c}(X_m)$, defined as the difference $Imp(X_m) - Imp(X_m|X_c = x_c)$, indeed suffers from the fact that $Imp(X_m)$ and $Imp(X_m|X_c=x_c)$ are estimated from different ensembles and thus do not explore the same conditionings in finite setting. Imp x_c also does not have the same guarantee as Imp $^{|x_c|}$ to find all contextdependent variables.

6.5 CONCLUSIONS AND FUTURE WORK

In this chapter, our first contribution is a formal framework defining and characterizing the dependence to a context variable of the relationship between the input variables and the output (Section 6.2). As a second contribution, we have proposed several novel adaptations of random forests-based variable importance scores that implement these definitions and characterizations and we have derived performance guarantees for these scores in asymptotic settings (Section 6.3). The relevance of these measures was illustrated on several artificial and real datasets (Section 6.4).

There remain several limitations to our framework that we would like to address as future works. All theoretical derivations in Sections 6.2 and 6.3 concern categorical input variables. It would be interesting to adapt our framework to continuous input variables, and also, probably with more difficulty, to continuous context variables. Finally, all theoretical derivations are based on forests of totally randomized trees (for which we have an asymptotic characterization). It would be interesting to also investigate non totally randomized tree algorithms (e.g., Breiman [2001]'s standard Random Forests method) that could provide better trade-offs in finite settings.

6.A DETAILS OF EXAMPLE 6.1

X_1	X ₂	X_{c}	Y
0		0	0
0 0	0 0	0 0	0
0	0	1	2
0 0 0 0	0	1	3
0	1	0	2
0	1	1 0 0	3
0	1	1	0
0	1	1	0
1 1 1	0	0	1
1	0 0	0 0 1	1
1	0	1	2
1	0	1	3
1	1	0	2
1	1	0	3
1 1 1	1	0 1 1	0 2 3 2 3 0 0 1 1 2 3 2 3 1 1
1	1	1	1

Table 6.A.1: Values of X_1 , X_2 , X_c and Y.

6.B PROOF OF THEOREM 6.1

Theorem. X_c is irrelevant to Y with respect to V iff all variables in V are context-independent to Y with respect to X_c (and V) and $I(Y; X_c) = 0$.

NECESSARY CONDITION.

Proof. If X_c is irrelevant to Y w.r.t. V, we have, by definition, that $I(Y; X_c|B) = 0$ for all subset $B \subseteq V$. Hence, we have $I(Y; X_c) = 0$ as a special case.

A variable $X_m \in V$ is context-independent if for all $B \subseteq V^{-m}$ and for all $x_c \in \mathcal{X}_c$, $b \in \mathcal{B}$, we have

$$I(Y; X_m|B = b, X_c = x_c) - I(Y; X_m|B = b) = 0.$$

Let us proof this:

$$\begin{split} I(Y; X_m | B = b, X_c = x_c) - I(Y; X_m | B = b) \\ &= H(Y | B = b, X_c = x_c) - H(Y | X_m, B = b, X_c = x_c) \\ &\hookrightarrow - H(Y | B = b) + H(Y | X_m, B = b) \end{split}$$

$$= H(Y|B = b) - H(Y|X_m, B = b)$$

$$\hookrightarrow -H(Y|B = b) + H(Y|X_m, B = b)$$

$$= 0,$$

where $H(Y|B=b,X_c=x_c)=H(Y|B=b)$ and $H(Y|X_m,B=b,X_c=x_c)=$ $H(Y|X_m, B = b)$ are consequences of $I(Y; X_c|B) = 0$ for all B if we assume that $p(B = b) \neq 0 \ (\forall b \in \mathcal{B}) \ \text{and} \ p(X_c = x_c, B = b) \neq 0 \ (\forall x_c \in \mathcal{X}_c \ \text{and} \ \forall b \in \mathcal{B}).$

SUFFICIENT CONDITION.

Proof. If all variables are context-independent, we have that for all $X_m \in V$, $B \subseteq P$ V^{-m} , $b \in \mathcal{B}$, and $x_c \in \mathcal{X}_c$:

$$I(Y; X_m | B = b, X_c = x_c) = I(Y; X_m | B = b).$$

By averaging the left- and right-hand sides of this equality over $P(B, X_c)$, we get:

$$I(Y; X_m | B, X_c) = I(Y; X_m | B).$$

From this, one can derive [Louppe et al., 2013]:

$$I(Y; X_c | B, X_m) = I(Y; X_c | B).$$

Since this equality is valid for all B, including $B = \emptyset$, and all X_m , we have that for all $B' \subseteq V$, $I(Y; X_c | B')$ can be reduced to $I(Y; X_c)$, which is equal to zero by hypothesis. The variable X_c is thus irrelevant to Y with respect to V.

6.C PROOF OF THEOREM 6.2

Theorem. A variable $X_m \in V$ is context-independent to Y with respect to X_c iff $\operatorname{Imp}^{|x_c|}(X_m) = 0$ for all x_c .

NECESSARY CONDITION.

Proof. By definition of context-independence, we have

$$I(Y; X_m | B = b, X_c = x_c) - I(Y; X_m | B = b) = 0$$

$$\forall B \subset V^{-m}, \forall x_c \in \mathcal{X}_c, \forall b \in \mathcal{B}.$$
(6.14)

Given that each term

$$|I(X_m; Y|B = b) - I(X_m; Y|B = b; X_c = x_c)|$$

of $Imp^{|x_c|}(X_m)$ (Equation (6.11)) is equal to 0, the sum is thus also equal to 0.

SUFFICIENT CONDITION.

Proof. Given the definition of $Imp^{|x_c|}(X_m)$:

$$Imp^{|x_c|}(X_m) = \sum_{k=0}^{p-1} \frac{1}{C_p^k} \frac{1}{p-k} \sum_{B \in \mathcal{P}_k(V^{-m})} \sum_{b \in \mathcal{B}} P(B = b)$$

$$\Leftrightarrow |I(X_m; Y|B = b) - I(X_m; Y|B = b; X_c = x_c)|,$$
(6.15)

appears to be a sum of positive terms (because of the absolute value). As in Theorem 6.1, we assume that probabilities are non-null and therefore, we have that the only way to have the sum equal to zero is to have each term of the sum equal to 0. Hence, we have $|I(X_m; Y|B = b) - I(X_m; Y|B = b; X_c = x_c)| = 0$ for all x_c , B and b which verifies the definition of context-independence for X_m .

6.D PROOF OF THEOREM 6.3

Theorem. If $|\operatorname{Imp}^{x_c}(X_m)| = \operatorname{Imp}^{|x_c|}(X_m)$ for a context-dependent variable X_m , then X_m is context-complementary if $Imp^{\kappa_c}(X_m) < 0$ and context-redundant if $Imp^{x_c}(X_m) > 0.$

Proof. The absolute value of a sum is less than or equal the sum of the absolute value of each terms. The equality is only verified when all terms are of the same sign. Therefore, the sign of $Imp^{x_c}(X_m)$ indicates the sign of all terms and thus verify either the context-complementarity if all terms are negative or the context-redundancy if all terms are positive.

6.E RESULTS FOR PROBLEM 3

		$Imp(X_m)$	$Imp(X_m X_c = x_c)$		$Imp^{ x_c }(X_m)$		$Imp_s^{x_c}(X_m)$	
m		-	$x_c = 0$	$x_c = 1$	$x_c = 0$	$x_c = 1$	$x_c = 0$	$x_c = 1$
0	age	0.2958	0.3386	0.2885	0.1382	0.1505	-0.0095	-0.0156
1	histologic-type	0.3522	0.1389	0.4366	0.2087	0.114	0.1988	-0.0569
2	degree-of-diffe	0.4413	0.4175	0.4208	0.1653	0.158	0.0561	0.0157
3	bone	0.2429	0.2502	0.2367	0.0933	0.0755	-0.0043	0.0165
4	bone-marrow	0.0192	0.0201	0.0148	0.0126	0.0101	0.0009	0.0041
5	lung	0.1627	0.2059	0.1370	0.1038	0.0949	-0.0259	0.0172
6	pleura	0.1485	0.1496	0.1015	0.0590	0.09	0.0313	0.0234
7	peritoneum	0.3184	0.3459	0.1979	0.0861	0.138	0.0147	0.0956
8	liver	0.2285	0.2138	0.2630	0.0786	0.1279	0.0375	-0.0602
9	brain	0.0465	0.0349	0.0548	0.0378	0.0254	0.0114	-0.0104
10	skin	0.0677	0.0362	0.0923	0.0314	0.0403	0.0252	-0.0133
11	neck	0.2215	0.0690	0.2582	0.1466	0.0692	0.1316	-0.0081
12	supraclavicular	0.1676	0.1915	0.1448	0.0845	0.067	-0.0198	0.0269
13	axillar	0.1393	0.1457	0.1068	0.0655	0.0629	-0.0067	0.0447
14	mediastinum	0.1838	0.2050	0.1716	0.1016	0.0806	-0.0059	0.0140
15	abdominal	0.2553	0.3296	0.1372	0.1346	0.1379	-0.0330	0.0898

Table 6.E.1: Importances as computed analytically using asymptotic formulas. The context is defined by the binary context feature Sex (Sex = 0 denotes female and Sex = 1 denotes *male*).

Overview

Dealing with datasets of very high dimension is a major challenge in machine learning. This chapter considers the problem of feature selection in applications where the memory is not large enough to contain all features. In this setting, we propose a novel tree-based feature selection approach that builds a sequence of randomised trees on small sub-samples of variables mixing both variables already identified as relevant by previous models and variables randomly selected among the other variables. As our main contribution, we provide an in-depth theoretical analysis of this method in infinite sample setting. In particular, we study its soundness with respect to common definitions of feature relevance and its convergence speed under various variable dependence scenarios. We also provide some preliminary empirical results highlighting the potential of this approach.

References: This chapter is an adapted version of the following publication:

A. Sutera, C. Châtel, G. Louppe, L. Wehenkel, and P. Geurts. Random subspace with trees for feature selection under memory constraints. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 929–937, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. URL http://proceedings.mlr.press/v84/sutera18a.html.

We do not reproduce Section 2 of this paper which provides background material already given in the preceding chapters of this manuscript.

7.1 MOTIVATION

We consider supervised learning and more specifically feature selection in applications where the memory is not large enough to contain all data. Such memory constraints can be due either to the large volume of available training data or to physical limits of the system on which training is performed (eg., mobile devices). A straightforward, but often efficient, way to handle such memory constraint is to build and average an ensemble of models, each trained on only a random subset of samples and/or features that can fit into memory. Such simple ensemble approaches have the advantage to be applicable to any batch learning algorithm, considered as a black-box, and they have been shown empirically to be very effective in terms of predictive performance, in particular when combined with trees, and even when samples and/or features are selected uniformly at random [see, eg., Chawla et al., 2004; Louppe and Geurts, 2012]. In particular, and independently of any considerations about memory constraints, feature subsampling has been shown in several works to be a very effective way to introduce randomization when building ensem-

bles of models [Ho, 1998; Kuncheva et al., 2010]. The idea of feature subsampling has also been investigated in the context of feature selection, where several authors have proposed to repeatedly apply a multivariate feature selection technique on random subsets of features and then to aggregate the results obtained on these subsets [see, eg., Dramiński et al., 2008; Lai et al., 2006; Konukoglu and Ganz, 2014; Nguyen et al., 2015; Dramiński et al., 2016].

In this chapter, focusing on feature subsampling, we adopt a simplistic setting where we assume that only q input features (among p in total, with typically $q \ll p$) can fit into memory. In this setting, we study ensembles of randomized decision trees trained each on a random subset of q features. In particular, we are interested in the properties of variable importance scores derived from these models and their exploitation to perform feature selection. In contrast to a purely uniform sampling of the features, we propose in Section 7.2 a modified sequential random subspace (SRS) approach that biases the random selection of the features at each iteration towards features already found relevant by previous models. As our main contribution, we perform in Section 7.3 an in-depth theoretical analysis of this method in infinite sample size condition. In particular, we show that (i) this algorithm provides some interesting asymptotic guarantees to find all (strongly) relevant variables, (ii) that accumulating previously found variables can reduce the number of trees needed to find relevant variables by several orders of magnitudes with respect to the standard random subspace method in some scenarios, and (iii) that these scenarios are relevant for a large class of (PC) distributions. As an important additional contribution, our analysis also sheds some new light on both the popular random subspace and random forests methods that are special cases of the SRS algorithm. Finally, Section 7.4 presents some preliminary empirical results with the approach on several artificial and real datasets.

SEQUENTIAL RANDOM SUBSPACE

In this chapter, we consider a simplistic memory-constrained setting where it is assumed that only q input features can fit into memory at once, with typically q small with respect to p. Under this hypothesis, Algorithm 7.1 describes the proposed sequential random subspace (SRS) algorithm to build an ensemble of randomized trees, which generalizes the Random Subspace (RS) method [Ho, 1998] (presented in Section 3.3.1). The idea of this method is to bias the random selection of the features at each iteration towards features that have already been found relevant by the previous trees. A parameter α is introduced that controls the degree of accumulation of previously identified features. When $\alpha = 0$, SRS reduces to the standard RS method. When $\alpha = 1$, all previously found features are kept while when $\alpha < 1$, some room in memory is left for randomly picked features, which ensures some permanent exploration of the feature space. Further randomization is introduced in the tree building step through the parameter $K \in [1, q]$, ie. the number of variables sampled at each tree node for splitting. Variable importance is assumed to be the MDI importance. This algorithm returns both an ensemble of trees and a subset F of variables, those that get an importance (significantly) greater than 0 in at least one tree of the ensemble. Importance scores for the variables can furthermore be derived from the final ensemble using Equation 4.2. In what follows, we will denote by $F_{q,T}^{K,\alpha}$ and $Imp_{q,T}^{K,\alpha}(X)$ resp. the set of features and the importance of feature Xobtained from an ensemble grown with SRS with parameters K, α , q and T.

Algorithm 7.1 Seguential Random Subspace algorithm

Inputs:

Data: Y the output and V, the set of all input variables (of size p).

Algorithm: q, the subspace size, and T the number of iterations, $\alpha \in [0,1]$, the percentage of memory devoted to previously found features.

Tree: K, the tree randomization parameter

Output: An ensemble of T trees and a subset F of features

Algorithm:

- 1. $F = \emptyset$
- 2. Repeat T times:
 - (a) Let $Q = R \cup C$, with R a subset of min{ $|\alpha q|, |F|$ } features randomly picked in F without replacement and C a subset of q - |R| features randomly selected in $V \setminus R$.
 - (b) Build a decision tree T from Q using randomization parameter K.
 - (c) Add to F all features from Q that get an importance greater than zero in \mathfrak{T} .

The modification of the RS algorithm is actually motivated by Propositions 5.1 and 5.2, stating that the relevance of high degree features can be determined only when they are analysed jointly with other relevant features of equal or lower degree. From this result, one can thus expect that accumulating previously found features will fasten the discovery of higher degree features on which they depend through some snowball effect. In the next section, we provide a theoretical asymptotic analysis of the SRS method that confirms and quantifies this effect.

Note that the SRS method can also be motivated from the perspective of accuracy. When $q \ll p$ and the number of relevant features r is also much smaller than the total number of features p (r \ll p), many trees with standard RS are grown from subsets of features that contain only very few, if any, relevant features and are thus expected not to be better than random guessing [Kuncheva et al., 2010]. In such setting, RS ensembles are thus expected not to be very accurate.

Example 7.1. With p = 10000, r = 10 and q = 50, the proportion of trees in a RS ensemble grown from only irrelevant variables is $C_{n-r}^q/C_{\mathfrak{p}}^q=0.95$.

With SRS (and $\alpha > 0$), we ensure that more and more relevant variables are given to the tree growing algorithm as iterations proceed and therefore we reduce the chance to include totally useless trees in the ensemble. Note however that in finite settings, there is a potential risk of overfitting when accumulating the variables. The parameter α thus controls a new bias-variance tradeoff and should be tuned appropriately. We will study the impact of SRS on accuracy empirically in Section 7.4.

7.3 THEORETICAL ANALYSIS

In this section, we carry out a theoretical analysis of the proposed method when seen as a feature selection technique. This analysis is performed in asymptotic sample size condition, assuming that all features, including the output, are discrete, and using Shannon entropy as the impurity measure. We proceed in two steps. First, we study the soundness of the algorithm, ie., its capacity to retrieve the relevant variables when the number of trees is infinite. Second, we study its convergence properties, ie. the number of trees needed to retrieve all relevant variables in different scenarios.

7.3.1 Soundness

Our goal in this section is to characterize the sets of features $F_{q,\infty}^{K,\alpha}$ that are identified by the SRS algorithm, depending on the value of its parameters q, α , and K, in an asymptotic setting, ie. assuming an infinite sample size and an infinite forest $(T = \infty)$. Note that in asymptotic setting, a variable is relevant as soon as its importance in one of the tree is strictly greater than zero and we thus have the following equivalence for all variables $X \in V$:

$$X \in F_{q,\infty}^{K,\alpha} \Leftrightarrow Imp_{q,\infty}^{K,\alpha}(X) > 0$$

Furthermore, in infinite sample size setting, irrelevant variables always get a zero importance and thus, whatever the parameters, we have the following property for all $X \in V$:

X irrelevant
$$\Rightarrow$$
 X \notin F_{q,\infty}^{K,\alpha} (and Imp_{q,\infty}^{K,\alpha}(X) = 0).

The method parameters thus only affect the number and nature of the relevant variables that can be found. Denoting by $r \leq p$ the number of relevant variables, we will analyse separately the case $r \leqslant q$ (all relevant variables can fit into memory) and the case r > q (all relevant variables can not fit into memory).

ALL RELEVANT VARIABLES CAN FIT INTO MEMORY $(r \leqslant q)$. consider the case of the RS method ($\alpha = 0$). In this case, Louppe et al. [2013] have shown the following asymptotic formula for the importances computed with totally randomized trees (K = 1):

$$Imp_{q,\infty}^{1,0}(X) = \sum_{k=0}^{q-1} \frac{1}{C_p^k} \sum_{B \in \mathcal{P}_k(V^{-m})} I(X;Y|B), \tag{7.1}$$

where $\mathcal{P}_k(V^{-m})$ is the set of subsets of $V^{-m} = V \setminus \{X_m\}$ of cardinality k. Given that all terms are positive, this sum will be strictly greater than zero if and only if there exists a subset $B \subseteq V$ of size at most q-1 such that $Y \not\perp \!\!\! \perp X \mid B \ (\Leftrightarrow I(X;Y\mid B)>0)$, or equivalently if deg(X) < q. When $\alpha = 0$, RS with K = 1 will thus find all and only the relevant variables of degree at most q-1. Given Proposition 5.1, the degree of a variable X can not be larger than r-1 and thus as soon as $r \leqslant q$, we have the guarantee that RS with K = 1 will find all and only the relevant variables. Actually, this result remains valid when $\alpha > 0$. Indeed, asymptotically, only relevant variables will be selected in the F subset by SRS and given that all relevant variables can fit into memory, cumulating them will not impact the ability of SRS to explore all conditioning subsets B composed of relevant variables. We thus have the following result:

Proposition 7.1. $\forall \alpha$, if $r \leq q$: $X \in \mathbb{F}_{q,\infty}^{1,\alpha}$ iff X is relevant.

In the case of non-totally randomized trees (K > 1), we lose the guarantee to find all relevant variables even when $r \leq q$. Indeed, there is potentially a masking effect due to K > 1 that might prevent the conditioning needed for a given variable to be relevant to appear in a tree branch. However, we have the following general result:

Theorem 7.2. $\forall \alpha, K$, if $r \leqslant q$: X strongly relevant $\Rightarrow X \in F_{q,\infty}^{K,\alpha}$

Proof. See Appendix 7.A

There is thus no masking effect possible for the strongly relevant features when K > 1 as soon as the number of relevant features is lower than q. For a given K, the features found by SRS will thus include all strongly relevant variables and some (when K > 1) or all (when K = 1) weakly relevant ones. It is easy to show that increasing K can only decrease the number of weakly relevant variables found. Using K = 1 will thus provide a solution for the **all-relevant** problem, while increasing K will provide a better and better approximation of the **minimal-optimal** problem in the case of strictly positive distributions (see Section 2.4.4 for definitions of those problems).

Interestingly, Theorem 7.2 remains true when q = p, ie., when forests are grown without any feature sampling. It thus extends Theorem 5.5 (also [Louppe et al., 2013, Theorem 3]) for arbitrary K in the case of standard random forests.

ALL RELEVANT VARIABLES CAN NOT FIT INTO MEMORY (r > q). When all relevant variables can not fit into memory, we do not have the guarantee anymore to explore all minimal conditionings required to find all (strongly or not) relevant variables, whatever the values of K and α . When $\alpha = 0$, we have the guarantee however to identify the relevant variables of degree strictly lower than q. When α > 1, some space in memory will be devoted to previously found variables that will introduce some further masking effect. We nevertheless have the following general results (without proof):

Proposition 7.3.

$$\forall X: \quad X \ \textit{relevant and} \\ deg(X) < (1-\alpha)q \Rightarrow X \in F_{q,\infty}^{1,\alpha}.$$

Proposition 7.4.

$$\forall K, X: \quad X \text{ strongly relevant and} \\ deg(X) < (1-\alpha)q \Rightarrow X \in F_{q,\infty}^{K,\alpha}.$$

In these propositions, $(1-\alpha)q$ is simply the amount of memory that always remains available for the exploration of variables not yet found relevant.

DISCUSSION. Results in this section show that SRS is a sound approach for feature selection as soon as either the memory is large enough to contain all relevant variables or the degree of the relevant variables is not too high. In this latter case, the approach will be able to detect all strongly relevant variables whatever its parameters (K and α) and the total number of features p. Of course, these parameters will have a potentially strong influence on the number of trees needed to reach convergence (see the next section) and the performance in finite setting.

7.3.2 Convergence

Results in the previous section show that accumulating relevant variables has no impact on the capacity at finding relevant variables asymptotically (when $r \leq q$). It has

however a potentially strong impact on the convergence speed of the algorithm, as measured for example by the expected number of trees needed to find all relevant variables. Indeed, when $\alpha = 0$ and $q \ll p$, the number of iterations/trees needed to find relevant variables of high degree can be huge as finding them requires to sample them together with all features in their conditioning. Given Proposition 2, we know that a minimum subset B such that $X \not\perp X \mid B$ for a relevant variable X contains only relevant variables. This suggests that accumulating previously found relevant features can improve significantly the convergence, as each time one relevant variable is found it increases the chance to find a relevant variable of higher degree that depends on it. In what follows, we will quantify the effect of accumulation on convergence speed in different best-case and worst-case scenarios and under some simplifications of the tree building procedure. We will conclude by a theorem highlighting the interest of the SRS method in the general class of PC distributions.

SCENARIOS AND ASSUMPTIONS. The convergence speed is in general very much dependent on the data distribution. We will study here the following three specific scenarios (where features $\{X_1, \dots, X_r\}$ are the only relevant features):

- Chaining: The only and minimal conditioning that makes variable X_i relevant is $\{X_1,\ldots,X_{i-1}\}$ (for $i=1,\ldots,r$). We thus have $\deg(X_i)=i-1$. This scenario should correspond to the most favorable situation for the SRS algorithm.
- Clique: The only and minimal conditioning that makes variable X_i relevant is $\{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_r\}$ (for $i = 1, \dots, r$). We thus have $deg(X_i) = r - 1$ for all i. This is a rather defavorable case for both RS and SRS since finding a relevant variable implies to draw all of them at the same iteration.
- Marginal-only: All variables are marginally relevant. We will furthermore make the assumption that these variables are all strongly relevant. They can not be masked mutually. This scenario is the most defavorable case for SRS (versus RS) since accumulating relevant variables is totally useless to find the other relevant variables and it should actually slow down the convergence as it will reduce the amount of memory left for exploration.

In Appendix 7.B.2, we provide explicit formulation of the expected number of iterations needed to find all r relevant features in the chaining and clique scenarios both when $\alpha = 0$ (RS) and $\alpha = 1$ (SRS). In Appendix 7.B.3, we provide order 1 Markov chains that model the evolution through the iterations of the number of variables found in the three scenarios when $\alpha = 0$ and $\alpha = 1$. These chains can be used to compute numerically the expected number of relevant variables found through the iterations (and in the case of the marginal-only setting, the expected number of iterations to find all variables). These derivations are obtained assuming $r \leqslant q$, K = q, and under the following additional simplifying assumptions.

Below, we compute analytically the average number of trees needed to find all relevant variables in the chaining and clique scenarios and we derive transition matrices of Markov chains that model the evolution of the number of variables found through the iterations in the three scenarios. These results are obtained assuming K = q and $r \leq q$, and with either $\alpha = 0$ (RS) or $\alpha = 1$ (SRS).

To make these derivations possible and independent of a particular data distribution, one needs furthermore to simplify the decision tree growing algorithm in the case of the chaining and clique scenarios. In what follows, trees are thus assumed to be grown such that a unique variable is selected at each tree level and this variable is selected at random among all variables X such that $Y \not\perp \!\!\! \perp X|B$ where B is the set of all variables tested at previous levels.

In the clique scenario, this assumption implies that only one variable of the clique will get a non-zero importance when all clique variables are selected at one iteration of RS/SRS (since only the last variable of the clique tested along a tree branch can get a non-zero score and this variable is the same in each branch given our tree growing assumption). This corresponds to a pessimistic scenario. Indeed, with standard unconstrained trees, several relevant variables could be found at one iteration given that the ordering of the variables, and thus the last variable of the clique tested, might differ from one tree branch to another. As a consequence, the tree growing assumption will lead to an overestimation of the number of trees needed to reach convergence. In the chaining scenario, the simplified tree growing algorithm implies that all relevant variables selected at one iteration of RS/SRS together with their minimal conditioning will get a non-zero importance. This corresponds this time to an optimistic scenario, as, with unconstrained trees, such variable might not be detected at one iteration depending on the exact data distribution. This will thus lead this time to an underestimation of the number of trees needed to reach convergence. Note however that, in both cases, these over/under-estimations will affect both RS and SRS in the same proportion and thus our assumption will not impact their relative performance.

Note that in the marginal-only scenario, given that all relevant variables are marginally and strongly relevant, they will always get a non-zero importance as soon as they are selected at one iteration. Our estimations below are thus not impacted by the simplification of the tree growing algorithm.

RESULTS AND DISCUSSION. Tables 7.3.2a, 7.3.2b, and 7.3.2c show the expected number of iterations needed to find all relevant variables for various configurations of the parameters p, q, and r, in the three scenarios. Figure 7.3.1 plots the expected number of variables found at each iteration both for RS and SRS in the three scenarios for some particular values of the parameters.

From these results, we can draw several conclusions. In all cases, expected times (ie., number of iterations/trees to find all relevant variables) depend mostly on the ratio $\frac{q}{p}$, not on absolute values of q and p. The larger this ratio, the faster the convergence. Parameter r has a strong impact on convergence speed in all three scenarios.

The most impressive improvements with SRS are obtained in the chaining hypothesis, where convergence is improved by several orders of magnitude (Table 7.3.2a and Figure 7.3.1a). At fixed p and q, the time needed by RS indeed grows exponentially with $r \ (\simeq (\frac{p}{q})^r \ \text{if} \ r \ll q)$, while time grows linearly with r for the SRS method ($\simeq r_q^p$ if $r \ll q$) (see Eq. (7.2) and (7.4) in Appendix 7.B.2).

In the case of cliques, both RS and SRS need many iterations to find all features from the clique (see Table 7.3.2b and Figure 7.3.1b). SRS goes faster than RS but the improvement is not as important as in the chaining scenario. This can be explained by the fact that SRS can only improve the speed when the first feature of the clique has been found. Since the number of iterations needed to find the r features from the clique for RS is close to r times the number of iterations needed to find one feature from the clique, SRS can only decrease at best the number of iterations by approximately a factor r (see Eq. (7.7) and (7.8) in Appendix 7.8.2).

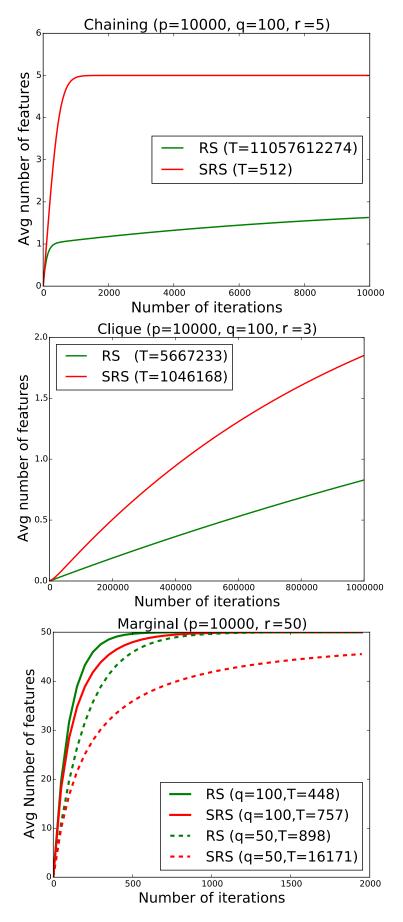


Figure 7.3.1: Evolution of the number of selected features in the different scenarios.

Table 7.3.1: Expected number of iterations needed to find all relevant variables for various configurations of parameters p, q and r with RS ($\alpha = 0$) and SRS ($\alpha = 1$) in the three scenarios.

Config (p,q,r)	RS	SRS	Config (p,q,r)	RS	SRS
10 ⁴ , 100, 1	100	100	10 ⁴ , 100, 1	100	100
10 ⁴ , 100, 2	10100	200	10 ⁴ , 100, 2	30300	10302
10 ⁴ , 100, 3	> 10 ⁶	301	10 ⁴ , 100, 3	5 · 10 ⁶	10 ⁶
10 ⁴ , 100, 5	> 10 ¹⁰	506	10 ⁴ , 100, 4	9 · 10 ⁸	108
10 ⁵ , 100, 3	> 109	3028	$10^4, 10^3, 4$	83785	11635
(a) Chaining.			(b)	Clique.	

Config (p,q,r) RS SRS $10^4, 100, 10$ 291 312 $10^4, 100, 50$ 448 757 104, 100, 90 2797 506 $10^4, 100, 100$ 1123 16187 25000, 100, 50 1123 1900

(c) Marginal-only.

In the marginal-only setting, SRS is actually slower than RS because the only effect of cumulating the variables is to leave less space in memory for exploration. The decrease of computing times is however contained when r is not too close to q(see Table 7.3.2c and Figure 7.3.1c).

Since we can obtain very significant improvement in the case of the chaining and clique scenarios and we only increase moderately the number of iterations in the marginal-only scenario (when r is not too close from q), we can reasonably expect improvement in general settings that mix these scenarios.

PC DISTRIBUTIONS AND CHAINING. Chaining is the most interesting scenario in terms of convergence improvement through variable accumulation. In this scenario, SRS makes it possible to find high degree relevant variables with a reasonable amount of trees, when finding these variables would be mostly untractable for RS. We provide below two theorems that show the practical relevance of this scenario in the specific case of PC distributions.

A PC distribution is defined as a strictly positive (P) distribution that satisfies the composition (C) property stated as follows Nilsson et al. [2007]:

Property 7.1. For any disjoint sets of variables R, S, T, U \subset V \cup {Y}:

$$S \perp\!\!\!\perp T | R \text{ and } S \perp\!\!\!\perp U | R \Rightarrow S \perp\!\!\!\perp T \cup U | R$$

The composition property prevents the occurrence of cliques and is preserved under marginalization. PC actually represents a rather large class of distributions that encompasses for example jointly Gaussian distributions and DAG-faithful distributions Nilsson et al. [2007].

The composition property allows to make Proposition 5.2 more stringent in the case of PC:

Proposition 7.5. Let B denote a minimal subset B such that Y \(\mathbb{L} \) X|B for a relevant variable X. If the distribution P over $V \cup \{Y\}$ is PC, then for all $X' \in B$, deg(X') < |B|.

Proof. Proposition 5.2 proves that the degree of all features in B is \leq |B| in the general case. Let us assume that there exists a feature $X' \in B$ of degree |B| in the case of PC distribution. Since this property remain true when the set of features V is reduced to a subset $V' = B \cup \{X\}$, the minimal B' of X' can only be $(B \setminus \{X_i\}) \cup \{X\}$. We thus have the following two properties:

$$Y \perp \!\!\! \perp X|B \setminus \{X'\}$$

$$Y \perp \!\!\! \perp X' | B' \setminus \{X\},$$

because B and B' are minimal. Together, by the composition property, they should imply that

$$Y \perp \{X, X_i\} | B \setminus \{X_i\},$$

which implies, by weak union: $Y \perp \!\!\! \perp X|B$, which contradicts the hypothesis.

In addition, one has the following result:

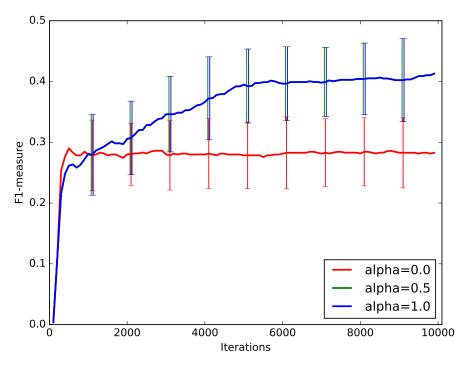
Theorem 7.6. For any PC distribution, let us assume that there exists a non empty minimal subset $B = \{X_1, \dots, X_k\} \subset V \setminus \{X\}$ of size k such that $X \not\perp\!\!\!\perp Y | B$ for a relevant variable X. Then, variables X_1 to X_k can be ordered into a sequence $\{X'_1,\ldots,X'_k\}$ such that $deg(X_i) < i$ for all i = 1, ..., k.

Proof. Let us denote by $\{X'_1, X'_2, \dots, X'_k\}$ the variables in B ordered according to their degree, ie., $\deg(X_i') \leq \deg(X_{i+1}')$, for i = 1, ..., k-1. Let us show that $\deg(X_i') < i$ for all $i=1,\ldots,k$. If this property is not true, then there exists at least one $X_i'\in B$ such that $\deg(X_i) \geqslant i$. Let us denote by I the largest i such that $\deg(X_i) \geqslant i$. Using a similar argument as in the proof of Proposition 7.5, there exists some minimal subset $B' \subseteq B \setminus \{X_1\}$ such that $Y \not\perp X_1 \mid B'$. Given that $\deg(X_1) \geqslant 1$, this subset B should contain 1 variables or more from B $\setminus \{X_1\}$. It thus contains at least one variable X_m with $l < m \le k$, and this variable is such that $deg(X_m) < m$. Given Proposition 7.5, if B' is minimal and contains $X_{\mathfrak{m}}$, then for a PC distribution, $deg(X_{\mathfrak{m}})$ should be strictly smaller than $|B'| \ge 1$, which contradicts the fact that X_m is after X_1 in the ordering and proves the theorem.

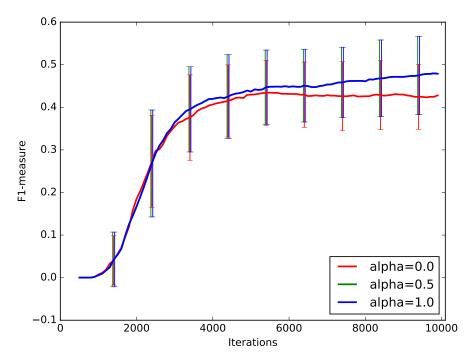
This theorem shows that, when the data distribution is PC, for all relevant variables of degree k, the k variables in its minimal conditioning form a chain of variables of increasing degrees (at worst). For such distribution, we thus have the guarantee that SRS finds all relevant variables with a number of iterations that grows almost only linearly with the maximum degree of relevant variables (see Eq.7.4 in Appendix 7.B.2), while RS would be unable to find relevant variables of even small degree.

7.4 EXPERIMENTS

Although our main contribution is the theoretical analysis in asymptotic setting of the previous section, we present here a few preliminary experiments in finite setting as a first illustration of the potential of the method. One of the main difficulties to implement the SRS algorithm as presented in Algorithm 7.1 is step 2(c) that decides which variable should be incorporated in F at each iteration. In infinite sample size setting, a variable with a non-zero importance in a single tree is guaranteed to be



(a) SRS with $q=0.05\times p$ on a dataset with p=50000 features and r=20 relevant features.



(b) SRS with $q=0.005 \times p$ on a dataset with p=50000 features and r=20 relevant features.

Figure 7.4.1: Evolution of the evaluation of the feature subset found by RS and SRS using the F1-measure computed with respect to relevant features. A higher value means that more relevant features have been found. This experiment was computed on an artificial dataset (similar to madelon) of 50000 features with 20 relevant features and for two sizes of memory.

truly relevant. Mutual informations estimated from finite samples however will always be greater than 0 even for irrelevant variables. One should thus replace step 2(c) by some statistical significance tests to avoid the accumulation of irrelevant variables that would jeopardize the convergence of the algorithm. In our experiments here, we use a random probe (ie., an artificially created irrelevant variable) to derive a statistical measure assessing the relevance of a variable Stoppiglia et al. [2003a]. Details about this test are given in Appendix 7.C.

Figure 7.4.1 evaluates the feature selection ability of SRS for three values of α (including $\alpha = 0$) and two memory sizes (250 and 2500) on an artificial dataset with 50000 features, among which only 20 are relevant (see Appendix 7.C for more details). The two plots show the evolution of the F1-score comparing the selected features (in F) with the truly relevant ones as a function of the number of iterations. As expected, SRS ($\alpha > 0$) is able to find better feature subsets than RS ($\alpha = 0$) for both memory sizes and both values of $\alpha > 0$.

Additional results are provided in Appendix 7.C that compare the accuracy of ensembles grown with SRS for different values of α and on 13 classification problems. These comparisons clearly show that accumulating the relevant variables is beneficial most of the time (eg., SRS with $\alpha = 0.5$ is significantly better than RS on 7 datasets, comparable on 5, and significantly worse on only 1). Interestingly, SRS ensembles with $\alpha = 0.5$ are also most of the time significantly better than ensembles of trees grown without memory constraint (see Appendix 7.C for more details).

7.5 CONCLUSIONS AND FUTURE WORK

Our main contribution is a theoretical analysis of the SRS (and RS) methods in infinite sample setting. This analysis showed that both methods provide some guarantees to identify all relevant (or all strongly relevant) variables as soon as the number of relevant variables or their degree is not too high with respect to the memory size. Compared to RS, SRS can reduce very strongly the number of iterations needed to find high degree variables in particular in the case of PC distributions. We believe that our results shed some new light on random subspace methods for feature selection in general as well as on tree-based methods, which should help designing better feature selection procedures.

Some preliminary experiments were provided that support the theoretical analysis, but more work is clearly needed to evaluate the approach empirically on controlled and real high-dimensional problems. We believe that the statistical test used to decide which feature to include in the relevant set should be improved with respect to our first implementation based on the introduction of a random probe. One drawback of the SRS method with respect to RS is that it can not be parallelized anymore because of its sequential nature. It would be interesting to design and study variants of the method that are allowed to grow parallel ensembles at each iteration instead of single trees. Finally, relaxing the main hypotheses of our theoretical analysis would be also of course of great interest.

7.A PROOF OF THEOREM 7.2

Theorem. $\forall \alpha, K, if r \leq q$: $X strongly relevant \Rightarrow X \in \mathbb{F}_{q,\infty}^{K,\alpha}$

Proof. By definition, X belonging to $F_{q,\infty}^{K,\alpha}$ means that there is at least one tree (grown with parameters q, α and K) in which X receives a strictly positive score for its split, i.e. such that Y depends on X conditionally to the variable assignement defined by the path from the root node to the node where X is used to split. Let us show that one such tree always exists whatever K and α when X is strongly relevant and $r \leq q$.

Within the infinite ensemble, let us consider only the trees grown using all r relevant variables (and q-r irrelevant ones randomly selected). Given that $r\leqslant q$ and given that only relevant features can be kept in memory, these trees are always explored whatever the value of $\alpha\geqslant 0$. Among these trees, let us furthermore only consider those such that irrelevant variables are tested in each branch only when all relevant variables (including X) are exhausted. These trees are always explored whatever the value of K. This derives from the fact that a relevant variable can always be picked with non zero probability at any tree node, except if all relevant variables have been tested above that node. Indeed, except in this latter case, the K tested variables can always include at least one relevant variable. If some relevant variable gets a non zero score, one relevant variable will be automatically used to split since irrelevant variables can only get zero scores. Even when all tested relevant variables get a zero score, one of them can still be selected instead of an irrelevant one given that tie are resolved by randomization.

Let us denote by τ_R the set of trees as just defined and let us show that X gets a non zero score in at least one tree in τ_R .

By definition 2 and property 1, X strongly relevant implies that there exists at least one assignement of values to all relevant variables but X such that conditionally to this assignement, Y is dependent on X. In each tree in τ_R , there is a path from the root node to a node where X is used to split that is compatible with this assignement. Let us assume that X always gets a zero score in all these compatible paths and show that this leads to a contradiction.

If all relevant variables are tested above X in a compatible path then X should receive a non zero score at its node, which would contradict our hypothesis. Thus, X can only be tested in a compatible path before all relevant variables have been tested. Given our hypothesis that X only gets zero scores, if X is used to split in one compatible path, then there exists another tree in τ_R with the same splits above X in the compatible path and with the split on X replaced by a split on another relevant variables (because of tie randomization or because of the randomization due to K < q). In this new tree, X is thus used to split at least one level below in the compatible path. Applying this argument recursively, one can thus show that there is at least one tree in τ_R where X is the last variable used to split in the compatible path. In this tree, X thus gets a non zero score, which contradicts the hypothesis and therefore concludes the theorem.

7.B **CONVERGENCE ANALYSIS**

Simplifying assumptions

Below, we compute analytically the average number of trees needed to find all relevant variables in the chaining and clique scenarios and we derive transition matrices of Markov chains that model the evolution of the number of variables found through the iterations in the three scenarios. These results are obtained assuming K = qand $r \leq q$, and with either $\alpha = 0$ (RS) or $\alpha = 1$ (SRS).

To make these derivations possible and independent of a particular data distribution, one needs furthermore to simplify the decision tree growing algorithm in the case of the chaining and clique scenarios. In what follows, trees are thus assumed to be grown such that a unique variable is selected at each tree level and this variable is selected at random among all variables X such that $Y \not\perp \!\!\! \perp X|B$ where B is the set of all variables tested at previous levels.

In the clique scenario, this assumption implies that only one variable of the clique will get a non-zero importance when all clique variables are selected at one iteration of RS/SRS (since only the last variable of the clique tested along a tree branch can get a non-zero score and this variable is the same in each branch given our tree growing assumption). This corresponds to a pessimistic scenario. Indeed, with standard unconstrained trees, several relevant variables could be found at one iteration given that the ordering of the variables, and thus the last variable of the clique tested, might differ from one tree branch to another. As a consequence, the tree growing assumption will lead to an overestimation of the number of trees needed to reach convergence. In the chaining scenario, the simplified tree growing algorithm implies that all relevant variables selected at one iteration of RS/SRS together with their minimal conditioning will get a non-zero importance. This corresponds this time to an optimistic scenario, as, with unconstrained trees, such variable might not be detected at one iteration depending on the exact data distribution. This will thus lead this time to an underestimation of the number of trees needed to reach convergence. Note however that, in both cases, these over/under-estimations will affect both RS and SRS in the same proportion and thus our assumption will not impact their relative performance.

Note that in the marginal-only scenario, given that all relevant variables are marginally and strongly relevant, they will always get a non-zero importance as soon as they are selected at one iteration. Our estimations below are thus not impacted by the simplification of the tree growing algorithm.

7.B.2 Average times

Let us denote by $T_{\text{chain}}^{RS}(i,p,q)$ (1 \leqslant i \leqslant r) the average number of iterations needed to find the feature X_i of degree i-1 and by $T_{chain}^{SRS}(i,p,q)$ the average number of iterations needed to find the same feature with the SRS algorithm (that forces the selection of already found relevant variables). Given our assumptions above, each tree will be able to identify all relevant variables X it gets as soon as it gets also the relevant variables in its minimal conditioning. Note that $T_{\text{chain}}^{RS/SRS}(i,p,q)$ can also be interpreted as the average time needed to find the first i relevant features, given that one can not find X_i without finding all features X_i with

 $1\leqslant j < i.\ T_{chain}^{RS/SRS}(r,p,q)$ also represents the average number of iterations needed to find all relevant variables under the chain assumption.

Theorem 7.7. Under our assumptions, the T_{chain}^{RS} function can be computed as follows:

$$T_{chain}^{RS}(i, p, q) = \prod_{l=0}^{i-1} \frac{p-l}{q-l}$$
 (7.2)

 $\textit{Proof.} \ \ \text{Indeed,} \ \ T^{RS}_{\text{chain}}(i,p,q) \ \ \text{is the mean of a geometric distributed random variance}$ able with a probability of success defined as the probability of drawing the i-1variables in X_i 's conditioning and X_i at the same time, which is given by:

$$\frac{\binom{p-i}{q-i}}{\binom{p}{q}} = \prod_{l=0}^{i-1} \frac{q-l}{p-l}.$$
 (7.3)

Theorem 7.8. Under the same assumption, $T_{chain}^{SRS}(i,p,q)$ can be computed as follows:

$$T_{chain}^{SRS}(i, p, q) = \sum_{l=0}^{i-1} \frac{p-l}{q-l} - (i-1)$$
 (7.4)

Proof. Let us show this by induction on i. The base case corresponds to i = 1. In this case, we have:

$$T_{chain}^{SRS}(1,p,q) = T_{chain}^{RS}(1,p,q) = \frac{p}{q},$$

which satisfies Eqn (7.4). Let us assume that Eqn. (7.4) is satisfied for i < i' and let us show that it is satisfied for i = i'. $T_{chain}^{SRS}(i',p,q)$ can be defined as follows:

$$\begin{split} T_{\text{chain}}^{SRS}(i',p,q) &= \frac{q}{p} T_{\text{chain}}^{SRS}(i'-1,p-1,q-1) + \\ &\qquad \qquad (1-\frac{q}{p})(1+T_{\text{chain}}^{SRS}(i',p,q)). \end{split} \tag{7.5}$$

One can indeed distinguish two cases:

- X_1 is selected at the first iteration (this happens with probability q/p): the average time needed to find feature $X_{i'}$ of degree i'-1 then becomes the time needed to find a feature of degree i'-2 when one is allowed to draw q-1features among p-1, which is $T_{chain}^{SRS}(i'-1, p-1, q-1)$
- X_1 is not selected at the first iteration (this happens with probability 1 q/p): in this case, the first iteration is useless and thus the number of iterations needed will be $1 + T_{chain}^{SRS}(i', p, q)$.

Eqn. (7.5) can be used to compute T^{SRS}_{chain} recursively:

$$T_{\text{chain}}^{SRS}(i',p,q) = T_{\text{chain}}^{SRS}(i'-1,p-1,q-1) + (\frac{p}{q}-1). \tag{7.6} \label{eq:7.6}$$

Deriving Eqn. (7.4) from Eqn. (7.6) is then straightforward, which concludes the proof by induction.

Eqn. (7.4) shows that the average time needed to find the i first features is equal to the sum of the time needed to find all features individually minus the number of features. This last term takes into account the fact that by chance, one might find several features at once.

CLIQUE. Let us denote by $\mathsf{T}^{RS}_{c1}(i,p,q)$ and $\mathsf{T}^{SRS}_{c1}(i,p,q)$, the average time needed to find i features (among r) from the clique respectively with the RS and the SRS algorithm. Given our assumptions above, when the tree growing algorithm is given all r relevant features, it will be able to identify one (and only one) feature from the clique at random. If it has already found i features from the clique, the chance to get a new one, when all r features are selected among the q ones, will thus be (r-i)/r, i.e., the probability to test one of the r-i not yet found features after all other r features from the clique.

Theorem 7.9.

$$\mathsf{T}_{\mathrm{cl}}^{\mathrm{RS}}(\mathfrak{i},\mathfrak{p},\mathfrak{q}) = \left(\prod_{l=0}^{r-1} \frac{\mathfrak{p}-l}{\mathfrak{q}-l}\right) \cdot \left(\sum_{l=0}^{\mathfrak{i}-1} \frac{r}{r-l}\right) \tag{7.7}$$

Proof. The first factor in Eqn.(7.7) is the inverse of the probability of selecting all r relevant features at once. Each term of the sum in the second factor corresponds to the inverse of the probability of testing a new relevant variables, not yet found, at the bottom of the tree. As discussed above, this probability is $\frac{r-1}{r}$ when we have already found 1 features from the clique.

Theorem 7.10.

$$T_{c1}^{SRS}(i,p,q) = \sum_{l=0}^{i-1} \frac{r}{r-l} \prod_{m=l}^{r-1} \frac{p-m}{q-m}$$
 (7.8)

Proof. Each term of the sum represents the average time needed to find a new clique feature given that we have already found I features. This time is equal to one over the probability of finding a new feature when we have already found l of them. This latter is the probability of selecting among q the r-l missing relevant features (i.e., $\prod_{m=1}^r \frac{q-m}{p-m})$ times the probability of testing one of the missing relevant features at the bottom of the tree (i.e., (r-1)/r).

When $i=1,\ T_{c1}^{SRS}(1,p,q)=T_{c1}^{RS}(1,p,q).$ Intuitively, it indeed takes the same time for the RS and the SRS algorithms to find the first relevant features. When i increases however, the SRS algorithm becomes faster and faster than the RS algorithm. Indeed, the RS algorithm always needs to find all r clique features, while the SRS one only needs to find the r-i missing relevant features.

Markov chain interpretation

Let us denote by $N_{+}^{X,Y}$ the number of variables found for t iterations, with $X=c,\,X=$ g, and X = m respectively for the chain hypothesis, the clique hypothesis and the marginal only hypothesis (as defined in the first section of this document) and Y = nand Y = s respectively for the RS and SRS algorithms. All these random variables follow order 1 Markov chains. The transition probabilities are provided below for each chain (without proof), under the assumptions given in Section 7.B.1.

CHAIN HYPOTHESIS.

$$P(N_{t}^{c,n} = l_{1} | N_{t}^{c,n} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ \frac{\binom{p-r}{q-l_{1}}}{\binom{p}{q}} & \text{if } l_{1} > l_{2} \\ 1 - \sum_{i=l_{2}+1}^{r} \frac{\binom{p-r}{q-i}}{\binom{p}{q}} & \text{if } l_{1} = l_{2} \end{cases}$$
 (7.9)

$$P(N_{t}^{c,s} = l_{1}|N_{t}^{c,s} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ \frac{\binom{p-r}{q-l_{1}}}{\binom{p-l_{2}}{q-l_{2}}} & \text{if } l_{1} > l_{2} \\ 1 - \sum_{i=l_{2}+1}^{r} \frac{\binom{p-r}{q-l_{2}}}{\binom{p-l_{2}}{q-l_{2}}} & \text{if } l_{1} = l_{2} \end{cases}$$

$$(7.10)$$

CLIQUE HYPOTHESIS.

$$P(N_{t}^{g,n} = l_{1} | N_{t}^{g,n} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ 1 - \frac{\binom{p-r}{q-r}}{\binom{p}{q}} \frac{r - l_{2}}{r} & \text{if } l_{1} = l_{2} \\ \frac{\binom{p-r}{q-r}}{\binom{p}{q}} \frac{r - l_{2}}{r} & \text{if } l_{1} = l_{2} + 1 \\ 0 & \text{if } l_{1} > l_{2} + 1 \end{cases}$$

$$(7.11)$$

$$P(N_{t}^{g,s} = l_{1}|N_{t}^{g,s} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ 1 - \frac{\binom{p-r}{q-l_{2}}}{\binom{p-l_{2}}{q-l_{2}}} \frac{r-l_{2}}{r} & \text{if } l_{1} = l_{2} \\ \frac{\binom{p-r}{q-l}}{\binom{p-l_{2}}{q-l_{2}}} \frac{r-l_{2}}{r} & \text{if } l_{1} = l_{2} + 1 \\ 0 & \text{if } l_{1} > l_{2} + 1 \end{cases}$$

$$(7.12)$$

MARGINAL ONLY HYPOTHESIS.

$$P(N_{t}^{m,n} = l_{1}|N_{t}^{m,n} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ \frac{\binom{r-l_{2}}{l_{1}-l_{2}}\binom{p-r+l_{2}}{q-l_{1}+l_{2}}}{\binom{p}{q}} & \text{if } l_{1} > l_{2} \\ \frac{\binom{p-r+l_{2}}{q}}{\binom{p}{q}} & \text{if } l_{1} = l_{2} \end{cases}$$
(7.13)

$$P(N_{t}^{m,s} = l_{1} | N_{t}^{m,s} = l_{2}) = \begin{cases} 0 & \text{if } l_{1} < l_{2} \\ \frac{\binom{r-l_{2}}{l_{1}-l_{2}}\binom{p-r}{q-l_{1}}}{\binom{p-l_{2}}{q-l_{2}}} & \text{if } l_{1} > l_{2} \\ \frac{\binom{p-r}{q-l_{2}}}{\binom{p-l_{2}}{q-l_{2}}} & \text{if } l_{1} = l_{2} \end{cases}$$

$$(7.14)$$

DETAILS FOR SECTION 7.4 7.c

In this section, we give more details about our practical implementation of SRS and performed experiments.

7.c.1 On the use of a random probe to distinguish relevant features from irrelevant features.

As explained in Section 7.4, we add an artificial irrelevant feature in data as a random probe. By comparison with that probe of importances scores, one can distinguish relevant features (better than the probe) from irrelevant features. Through iterations, we can compute a p-value score which is the percentage of times a variable

Dataset	# samples	# features
arcene	100	10000
breast2	295	24496
cina0	16033	132
isolet	7797	617
madelon	2000	500
marti0	500	1024
reged0	500	999
secom	1567	591
mnist	70000	784
mnist3v8	13966	784
mnist4v9	13782	784
sido0	12678	4932
tis	13375	927

Table 7.C.1: Dataset specifications

has been better than the probe. If the p-value is above a given threshold β then the feature is likely relevant. Moreover, a variable has to be sampled more than L times in Q sets to insure that the p-value is reliable. Then at each iteration, the variables that satisfy the two criteria are added to F. In the following experiments, we choose arbitrarily L = 10 and $\beta = 95\%$.

7.c.2 On the datasets and on the protocol

We evaluate the accuracy of all these methods on a list of both artificial and real classifications problems (all but madelon are real data) described in Table 7.C.1 and publicly available in the UCI machine learning repository Lichman [2013]. For each dataset, we separate it into two random partitions of the same size (i.e., the same number of samples) to have a training set and a test set. There is no optimization of the parameters. For all datasets, the procedure was repeated 50 times, using the same random partitions between all methods. Following results are averages over those 50 runs.

7.c.3 Detailed results

Table 7.C.2 is average accuracy scores obtained on all datasets for each method for some parameters. We consider different sizes of memory (i.e., parameter q) and different value for the parameter α for the SRS algorithm. This allows to consider every behaviour of the SRS algorithm : without memory ($\alpha = 0$) which is equivalent to the Random Subspace method, with a full memory ($\alpha = 1$) and a non-full memory ($\alpha = 0.5$). For both methods (RS and SRS), a single extra-tree is build at each iteration. The randomization parameter of the extra-tree is set to its maximal value (ie., all features). For the tree-based ensemble methods, we consider different values for the randomization parameter. This parameter reduces the ability to consider the whole dataset in once and in that it relates in a way to the size of the memory of SRS. We choose for that parameter values of 0.01, 0.1 and 1 corresponding to considering respectively 1%, 10%, 100% of all features at each node.

		SRS									Tree	-based ens	emble me	thods	
		q=0.01			q=0.05			q=0.1		RF ET					
					α						Ra	ndomizatio	n paramet	er K	
	0.0	0.5	1.0	0.0	0.5	1.0	0.0	0.5	1.0	0.01	0.1	1	0.01	0.1	1
arcene	0.743	0.717	0.717	0.743	0.743	0.743	0.732	0.732	0.732	0.717	0.706	0.678	0.739	0.729	0.701
breast2	0.649	0.647	0.647	0.651	0.651	0.650	0.654	0.654	0.654	0.646	0.649	0.649	0.650	0.654	0.651
cina0	0.755	0.755	0.777	0.809	0.929	0.873	0.931	0.933	0.921	0.933	0.939	0.939	0.931	0.934	0.934
isolet	0.906	0.899	0.336	0.944	0.945	0.766	0.949	0.950	0.817	0.936	0.940	0.912	0.943	0.951	0.943
madelon	0.558	0.689	0.745	0.639	0.858	0.861	0.673	0.845	0.845	0.620	0.700	0.754	0.608	0.690	0.815
marti0	0.881	0.881	0.881	0.874	0.874	0.874	0.870	0.870	0.870	0.878	0.870	0.866	0.879	0.868	0.854
reged0	0.880	0.966	0.939	0.885	0.974	0.974	0.898	0.974	0.974	0.882	0.963	0.960	0.881	0.948	0.978
secom	0.935	0.935	0.930	0.935	0.931	0.931	0.934	0.932	0.932	0.935	0.933	0.929	0.935	0.930	0.928
mnist	0.564	0.823	0.525	0.959	0.966	0.905	0.968	0.970	0.938	0.964	0.966	0.953	0.966	0.971	0.968
mnist3v8	0.910	0.941	0.828	0.980	0.986	0.958	0.987	0.989	0.975	0.980	0.985	0.978	0.981	0.988	0.987
mnist4v9	0.889	0.957	0.848	0.981	0.986	0.960	0.986	0.988	0.974	0.983	0.984	0.974*	0.985	0.987	0.984*
sido0	0.970	0.972	0.953	0.973	0.968	0.968	0.974	0.969	0.969	0.972	0.973	0.973*	0.973	0.974	0.960*
tis	0.751	0.751	0.757	0.753	0.887	0.888	0.844	0.917	0.915	0.854	0.916	0.913*	0.856	0.906	0.914*

Table 7.C.2: Average accuracy scores for all methods with specified parameters on original datasets. SRS and RS were computed with 10000 iterations and RF/ET with 10000 trees.

	RS	SRS		E	T
	q = 0.1	q = 0.1	q = 0.1	k = 0.1	k = 1.0
q = 0.1		$\alpha = 0.5$	$\alpha = 1.0$		
RS	_	1/5/7	6/4/3	7/2/4	6/3/4
${\tt SRS}_{\alpha=0.5}$	7/5/1	_	5/8/0	9/2/2	10/2/1
$SRS_{\alpha=1.0}$	3/4/6	0/8/5	_	5/2/6	6/2/5

	RS	SRS	
	q = 0.01	q = 0.01	q = 0.01
q = 0.01		$\alpha = 0.5$	$\alpha = 1.0$
RS	_	2/2/9	5/2/6
$SRS_{\alpha=0.5}$	9/2/2	_	7/3/3
$SRS_{\alpha=1.0}$	6/2/5	3/3/7	_

(a) $q = 0.1 \times p$ (b) $q = 0.01 \times p$

Table 7.C.3: Pairwise t-test (with a significance level of 0.05) comparisons: each element on line i and column j of the table in terms of Win/Draw/Loss is the result of the comparison for method i vs. method j: the tree values indicate respectively on how many datasets method $\mathfrak i$ is significantly better / not significantly different / significantly worse than method j. All methods were computed with 10000 iterations or trees on all 14 datasets (from Table 7.C.1) with parameters specified on columns. In **bold** when the first value is greater than other values.

Overview

This chapter considers a specific machine learning task consisting in reconstructing a network from data. We first present the principle of GENIE3, originally designed to infer gene regulatory networks from samples of gene expression levels. Then we propose a simple yet effective solution to the problem of connectome inference in calcium imaging data. The proposed algorithm consists of two steps. First, processing the raw signals to detect neural peak activities. Second, inferring the degree of association between neurons from partial correlation statistics. Section 8.3 summarises the methodology that led us to win the Connectomics Challenge, proposes a simplified version of our method, and finally compares our result with respect to other inference methods.

References: Section 8.3 reproduces the following publication:

A. Sutera, A. Joly, V. François-Lavet, A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In *Neural Connectomics Workshop*, pages 23–35, 2015.

These results have also been published afterwards as a book chapter: A. Sutera, A. Joly, V. François-Lavet, Z. A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In *Neural Connectomics Challenge*, pages 23–36. Springer, 2017.

Note that Section 8.3.1 was not in the original publication and aims at putting the proposed method in perspective with tree-based network inference techniques.

8.1 MOTIVATION

In systems biology, networks provide a natural representation for complex feature interactions (where features are biological entities such as genes, proteins, ...) [Schrynemackers et al., 2013]. Network inference consists in the reconstruction of such biological networks from high-throughput data [De Smet and Marchal, 2010]. Concretely, given a set of p input variables $V = \{X_1, \ldots, X_p\}$, it aims at inferring (or completing¹) a directed graph with p nodes, where each node represents a variable, and an edge directed from one variable X_i to another variable X_j indicates a direct (causal) influence of X_i on X_j [Huynh-Thu et al., 2010; Louppe, 2014]. Sometimes, targeted networks are undirected and only represent interactions (i.e., conditional dependencies) between variables without any causal interpretation of the edge di-

¹Some network inference techniques use a priori knowledge including known interactions.

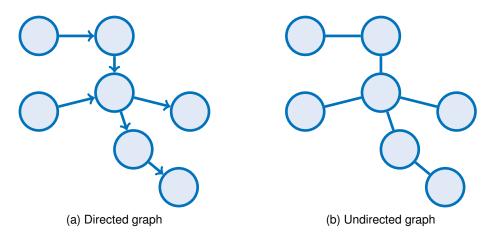


Figure 8.1.1: Examples of inferred networks.

rection [Louppe, 2014]. Figure 8.1.1 illustrates two possible networks: Figure 8.1.1a is a network of (causal) influences represented through a directed graph while Figure 8.1.1b is a network of statistical dependencies represented by an undirected graph.

Biological network inference consists in reconstructing a network in which biological entities interact (e.g., genes, proteins, cells, neurons, ...) [Tieri et al., 2016] and two applications in particular will be of interest in the rest of this chapter.

In genomics, gene regulatory networks represent interactions between genes and transcription factors² [Huynh-Thu, 2012; Louppe, 2014; Tieri et al., 2016]. The inference of such a network is based on gene expression levels.

In neuroscience, the connectome represents the neural connectivity, i.e., the interaction between neurons [de Abril et al., 2018; Panagopoulos, 2018]. Inferring the connectome from neural activity gives insights on effective brain structure. The effective brain structure gathers the structural (anatomical) connectivity (referring to physical connections between neurons, i.e., synapses) and the functional connectivity (referring to patterns of neuron activation regardless of spatiality, that are specific to a brain function and may change over time) and represent directional effects of neural elements on others [Sporns, 2007]. The activation of a neuron (i.e., an action potential) is characterised by a sudden change in membrane potential by opening Ca channels for instance [Simons, 1988; Tian et al., 2009]. Calcium imaging can thus be used to record the neuronal activity by means of fluorescent marker. Calcium fluorescent levels are converted into neural activation times series that are in turn used for connectome inference [Panagopoulos, 2018]

The interest in (biological) network inference has lead to many studies in the literature³ involving models based on statistical measures (e.g., mutual information and (cross- and partial-)correlation) and probabilistic models (e.g., Bayesian networks and Gaussian graphical models) for example.

²Transcription factors are proteins that regulates gene expression [Huynh-Thu, 2012].

³Exhaustive lists of methods are given in [de Abril et al., 2018; Panagopoulos, 2018] for connectome inference and in [Huynh-Thu et al., 2010; Marbach et al., 2012] for gene regulatory inference.

8.2 TREE-BASED NETWORK INFERENCE BASED ON VARIABLE IMPORTANCES

Tree-based models have been also developed for network inference because they advantageously do not make any assumption about the target function, deal with non-linearity and take into account feature dependencies [Huynh-Thu et al., 2010; Schrynemackers et al., 2015]. Both supervised and unsupervised approaches have been proposed for network inference and aim at deriving a score expressing the confidence for a pair of nodes to interact. In (tree-based) supervised approaches, a (tree-based) supervised model is usually constructed using a partial knowledge of the network and then used to assess the remaining untested pairs [Schrynemackers et al., 2013]. In tree-based unsupervised methods, variable importances are derived and used to estimate the degree of association between two variables [Huynh-Thu et al., 2010; Louppe, 2014]. We focus here on these latter methods.

8.2.1 **GENIE3**

GENIE3 [Huynh-Thu et al., 2010] is an approach that aims at inferring a network of p nodes by decomposing it into p independent supervised learning problems.

Given a set of variables $V = \{X_1, \dots, X_p\}$, the i^{th} sub-problem consists in learning a tree-based ensemble method (e.g., Random Forests or Extra-Trees) in order to predict the value of the variable X_i from all remaining p-1 variables X_j (with $j \neq i$). The contribution of X_i in the prediction of X_i gives an indication of the confidence level $p_{i,i}$ for the putative edge from X_i to X_i in the network (i.e., the degree of association between node j and node i). Aggregating the confidence levels of all pairs of nodes allows to reconstruct the whole network by selecting the top-ranked interactions (i.e., above a given threshold of confidence level) for example.

In the case of tree-based ensemble models, confidence scores are given by the variable importance scores⁴. However, the aggregation of importance scores resulting from the p sub-problems should be done cautiously if the variables are (i) of different scale, (ii) of different variability, or (iii) vary in the number of categories.

Indeed, Huynh-Thu et al. [2010] and Louppe [2014] point out⁵ a positive bias in the upper bound of (the sum of) all variable importances which depends on the target variable. In other words, if variables differ from each other on (i), (ii) or (iii), importance scores are not directly comparable without an appropriate normalisation⁶ before their aggregation.

8.2.2 Direct interaction

Network inference only considers direct connections between variables. Indirect effect must therefore be filtered out. However, neither relevance/usefulness nor importance score rankings may help to discriminate direct effects from indirect ones. Indeed, in all generality, variable importances do not guarantee that the importance of a feature indirectly related to the target has a lower importance score than any

 $^{^4}p_{j,i}$ is given by the importance of X_j in the sub-problem in which X_i is the target variable.

⁵It can also be retrieved in Chapter 4, especially from Equation 4.11.

⁶Let us consider a model learnt on a learning set LS using Shannon entropy (respectively, variance/gini index) as impurity measure and MDI importance scores. One may normalise the target variable by its entropy (respectively, variance) estimated on LS so that all variables have unit-entropy (resp., unit-variance) making importance scores comparable to each others. With MDA, one should consider a normalised accuracy metric.

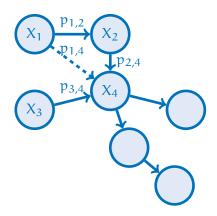


Figure 8.2.1: Direct interaction may be outscored by indirect ones. Solid arrows represent direct interactions while dashed arrows represent indirect effects. p_{1,3} may be numerically higher than $p_{3,4}$ while being associated to an indirect effect.

other feature directly related to the target. Moreover, one can imagine that several paths actually connect an input feature to the target (e.g., $X_1 \rightarrow X_2 \rightarrow X_4$ and $X_1 \rightarrow X_4$ in Figure 8.2.1) and so importance scores may reflect simultaneously direct and indirect interactions.

Regardless of the degree of association between nodes, one should aim at discovering the Markov boundary of a node in order to identify all its direct neighbours [Aliferis et al., 2010]. With strictly positive distributions (as seen in Section 2.4.2), it corresponds to identify all strongly relevant variables.

By definition, a strongly relevant feature X_m is such that $I(X_m; Y|V^{-m}) > 0$. Therefore, importance measures can, in theory, be adjusted to only detect strongly relevant features by only considering the deepest level of fully developed trees (corresponding to $B = V^{-m}$):

$$Imp_{\infty}^{mdi,last}(X_m) = \sum_{\nu^{-m}} P(V^{-m} = \nu^{-m}) I(X_m; Y | V^{-m} = \nu^{-m})$$
 (8.1)

where the sum of v^{-m} is a sum of over all possible value configurations of the set of variables V^{-m} . A feature such that $Imp_{\infty}^{mdi,last}(X_m) > 0$ is strongly relevant. However, the deepest level of fully developed trees is also the one where impurity decrease is estimated with the less samples and thus not reliable under other circumstances than infinite sample size.

In practice, a good heuristic to filter out as much as much as possible indirect interactions (and thus mainly focusing on strongly relevant variables) is to accentuate the masking effect (as strongly relevant variables can not be masked) by setting K > 1 (ideally, K = p) [Louppe, 2014]. Additionally, an adequate stopping criterion may help to mitigate impurity miss-estimation effects (by avoiding estimation on too few samples) [Louppe, 2014].

8.2.3 Edge orientation

In GENIE3, there is no explicit edge orientation despite that the importance score is usually asymmetrical $(p_{i,j} \neq p_{j,i})$ in opposition to symmetrical measures such as correlation or mutual information. In the method, only edges with confident scores above a given threshold are considered. Only one confidence score can be above

⁷In their experiments, Huynh-Thu et al. [2010] set the threshold such that the number of inferred edges corresponds to the number of edges in the gold standard.

the threshold implying a seemingly edge orientation while, on the contrary, both $p_{i,i}$ and $p_{i,i}$ can be kept making the edge undirected. Huynh-Thu et al. [2010] further analyse the ability of GENIE3 to correctly deduce the edge orientation including a comparison between $p_{i,j}$ and $p_{j,i}$ (i.e., $p_{i,j} > p_{j,i}$ implying the edge $i \to j$, or vice versa). Despite relatively symmetrical inferred networks (i.e., only few edges are only directed), GENIE3 seems to infer fairly correctly the edge orientation at least considering $p_{i,j} > p_{j,i}$ when an edge is such that $i \to j$. More recently, Bloebaum et al. [2018] investigate the asymmetry in the mean-squared errors of predicting the cause from the effect and the effect from the cause in order to determine the causal direction between two variables. Such researches are promising to infer edge orientation from observational data.

CONNECTOMICS CHALLENGE

8.3.1 Preamble

In the previous section, we introduced network inference and GENIE3, a tree-based method to infer a gene regulatory network. We also presented the questions of the direct effect identification and edge orientation.

Based on variable importances, GENIE3 provides excellent results in the context of gene regulatory network inference (best performer in the DREAM4 In Silico Multifactorial challenge in 2009 and in the DREAM5 Network Inference challenge in 2010). We however noticed that variable importances as usually used do not filter out indirect effects. Edge orientation seems promising but GENIE3-inferred networks are relatively symmetric and only few edges are undoubtedly oriented.

This section summarises contributions made in the scope of the competition "Neural Connectomics Challenge" organised in the context of 2014 ECML/PKDD conference [Battaglia et al., 2017], consisting in inferring a connectome from fluorescent calcium data. In what follows, we present our solution, which was the winning solution of the challenge.

In the context of connectome inference, neural networks seem to consist of fewer edges than gene regulatory networks (proportionally to the number of nodes). Subsequently, the number of indirect effects should be higher and thus it is even more crucial to identify direct interactions (actual edges). The edge orientation is however comparable with gene regulatory network.

The GENIE3 approach suggests to decompose the inference of a network of p nodes into p independent sub-problems. In order to identify direct effects, one should consider ensemble methods with fully developed trees as the learning algorithm for each sub-problem.

At first sight, GENIE3 seems to be a good candidate for connectome inference. We however noticed that running the learning algorithm p times (typically, p = 1000) was computationally too expensive under time constraints pertaining to a machine learning challenge. We therefore opt for another learning algorithm - based on partial correlation - that is computationally advantageous⁸. Conversely with GENIE3, partial correlation based approach aims at finding explicitly only direct interactions.

⁸Especially for a fast development and parameter tuning.

8.3.2 Connectome inference

The human brain is a complex biological organ made of about 100 billion of neurons, each connected to, on average, 7,000 other neurons [Pakkenberg et al., 2003]. Unfortunately, direct observation of the connectome, the wiring diagram of the brain, is not yet technically feasible. Without being perfect, calcium imaging currently allows for real-time and simultaneous observation of neuron activity from thousands of neurons, producing individual time-series representing their fluorescence intensity. From these data, the connectome inference problem amounts to retrieving the synaptic connections between neurons on the basis of the fluorescence time-series. This problem is difficult to solve because of experimental issues, including masking effects (i.e., some of the neurons are not observed or confounded with others), the low sampling rate of the optical device with respect to the neural activity speed, or the slow decay of fluorescence.

Formally, the connectome can be represented as a directed graph G = (V, E), where V is a set of p nodes representing neurons, and $E \subseteq \{(i,j) \in V \times V\}$ is a set of edges representing direct synaptic connections between neurons. Causal interactions are expressed by the direction of edges: $(i,j) \in E$ indicates that the state of neuron j might be caused by the activity of neuron i. In those terms, the connectome inference problem is formally stated as follows: Given the sampled observations $\{x_i^t \in \mathbb{R} | i \in V, t = 1, ..., T\}$ of p neurons for T time intervals, the goal is to infer the set E of connections in G.

In this section, we present a simplified - and almost as good - version of the winning method⁹ of the Connectomics Challenge¹⁰, as a simple and theoretically grounded approach based on signal processing techniques and partial correlation statistics. The rest of this chapter is structured as follows: Section 8.3.3 describes the signal processing methods applied on fluorescent calcium time-series; Section 8.3.4 then presents the proposed approach and its theoretical properties; Section 8.3.5 provides an empirical analysis and comparison with other network inference methods, while finally, in Section 8.3.6 we discuss our work and provide further research directions. Additionally, Appendix 8.A further describes, in full detail, our actual winning method which gives slightly better results than the method presented in this paper, at the cost of parameter tuning. Appendix 8.B provides supplementary results on other datasets.

8.3.3 Signal processing

Under the simplifying assumption that neurons are on-off units, characterised by short periods of intense activity, or peaks, and longer periods of inactivity, the first part of our algorithm consists of cleaning the raw fluorescence data. More specifically, time-series are processed using standard signal processing filters in order to : (i) remove noise mainly due to fluctuations independent of calcium, calcium fluctuations independent of spiking activity, calcium fluctuations in nearby tissues that have been mistakenly captured, or simply by the imaging process; (ii) to account for fluorescence low decay; and (iii) to reduce the importance of high global activity in the network. The overall process is illustrated in Figure 8.3.1.

⁹Code available at https://github.com/asutera/kaggle-connectomics

¹⁰http://connectomics.chalearn.org

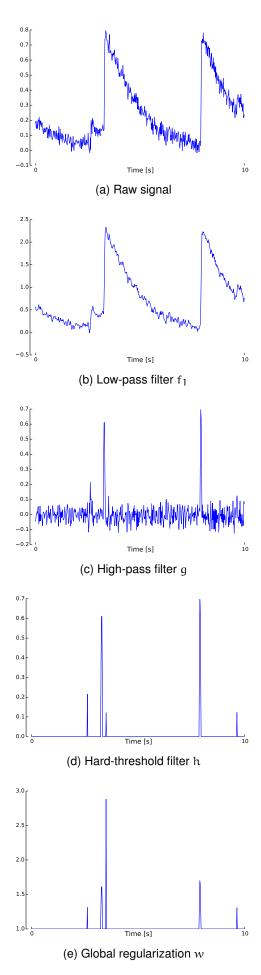


Figure 8.3.1: Signal processing pipeline for extracting peaks from the raw fluorescence data.

As Figure 8.3.1a shows, the raw fluorescence signal is very noisy due to light scattering artifacts that usually affect the quality of the recording [Lichtman and Denk, 2011]. Accordingly, the first step of our pipeline is to smooth the signal, using one of the following low-pass filters for filtering out high frequency noise:

$$f_1(x_i^t) = x_i^{t-1} + x_i^t + x_i^{t+1}, \tag{8.2}$$

$$f_2(x_i^t) = 0.4x_i^{t-3} + 0.8x_i^{t-2} + x_i^{t-1} + x_i^t.$$
(8.3)

These filters are standard in the signal processing field [Kaiser and Reed, 1977; Oppenheim et al., 1983]. For the purposes of illustration, the effect of the filter f₁ on the signal is shown in Figure 8.3.1b.

Furthermore, short spikes, characterized by a high frequency, can be seen as an indirect indicator of neuron communication, while low frequencies of the signal mainly correspond to the slow decay of fluorescence. To have a signal that only has high magnitude around instances where the spikes occur, the second step of our pipeline transforms the time-series into its backward difference

$$g(x_i^t) = x_i^t - x_i^{t-1},$$
 (8.4)

as shown in Figure 8.3.1c.

To filter out small variations in the signal obtained after applying the function q, as well as to eliminate negative values, we use the following hard-threshold filter

$$h(x_i^t) = x_i^t \mathbb{1}(x_i^t \geqslant \tau) \text{ with } \tau > 0, \tag{8.5}$$

yielding Figure 8.3.1d where τ is the threshold parameter and 1 is the indicator function. As can be seen, the processed signal only contains clean spikes.

The objective of the last step of our filtering procedure is to decrease the importance of spikes that occur when there is high global activity in the network with respect to spikes that occur during normal activity. Indeed, we have conjectured that when a large part of the network is firing, the rate at which observations are made is not high enough to be able to detect interactions, and that it would therefore be preferable to lower their importance by changing their magnitude appropriately. Additionally, it is well-known that neurons may also spike because of a high global activity [Stetter et al., 2012]. In such context, detecting pairwise neuron interactions from the firing activity is meaningless. As such, the signal output by h is finally applied to the following function

$$w(x_i^t) = (x_i^t + 1)^{1 + \frac{1}{\sum_j x_j^t}},$$
(8.6)

whose effect is to magnify the importance of spikes that occur in cases of low global activity (measured by $\sum_i x_i^t$), as observed, for instance, around t = 4s in Figure 8.3.1e. Note the particular case where there is no activity, i.e., $\sum_i x_i^t = 0$, is solved by setting $w(x_i^t) = 1$.

To summarise, the full signal processing pipeline of our simplified approach is defined by the composed function $w \circ h \circ g \circ f_1$ (resp. f_2). When applied to the raw signal of Figure 8.3.1a, it outputs the signal shown in Figure 8.3.1e.

8.3.4 Connectome inference from partial correlation statistics

Our procedure to infer connections between neurons first assumes that the (filtered) fluorescence concentrations of all p neurons at each time point can be modelled as a set of random variables $X = \{X_1, \dots, X_p\}$ that are independently drawn from the same time-invariant joint probability distribution P_X . As a consequence, our inference method does not exploit the time-ordering of the observations (although time-ordering is exploited by the filters).

Given this assumption, we then propose to use as a measure of the strength of the connection between two neurons i and j, the Partial correlation coefficient $p_{i,j}$ between their corresponding random variables X_i and X_i , defined by:

$$p_{i,j} = -\frac{\sum_{i,j}^{-1}}{\sqrt{\sum_{i,j}^{-1} \sum_{j,j}^{-1}}},$$
(8.7)

where Σ^{-1} , known as the precision or concentration matrix, is the inverse of the covariance matrix Σ of X. Assuming that the distribution P_X is a multivariate Gaussian distribution $\mathcal{N}(\mu, \Sigma)$, it can be shown that $p_{i,j}$ is zero if and only if X_i and X_j are independent given all other variables in X, i.e., $X_i \perp X_j | X^{-i,j}$ where $X^{-i,j} = X \setminus \{X_i, X_j\}$. Partial correlation (illustrated by Figure 8.3.2) thus measures conditional dependencies between variables; therefore it should naturally only detect direct associations between neurons and filter out spurious indirect effects. The interest of partial correlation as an association measure has already been shown for the inference of gene regulatory networks [De La Fuente et al., 2004; Schäfer and Strimmer, 2005]. Note that the partial correlation statistic is symmetric (i.e. $p_{i,j} = p_{j,i}$). Therefore, our approach cannot identify the direction of the interactions between neurons. We will see in Section 8.3.5 why this only slightly affects its performance, with respect to the metric used in the Connectomics Challenge.

Practically speaking, the computation of all $p_{i,j}$ coefficients using Equation 8.7 requires the estimation of the covariance matrix Σ and then computing its inverse. Given that typically we have more samples than neurons, the covariance matrix can be inverted in a straightforward way. We nevertheless obtained some improvement by replacing the exact inverse with an approximation using only the M first principal components [Bishop, 2006] (with M = 0.8p in our experiments, see Appendix 8.C).

Finally, it should be noted that the performance of our simple method appears to be quite sensitive to the values of parameters (e.g., choice of f₁ or f₂ or the value of the threshold τ) in the combined function of the filtering and inferring processes. One approach, further referred to as Averaged Partial correlation statistics, for improving its robustness is to average correlation statistics over various values of the parameters, thereby reducing the variance of its predictions. Further details about parameter selection are provided in Appendix 8.A.

8.3.5 Experiments

DATA AND EVALUATION METRICS. We report here experiments on the normal-1,2,3, and 4 datasets provided by the organisers of the Connectomics Challenge (see Appendix 8.B for experiments on other datasets). Each of these datasets is obtained from the simulation [Stetter et al., 2012] of different neural networks of 1,000 neurons and approximately 15,000 edges (i.e., a network density of about 1.5%). Each neuron is described by a calcium fluorescence time-series of length T = 179500. All inference methods compared here provide a ranking of all pairs of neurons according to some association score. To assess the quality of this ranking, we compute both ROC and precision-recall curves against the ground-truth

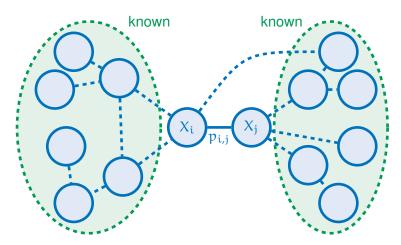


Figure 8.3.2: Partial correlation coefficient $p_{i,j}$ measures the degree of direct association between X_i and X_j given all other nodes (in green areas).

network, which are represented by the area under the curves and respectively denoted AUROC and AUPRC. Only the AUROC score was used to rank the challenge participants, but the precision-recall curve has been shown to be a more sensible metric for network inference, especially when network density is small (see e.g., Schrynemackers et al. [2013]). Since neurons are not self-connected in the ground-truth networks (i.e., $(i,i) \not\in E, \forall i \in V$), we have manually set the score of such edges to the minimum possible association score before computing ROC and PR curves.

EVALUATION OF THE METHOD. The top of Table 8.3.1 reports AUROC and AUPRC for all four networks using, in each case, partial correlation with different filtering functions. Except for the last two rows that use PCA, the exact inverse of the covariance matrix was used in each case. These results clearly show the importance of the filters. AUROC increases in average from 0.77 to 0.93. PCA does not really affect AUROC scores, but it significantly improves AUPRC scores. Taking the average over various parameter settings gives an improvement of 10% in AUPRC but only a minor change in AUROC. The last row ("Full method") shows the final performance of the method specifically tuned for the challenge (see Appendix 8.A for all details). Although this tuning was decisive to obtain the best performance in the challenge, it does not significantly improve either AUROC or AUPRC.

COMPARISON WITH OTHER METHODS. At the bottom of Table 8.3.1, we provide as a comparison the performance of three other methods: standard (Pearson) correlation (PC), generalised transfer entropy (GTE), and GENIE3. ROC and PR curves on the *normal-2* network are shown for all methods in Figure 8.3.3. Pearson correlation measures the unconditional linear (in)dependence between variables and it should thus not be able to filter out indirect interactions between neurons. GTE [Stetter et al., 2012] was proposed as a baseline for the challenge. This method builds on Transfer Entropy to measure the association between two neurons. Unlike our approach, it can predict the direction of the edges. GENIE3 [Huynh-Thu et al., 2010] is a gene regulatory network inference method that was the best performer in the DREAM5 challenge [Marbach et al., 2012] (more details are given in Section 8.2.1). When transposed to neural networks, this method uses the importance score of variable X_i in a Random Forest model trying to predict X_j from all variables in $X \setminus X_j$ as a confidence score for the edge going from neuron i to neuron

		AUF	ROC			AUF	PRC	
Method \ normal-	1	2	3	4	1	2	3	4
No filtering	0.777	0.767	0.772	0.774	0.070	0.064	0.068	0.072
$h \circ g \circ f_1$	0.923	0.925	0.923	0.922	0.311	0.315	0.313	0.304
$w \circ h \circ g \circ f_1$	0.931	0.929	0.928	0.926	0.326	0.323	0.319	0.303
+ PCA	0.932	0.930	0.928	0.926	0.355	0.353	0.350	0.333
Averaging	0.937	0.935	0.935	0.931	0.391	0.390	0.385	0.375
Full method	0.943	0.942	0.942	0.939	0.403	0.404	0.398	0.388
PC	0.886	0.884	0.891	0.877	0.153	0.145	0.170	0.132
GTE	0.890	0.893	0.894	0.873	0.171	0.174	0.197	0.142
GENIE3	0.892	0.891	0.887	0.887	0.232	0.221	0.237	0.215

Table 8.3.1: Top: Performance on normal-1,2,3,4 with partial correlation and different filtering functions. Bottom: Performance on *normal-1,2,3,4* with different methods.

j. However, to reduce the computational cost of this method, we had to limit each tree in the Random Forest model to a maximum depth of 3. This constraint has a potentially severe effect on the performance of this method with respect to the use of fully-grown trees. PC and GENIE3 were applied to the time-series filtered using the functions $w \circ h \circ g$ and $h \circ g \circ f_1$ (which gave the best performance), respectively. For GENIE3, we built 10,000 trees per neuron and we used default settings for all other parameters (except for the maximal tree depth). For GTE, we reproduced the exact same setting (conditioning level and pre-processing) that was used by the organisers of the challenge.

Partial correlation and averaged partial correlation clearly outperform all other methods on all datasets (see Table 8.3.1 and Appendix 8.B). The improvement is more important in terms of AUPRC than in terms of AUROC. As expected, Pearson correlation performs very poorly in terms of AUPRC. GTE and GENIE3 work much better, but these two methods are nevertheless clearly below partial correlation. Among these two methods, GTE is slightly better in terms of AUROC, while GENIE3 is significantly better in terms of AUPRC. Given that we had to limit this latter method for computational reasons, these results are very promising and a comparison with the full GENIE3 approach is certainly part of our future works.

The fact that our method is unable to predict edge directions does not seem to be a disadvantage with respect to GTE and GENIE3. Although partial correlation scores each edge, and its opposite, similarly, it can reach precision values higher than 0.5 (see Figure 8.3.3(b)), suggesting that it mainly ranks high pairs of neurons that interact in both directions. It is interesting also to note that, on normal-2, a method that perfectly predicts the undirected network (i.e., that gives a score of 1 to each pair (i,j) such that $(i,j) \in E$ or $(j,i) \in E$, and 0 otherwise) already reaches an AUROC as high as 0.995 and an AUPRC of 0.789.

8.3.6 Conclusion for connectome inference

In this section, we outlined a simple but efficient methodology for the problem of connectome inference from calcium imaging data. Our approach consists of two steps: (i) processing fluorescence data to detect neural peak activities; and (ii) inferring the degree of association between neurons from partial correlation statistics. Its simpli-

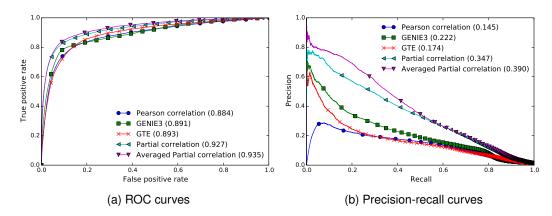


Figure 8.3.3: ROC (left) and PR (right) curves on *normal-2* for the compared methods. Areas under the curves are reported in the legend.

fied variant outperforms other network inference methods while its optimized version proved to be the best method on the Connectomics Challenge. Given its simplicity and good performance, we therefore believe that the methodology presented in this work would constitute a solid and easily-reproducible baseline for further work in the field of connectome inference.

8.A DESCRIPTION OF THE "FULL METHOD"

This section provides a detailed description of the method specifically tuned for the Connectomics Challenge. We restrict our description to the differences with respect to the simplified method presented in the main paper. Most parameters were tuned so as to maximize AUROC on the *normal-1* dataset and our design choices were validated by monitoring the AUROC obtained by the 145 entries we submitted during the challenge. Although the tuned method performs better than the simplified one on the challenge dataset, we believe that the tuned method clearly overfits the simulator used to generate the challenge data and that the simplified method should work equally well on new independent datasets. We nevertheless provide the tuned method here for reference purposes. Our implementation of the tuned method is available at https://github.com/asutera/kaggle-connectomics.

This appendix is structured as follows: Section 8.A.1 describes the differences in terms of signal processing. Section 8.A.2 then provides a detailed presentation of the averaging approach. Section 8.A.3 presents an approach to correct the $p_{i,j}$ values so as to take into account the edge directionality. Finally, Section 8.A.4 presents some experimental results to validate the different steps of our proposal.

8.A.1 Signal processing

In Section 8.3.3, we introduced four filtering functions (f, g, h, and w) that are composed in sequence (i.e., $w \circ h \circ g \circ f$) to provide the signals from which to compute partial correlation statistics. Filtering is modified as follows in the tuned method:

 In addition to f₁ and f₂ (Equations 8.2 and 8.3), two alternative low-pass filters f₃ and f₄ are considered:

$$f_3(x_i^t) = x_i^{t-1} + x_i^t + x_i^{t+1} + x_i^{t+2},$$
(8.8)

$$f_4(x_i^t) = x_i^t + x_i^{t+1} + x_i^{t+2} + x_i^{t+3}.$$
(8.9)

 An additional filter r is applied to smoothe differences in peak magnitudes that might remain after the application of the hard-threshold filter h:

$$\mathbf{r}(\mathbf{x}_{i}^{t}) = (\mathbf{x}_{i}^{t})^{c}, \tag{8.10}$$

with c = 0.9.

Filter w is replaced by a more complex filter w* defined as:

$$w^*(x_i^t) = (x_i^t + 1)^{\left(1 + \frac{1}{\sum_j x_j^t}\right)^{k(\sum_j x_j^t)}}$$
(8.11)

where the function k is a piecewise linear function optimised separately for each filter f_1 , f_2 , f_3 and f_4 (see the implementation for full details). Filter w in the simplified method is a special case of w^* with $k(\sum_j x_j^t) = 1$.

The pre-processed time-series are then obtained by the application of the following function: $w^* \circ r \circ h \circ g \circ f_i$ (with i = 1, 2, 3, or 4).

Weighted average of partial correlation statistics 8.A.2

As discussed in Section 8.3.4, the performance of the method (in terms of AUROC) is sensitive to the value of the parameter τ of the hard-threshold filter h (see Equation 8.5), and to the choice of the low-pass filter (among $\{f_1, f_2, f_3, f_4\}$). As in the simplified method, we have averaged the partial correlation statistics obtained for all the pairs $(\tau, low-pass filter) \in \{0.100, 0.101, \dots, 0.210\} \times \{f_1, f_2, f_3, f_4\}.$

Filters f₁ and f₂ display similar performances and thus were given similar weights (i.e., resp. 0.383 and 0.345). These weights were chosen equal to the weights selected for the simplified method. In contrast, filters f₃ and f₄ turn out, individually, to be less competitive and were therefore given less importance in the weighted average (i.e., resp. 0.004 and 0.268). Yet, as further shown in Section 8.A.4, combining all 4 filters proves to marginally improve performance with respect to using only f₁ and f_2 .

8.A.3 Prediction of edge orientation

Partial correlation statistics is a symmetric measure, while the connectome is a directed graph. It could thus be beneficial to try to predict edge orientation. In this section, we present an heuristic that modifies the p_{ii} computed by the approach described before which takes into account directionality.

This approach is based on the following observation. The rise of fluorescence of a neuron indicates its activation. If another neuron is activated after a slight delay, this could be a consequence of the activation of the first neuron and therefore indicates a directed link in the connectome from the first to the second neuron. Given this observation, we have computed the following term for every pair (i, j):

$$s_{i,j} = \sum_{t=1}^{T-1} \mathbb{1}((x_j^{t+1} - x_i^t) \in [\phi_1, \phi_2])$$
 (8.12)

that could be interpreted as an image of the number of times that neuron i activates neuron j. ϕ_1 and ϕ_2 are parameters whose values have been chosen in our experiments equal to 0.2 and 0.5, respectively. Their role is to define when the difference between x_i^{t+1} and x_i^t can indeed be assimilated to an event for which neuron i activates neuron j.

Afterwards, we have computed the difference between $s_{i,j}$ and $s_{j,i}$, that we call $z_{i,j}$, and used this difference to modify $p_{i,j}$ and $p_{j,i}$ so as to take into account directionality. Naturally, if $z_{i,j}$ is greater (smaller) than 0, we may conclude that should there be an edge between i and j, then this edge would have to be oriented from i to j (j to i).

This suggests the new association matrix r:

$$r_{i,j} = 1(z_{i,j} > \phi_3) * p_{i,j}$$
 (8.13)

where $\phi_3 > 0$ is another parameter. We discovered that this new matrix r was not providing good results, probably due to the fact that directivity was not rewarded well enough in the challenge.

This has lead us to investigate other ways for exploiting the information about directionality contained in the matrix z. One of those ways that gave good performance was to use as an association matrix:

$$q_{i,j} = weight * p_{i,j} + (1 - weight) * z_{i,j}$$
 (8.14)

with weight chosen close to 1 (weight = 0.997). Note that with values for weight close to 1, matrix q only uses the information to a minimum about directivity contained in z to modify the partial correlation matrix p. We tried smaller values for weight but those provided poorer results.

It was this association matrix $q_{i,j}$ that actually led to the best results of the challenge, as shown in Table 8.A.2 of Section 8.A.4.

8.A.4 Experiments

ON THE INTEREST OF LOW-PASS FILTERS f_3 AND f_4 . As reported in Table 8.A.1, averaging over all low-pass filters leads to better AUROC scores than averaging over only two low-pass filters, i.e., f₁ and f₂. However this slightly reduces AUPRC.

Table 8.A.1: Performance on normal-1, 2, 3, or 4 with partial correlation with different averaging approaches.

	AUROC					AUF	PRC	
Averaging \ normal-	1	2	3	4	1	2	3	4
with f ₁ , f ₂	0.937	0.935	0.935	0.931	0.391	0.390	0.385	0.375
with f ₁ , f ₂ , f ₃ , f ₄	0.938	0.936	0.936	0.932	0.391	0.389	0.385	0.374

ON THE INTEREST OF USING MATRIX Q RATHER THAN p TO TAKE INTO AC-Table 8.A.2 compares AUROC and AUPRC with or with-COUNT DIRECTIVITY. out correcting the p_{i,i} values according to Equation 8.14. Both AUROC and AUPRC are (very slightly) improved by using information about directivity.

Table 8.A.2: Performance on normal-1,2,3,4 of "Full Method" with and without using information about directivity.

		AUF	ROC			AUF	PRC	
Full method \ normal-	1	1 2 3 4				2	3	4
Undirected	0.943	0.942	0.942	0.939	0.403	0.404	0.398	0.388
Directed	0.944	0.943	0.942	0.940	0.404	0.405	0.399	0.389

8.B SUPPLEMENTARY RESULTS

In this appendix we report the performance of the different methods compared in the paper on 6 additional datasets provided by the Challenge organisers. These datasets, corresponding each to networks of 1,000 neurons, are similar to the normal datasets except for one feature:

LOWCON: Similar network but on average with a lower number of connections per neuron.

HIGHCON: Similar network but on average with a higher number of connections per neuron.

LOWCC: Similar network but on average with a lower clustering coefficient.

HIGHCC: Similar network but on average with a higher clustering coefficient.

NORMAL-3-HIGHRATE: Same topology as normal-3 but with a higher firing frequency, i.e., with highly active neurons.

NORMAL-4-LOWNOISE: Same topology as normal-4 but with a better signal-tonoise ratio.

The results of several methods applied to these 6 datasets are provided in Table 8.B.1. They confirm what we observed on the *normal* datasets. Average partial correlation and its tuned variant, i.e., "Full method", clearly outperform other network inference methods on all datasets. PC is close to GENIE3 and GTE, but still slightly worse. GENIE3 performs better than GTE most of the time. Note that the "Full method" reported in this table does not use Equation 8.14 to slightly correct the values of $p_{i,j}$ to take into account directivity.

Table 8.B.1: Performance (top: AUROC, bottom: AUPRC) on specific datasets with different methods.

	AUROC							
Method \ normal-	lowcon	highcon	lowcc	highcc	3-highrate	4-lownoise		
Averaging	0.947	0.943	0.920	0.942	0.959	0.934		
Full method	0.955	0.944	0.925	0.946	0.961	0.941		
PC	0.782	0.920	0.846	0.897	0.898	0.873		
GTE	0.846	0.905	0.848	0.899	0.905	0.879		
GENIE3	0.781	0.924	0.879	0.902	0.886	0.890		
			ı	AUPRC				
Averaging	0.320	0.429	0.262	0.478	0.443	0.412		
Full method	0.334	0.413	0.260	0.486	0.452	0.432		
PC	0.074	0.218	0.082	0.165	0.193	0.135		
GTE	0.094	0.211	0.081	0.165	0.210	0.144		
GENIE3	0.128	0.273	0.116	0.309	0.256	0.224		

8.C ON THE SELECTION OF THE NUMBER OF PRINCIPAL COMPONENTS

The (true) network, seen as a matrix, can be decomposed through a singular value decomposition (SVD) or principal component analysis (PCA), so as to respectively determine a set of independent linear combinations of the variable [Alter et al., 2000], or a reduced set of linear combinations combine, which then maximize the explained variance of the data [Jolliffe, 2005]. Since SVD and PCA are related, they can be defined by the same goal: both aim at finding a reduced set of neurons, known as components, whose activity can explain the rest of the network.

The distribution of component eigen values obtained from PCA and SVD decompositions can be studied by sorting them in descending order of magnitude, as illustrated in Figure 8.C.1. It can be seen that some component eigen values are zero, implying that the behaviour of the network could be explained by a subset of neurons because of the redundancy and relations between the neurons. For all datasets, the eigen value distribution is exactly the same.

In the context of the challenge, we observe that only 800 components seem to be necessary and we exploit this when computing partial correlation statistics. Therefore, the value of the parameter M is immediate and should be clearly set to 800 (= 0.8p).

Note that if the true network is not available, similar decomposition analysis could be carried on the inferred network, or on the data directly.

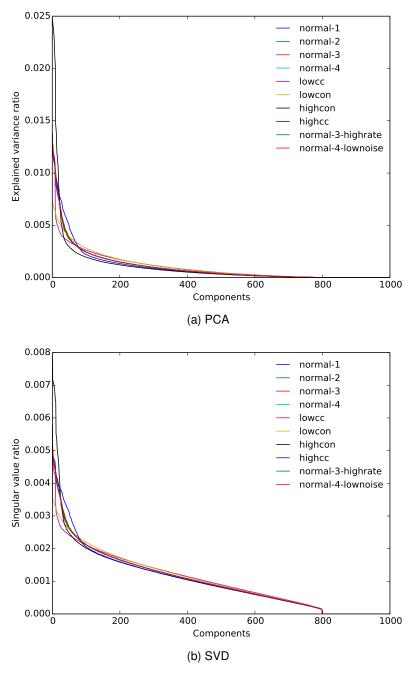


Figure 8.C.1: Explained variance ratio by number of principal components (left) and singular value ratio by number of principal components (right) for all networks.

9

Overview

The objective of this thesis was to better understand and characterise the properties of tree-based feature importance measures. Indeed, despite numerous works from either empirical or theoretical points of view, these importance measures are not yet fully understood. We are convinced that a more in depth understanding of the various properties of those measures would help to foster the scientific community to more systematically exploit these measures within the context of a wide variety of problems and methods. Within this context, we have mostly focused our study on the so-called *Mean Decrease of Impurity* type of importance measure. In this chapter we summarise our findings and discuss directions for further research.

9.1 MAIN FINDINGS

In the first part of this thesis we gave the background for the subsequent chapters. In particular, we introduced various notions of feature relevance and redundancy between features and described various feature selection problems in Chapter 2, and we presented all relevant notions and algorithms pertaining to tree-based methods in Chapter 3.

Our first step towards a better understanding of tree-based feature importance measures consisted of a survey of the literature about this topic, provided in Chapter 4. We proposed a framework of the MDA approach that is not tree-specific. In asymptotic conditions, i.e. infinite sample size N and infinite ensemble size N_T, we gathered analytical formulations of both MDA and MDI importance measures and highlighted their main properties, in particular in the presence of correlated or redundant features. From a more practical point of view, we discussed their main biases. Despite many desirable properties, it emerged from empirical analyses that tree-based parameters, and feature characteristics and dependencies, may strongly impact the measured importance scores. In particular, the split randomisation parameter K (i.e., the number of features considered at each node as split variable candidates) introduces the so-called masking effect, that prevents some relevant features to appear as important to the eyes of a random forest model. We also noticed a preference for smaller groups of correlated features, worth to take into account in the context of high dimensional applications where features often come in groups of correlated features of variable sizes. All those observations should help to analyse more cautiously importance scores.

Another downside of tree-based importance measures is that they do not provide an explicit way to distinguish important features from non-important ones, for example by providing meaningful thresholds on feature importances. However, many approaches have been proposed to circumvent this issue. In particular, we investigated permutation schemes and pointed out that the conditional permutation scheme proposed in Strobl et al. [2008] focuses only on strongly relevant features, those conveying unique information about the output variable (in contrast with weakly relevant features).

In Chapter 5 we focused on the MDI importance measure, and extended its characterisation from totally randomised trees (K = 1) to more realistic random forest algorithms in asymptotic conditions. While all relevant features receive positive MDI values when using totally randomised trees, non-totally randomised trees (K > 1)guarantee non zero importance scores only for strongly relevant features. Depending on the value of K, more or fewer weakly relevant features might be missed. In case of non-totally developed trees, we related these properties to the maximal tree depth, the feature degree of interaction, and the number of relevant features.

However, in all these theoretical analyses, trees with multiway splits were considered, while in practice binary trees are generally preferred. We therefore transposed results obtained for multiway trees to binary ones. Relaxing the asymptotic conditions, we discussed the implications of a finite setting. Limiting the size of the forest increases the number of "missed" features. In addition to masked features, some features are never evaluated, or not often enough to estimate their true importance. Importances derived from a finite sample suffer from a positive bias that makes all features, including the irrelevant ones, have strictly positive importances.

In many problems, feature selection is usually more complicated than identifying a single subset of input features that would together explain the output. Therefore we proposed in Chapter 6 a methodological contribution that takes into account the context (i.e., the circumstances that form the setting for the experiment) in feature importance evaluations. The characterisation considers both contextual and noncontextual relevance of features and is based on several importance scores derived from tree-based methods. This approach was also illustrated on two artificial and two real biomedical problems.

When facing high-dimensional datasets, most approaches suffer from the curse of dimensionality. Chapter 7 proposed an improved tree-based method that handles large datasets while being computationally tractable and able to identify relevant features efficiently. Marginally relevant features can be easily identified, even by univariate approaches, however some features are only relevant in the context of others, and their identification requires sophisticated methods that handle feature dependencies. The key idea of our method is that all features that make the others appear as relevant are necessarily relevant too. We used this simple result to propose a sequential approach that keeps in memory some already identified relevant features to speed up the identification of others. We observed that this approach is particularly interesting in case of highly dependent and structured features.

The last chapter of this thesis is devoted to a specific machine learning task consisting in reconstructing a network from data. The first part of Chapter 8 recalled the principle of GENIE3, a tree-based network inference technique designed a few years ago in our research group. Then, we proposed a method to infer the connectome from calcium imaging data using partial correlation statistics and we put it in perspective with GENIE3-like techniques.

9.2 LIMITATIONS AND FUTURE WORK

We believe that this thesis provides additional steps towards a better understanding of tree-based feature importance measures. However, there still remain several limitations to the frameworks proposed along this thesis that are all potential directions of improvements.

EXTENDING OUR CHARACTERISATION OF THE MDI IMPORTANCE TO CON-TINUOUS FEATURES

All theoretical derivations from Chapters 4, 5 and 6 concern categorical input variables which are the keystones of our characterisation of the measure. It would be interesting to adapt our framework to continuous input variables, and also, probably with more difficulty, to continuous context variables.

FEATURE IMPORTANCE ESTIMATION IN NON ASYMPTOTIC CONDITIONS

Another key assumption of our characterisations was asymptotic conditions. In practice sample sizes are finite, as are the number of trees in an ensemble. We believe that a very significant step towards a full understanding of importance measures would be the derivation of statistical distributions of importance scores depending on feature characteristics as well as sample and ensemble size.

MDA VS. MDI

The study conducted in Chapter 4 initiated a comparison between the two treebased importance measures. Both methods can be used for classification and regression problems and yield similar results while being intrinsically different on several aspects. In what follows, we give an outline of some elements of comparison that may be subject to future studies.

MDA exploits out-of-bag (OOB) samples (i.e., not used to learn the model) to compute an error-rate evaluating the impact of the removal (by permutation) of a feature. MDI assesses the importance of a feature based on its average contribution in the impurity reduction in the tree-ensemble learning. Unlike MDA, MDI therefore uses the same samples for learning the model and evaluating the importance of features. Future research could examine if importance scores evaluated in this way and those computed using MDI on independent samples (e.g., OOB samples or a holdout test set) are similar.

Another point of comparison is that MDA depends explicitly on the loss function used, whereas MDI depends explicitly on the impurity measure used. However, MDA also depends indirectly on the tree structure and hence on the impurity measure used to grow the tree. Permuting a feature that is not used in the tree model obviously does not impact the OOB error-rate of this tree. This suggests that both importance measures will identify approximately the same set of important features. Some differences can however be pointed out. Let us consider a two-class classification problem. A feature that slightly changes the output value distributions in the tree leaves would be seen as important by the MDI importance measure. If this change is too subtle to change the predicted class of the tree, the MDA importance

using a loss function that is not sensitive enough (e.g., the zero-one loss function) would miss such a feature. In presence of two variables in a XOR configuration with respect to the output, MDI is only able to identify one feature (the second one) per tree, whereas MDA can detect the importance of both features in a single tree. Indeed, permuting the values of the first variable used induces that samples reaching nodes using the second variable are mixed up. This necessarily impacts the errorrate of the tree and thus makes the first variable appear as important in the eyes of MDA.

Future studies could investigate if it is possible to move both importance closer to each other by considering some specific loss function (e.g., that would have the same properties as the impurity measure used).

EMPIRICAL EVALUATION AND IMPROVEMENTS OF THE SEQUENTIAL RAN-DOM SUBSPACE METHOD

Despite a strong theoretical motivation, some more work is clearly needed to evaluate the Sequential Random Subspace method empirically, on controlled and real high-dimensional problems. In this context, it would probably be necessary to overcome one of the main drawbacks of the sequential random subspace method with respect to the random subspace method which is that it can not be parallelised. One possible approach is to grow ensemble of trees at each iteration instead of single trees. Such a variant is clearly an improvement for our proposed method. We finally believe that another possible improvement to our algorithm is the statistical test based on the introduction of a random probe used to decide which feature to include in the relevant set.

9.3 OPEN RESEARCH QUESTIONS

Alongside future work resulting directly from this thesis, we propose in this section some new (open-)research questions that go beyond the scope of this thesis but that should be investigated by further studies to complete the understanding of treebased feature importance measures.

FEATURE IMPORTANCE CHARACTERISATION FOR TREE BOOSTING

Boosting approaches were not discussed in this thesis. They however constitute powerful and well performing ensemble algorithms. Concretely, in tree-based boosting ensemble methods, trees are not built independently but sequentially in order to correct predictions of previous trees. Each tree is therefore weighted according to its contribution to the total model performance. Given the state-of-the-art performance of these methods, it would be interesting to compute feature importances from these ensembles of trees by extending our formulation to take into account tree weights, and examine to what extent asymptotic guarantees are still valid.

IMPROVING INTERPRETABILITY OF OTHER STATE-OF-THE-ART MACHINE LEARNING MODELS

This thesis was devoted to tree-based methods only. However, all machine learning algorithms could benefit from more interpretability of their induced models, particularly deep learning methods. Taking inspiration from feature importance derived from tree-based methods, it would interesting to evaluate the MDA approach (that is not tree-specific) on other machine learning algorithms and compare it to other importance measures (e.g., individual feature importance measures that assess the importance of features for a single prediction, and which were not discussed in this thesis).

CAUSALITY

The main advantage of our partial correlation approach is to filter out indirect links that the tree-based inference method GENIE3 is unable to do. Future work might investigate the relationship between direct links and strong relevance, and evaluate to what extent it may be possible to reduce the number of indirect links that are actually kept in the final reconstructed network. Causality in tree-based methods has been considered in only few works (see, e.g., [Li et al., 2017]) and still remains an open question to date.

- T. Abeel, T. Helleputte, Y. Van de Peer, P. Dupont, and Y. Saeys. Robust biomarker identification for cancer diagnosis with ensemble feature selection methods. *Bioinformatics*, 26(3):392–398, 2009. (Cited on page 55.)
- C. F. Aliferis, I. Tsamardinos, and A. Statnikov. Hiton: a novel markov blanket algorithm for optimal variable selection. In *AMIA Annual Symposium Proceedings*, volume 2003, page 21. American Medical Informatics Association, 2003. (Cited on page 36.)
- C. F. Aliferis, A. Statnikov, I. Tsamardinos, S. Mani, and X. D. Koutsoukos. Local causal and markov blanket induction for causal discovery and feature selection for classification part i: Algorithms and empirical evaluation. *Journal of Machine Learning Research*, 11(Jan):171–234, 2010. (Cited on pages 37, 47, and 188.)
- H. Almuallim and T. G. Dietterich. Efficient algorithms for identifying relevant features. In *Proc. of the 9th Canadian Conference on Artificial Intelligence*, pages 38–45. Citeseer, 1991a. (Cited on pages 49 and 51.)
- H. Almuallim and T. G. Dietterich. Learning with many irrelevant features. In *AAAI*, volume 91, pages 547–552. Citeseer, 1991b. (Cited on pages 33, 49, and 51.)
- H. Almuallim and T. G. Dietterich. Learning boolean concepts in the presence of many irrelevant features. *Artificial Intelligence*, 69(1-2):279–305, 1994. (Cited on pages 49 and 51.)
- G. Altay, M. Asim, F. Markowetz, and D. E. Neal. Differential c3net reveals disease networks of direct physical interactions. *BMC bioinformatics*, 12(1):296, 2011. (Cited on page 54.)
- O. Alter, P. O. Brown, and D. Botstein. Singular value decomposition for genome-wide expression data processing and modeling. *Proceedings of the National Academy of Sciences*, 97(18):10101–10106, 2000. (Cited on page 200.)
- A. Altmann, L. Toloşi, O. Sander, and T. Lengauer. Permutation importance: a corrected feature importance measure. *Bioinformatics*, 26(10):1340–1347, 2010. (Cited on pages 109, 110, and 111.)
- D. Amaratunga, J. Cabrera, and Y.-S. Lee. Enriched random forests. *Bioinformatics*, 24(18):2010–2014, 2008. (Cited on page 112.)
- C. Ambroise and G. J. McLachlan. Selection bias in gene extraction on the basis of microarray gene-expression data. *Proceedings of the national academy of sciences*, 99(10):6562–6566, 2002. (Cited on pages 44 and 54.)
- Y. Amit and D. Geman. Shape quantization and recognition with randomized trees. *Neural computation*, 9(7):1545–1588, 1997. (Cited on page 77.)

- C. V. Ananth and E. F. Schisterman. Confounding, causality, and confusion: the role of intermediate variables in interpreting observational studies in obstetrics. *American journal of obstetrics and gynecology*, 217(2):167–175, 2017. (Cited on pages 25 and 27.)
- K. Archer and R. Kimes. Empirical characterization of random forest variable importance measures. *Computational Statistics & Data Analysis*, 52(4):2249–2260, 2008. (Cited on pages 85, 99, 101, and 108.)
- A. Arnold, R. Nallapati, and W. W. Cohen. A comparative study of methods for transductive transfer learning. In *Data Mining Workshops, 2007. ICDM Workshops 2007. Seventh IEEE International Conference on*, pages 77–82. IEEE, 2007. (Cited on page 31.)
- L. Auret and C. Aldrich. Empirical comparison of tree ensemble variable importance measures. *Chemometrics and Intelligent Laboratory Systems*, 105(2):157–170, 2011. (Cited on pages 85, 86, 99, 100, and 102.)
- D. Battaglia, I. Guyon, V. Lemaire, J. Orlandi, B. Ray, and J. Soriano, editors. *Neural Connectomics Challenge*. Springer, 2017. (Cited on page 189.)
- J. Beirlant, E. J. Dudewicz, L. Györfi, and E. C. Van der Meulen. Nonparametric entropy estimation: An overview. *International Journal of Mathematical and Statistical Sciences*, 6(1):17–39, 1997. (Cited on page 107.)
- M. Belgiu and L. Drăguţ. Random forest in remote sensing: A review of applications and future directions. *ISPRS Journal of Photogrammetry and Remote Sensing*, 114:24–31, 2016. (Cited on page 91.)
- D. A. Bell and H. Wang. A formalism for relevance and its application in feature subset selection. *Machine learning*, 41(2):175–195, 2000. (Cited on page 35.)
- G. Biau. Analysis of a random forests model. *Journal of Machine Learning Research*, 13(Apr):1063–1095, 2012. (Cited on page 79.)
- G. Biau and E. Scornet. A random forest guided tour. *Test*, 25(2):197–227, 2016. (Cited on page 90.)
- G. Biau, L. Devroye, and G. Lugosi. Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, 9(Sep):2015–2033, 2008. (Cited on pages 79 and 92.)
- C. M. Bishop. *Pattern recognition and machine learning*, volume 1. Springer New York, 2006. (Cited on page 193.)
- P. Bloebaum, D. Janzing, T. Washio, S. Shimizu, and B. Schoelkopf. Cause-effect inference by comparing regression errors. In *International Conference on Artificial Intelligence and Statistics*, pages 900–909, 2018. (Cited on page 189.)
- A. L. Blum and P. Langley. Selection of relevant features and examples in machine learning. *Artificial intelligence*, 97(1-2):245–271, 1997. (Cited on pages 32, 33, 47, 48, 49, 50, and 51.)
- V. Bolón-Canedo, N. Sánchez-Maroño, and A. Alonso-Betanzos. Recent advances and emerging challenges of feature selection in the context of big data. Knowledge-Based Systems, 86:33–45, 2015. (Cited on page 54.)

- V. Botta. A walk into random forests: adaptation and application to Genome-Wide Association Studies. PhD thesis, Université de Liège, Liège, Belgique, 2013. (Cited on page 60.)
- V. Botta, G. Louppe, P. Geurts, and L. Wehenkel. Exploiting snp correlations within random forest for genome-wide association studies. *PloS one*, 9(4):e93379, 2014. (Cited on page 99.)
- A.-L. Boulesteix and M. Slawski. Stability and aggregation of ranked gene lists. *Briefings in bioinformatics*, 10(5):556–568, 2009. (Cited on page 102.)
- A.-L. Boulesteix, A. Bender, J. Lorenzo Bermejo, and C. Strobl. Random forest gini importance favours snps with large minor allele frequency: impact, sources and recommendations. *Briefings in Bioinformatics*, 13(3):292–304, 2011. (Cited on page 106.)
- A.-L. Boulesteix, S. Janitza, J. Kruppa, and I. R. König. Overview of random forest methodology and practical guidance with emphasis on computational biology and bioinformatics. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(6):493–507, 2012. (Cited on page 106.)
- O. Bousquet. Transductive learning: Motivation, models, algorithms. *University of New Mexico, Albuquerque, USA*, 2002. (Cited on page 31.)
- C. Boutilier, N. Friedman, M. Goldszmidt, and D. Koller. Context-specific independence in bayesian networks. In *Proceedings of the Twelfth International Conference on Uncertainty in Artificial Intelligence*, UAI'96, pages 115–123, San Francisco, CA, USA, 1996. Morgan Kaufmann Publishers Inc. ISBN 1-55860-412-X. URL http://dl.acm.org/citation.cfm?id=2074284.2074298. (Cited on page 150.)
- U. M. Braga-Neto and E. R. Dougherty. Is cross-validation valid for small-sample microarray classification? *Bioinformatics*, 20(3):374–380, 2004. (Cited on page 54.)
- L. Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996a. (Cited on pages 74 and 75.)
- L. Breiman. Heuristics of instability and stabilization in model selection. *The Annals of Statistics*, pages 2350–2383, 1996b. (Cited on page 71.)
- L. Breiman. Out-of-bag estimation, 1996c. (Cited on page 81.)
- L. Breiman. Some infinity theory for predictor ensembles. Technical report, Technical Report 579, Statistics Dept. UCB, 2000. (Cited on page 79.)
- L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001. (Cited on pages 57, 77, 79, 85, 88, 89, 90, 117, 126, 150, and 160.)
- L. Breiman. Manual on setting up, using, and understanding random forests v3. 1. Statistics Department University of California Berkeley, CA, USA, 1, 2002. (Cited on pages 79, 81, 85, and 88.)
- L. Breiman. Consistency for a simple model of random forests. Technical report, Berkeley, 2004. (Cited on page 79.)

- L. Breiman and A. Cutler. Random forests manual v4. In *Technical report*. UC Berkel, 2003. (Cited on pages 79, 81, 85, 89, and 117.)
- L. Breiman and A. Cutler. Random forests—classification manual. *URL http://www.math. usu. edu/~adele/forests*, 2008. (Cited on page 109.)
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. Classification and Regression Trees. Wadsworth and Brooks, Monterey, CA, 1984. (Cited on pages 57, 63, 64, 65, 67, 68, 71, 72, 86, 88, 106, 121, 136, 137, 156, and 241.)
- G. Brown. A new perspective for information theoretic feature selection. In *International conference on artificial intelligence and statistics*, pages 49–56, 2009. (Cited on page 149.)
- G. Brown, A. Pocock, M.-J. Zhao, and M. Luján. Conditional likelihood maximisation: a unifying framework for information theoretic feature selection. *The Journal of Machine Learning Research*, 13(1):27–66, 2012. (Cited on pages 47, 48, and 149.)
- A. Bureau, J. Dupuis, K. Falls, K. L. Lunetta, B. Hayward, T. P. Keith, and P. Van Eerdewegh. Identifying snps predictive of phenotype using random forests. *Genetic Epidemiology: The Official Publication of the International Genetic Epidemiology Society*, 28(2):171–182, 2005. (Cited on page 99.)
- C. Cardie. Using decision trees to improve case-based learning. In *Proceedings* of the tenth international conference on machine learning, pages 25–32, 1993. (Cited on page 48.)
- B. Carlson. Snps-a shortcut to personalized medicine. *Genetic Engineering & Biotechnology News*, 28(12):12–12, 2008. (Cited on page 106.)
- G. Chandrashekar and F. Sahin. A survey on feature selection methods. *Computers & Electrical Engineering*, 40(1):16–28, 2014. (Cited on pages 47, 48, 50, and 51.)
- N. V. Chawla, L. O. Hall, K. W. Bowyer, and W. P. Kegelmeyer. Learning ensembles from bites: A scalable and accurate approach. *J. Mach. Learn. Res.*, 5:421–451, Dec. 2004. ISSN 1532-4435. (Cited on page 165.)
- T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012. (Cited on pages 34, 42, and 239.)
- T. M. Cover and J. M. Van Campenhout. On the possible orderings in the measurement selection problem. *IEEE Trans. Systems, Man, and Cybernetics*, 7(9): 657–661, 1977. (Cited on page 47.)
- A. Cutler and G. Zhao. Pert-perfect random tree ensembles. *Computing Science and Statistics*, 33:490–497, 2001. (Cited on page 77.)
- D. R. Cutler, T. C. Edwards, K. H. Beard, A. Cutler, K. T. Hess, J. Gibson, and J. J. Lawler. Random forests for classification in ecology. *Ecology*, 88(11):2783–2792, 2007. (Cited on page 91.)
- I. M. de Abril, J. Yoshimoto, and K. Doya. Connectivity inference from neural recording data: Challenges, mathematical bases and research directions. *Neural Networks*, 2018. (Cited on page 186.)

- A. De La Fuente, N. Bing, I. Hoeschele, and P. Mendes. Discovery of meaningful associations in genomic data using partial correlation coefficients. *Bioinformatics*, 20(18):3565–3574, 2004. (Cited on page 193.)
- R. De Smet and K. Marchal. Advantages and limitations of current network inference methods. *Nature Reviews Microbiology*, 8(10):717, 2010. (Cited on pages 54 and 185.)
- M. Del Campo, B. Mollenhauer, A. Bertolotto, S. Engelborghs, H. Hampel, A. H. Simonsen, E. Kapaki, N. Kruse, N. Le Bastard, S. Lehmann, et al. Recommendations to standardize preanalytical confounding factors in alzheimer's and parkinson's disease cerebrospinal fluid biomarkers: an update. *Biomarkers in medicine*, 6(4):419–430, 2012. (Cited on page 27.)
- H. Deng and G. Runger. Feature selection via regularized trees. In *Neural Networks* (*IJCNN*), *The 2012 International Joint Conference on*, pages 1–8. IEEE, 2012. (Cited on page 112.)
- H. Deng and G. Runger. Gene selection with guided regularized random forest. *Pattern Recognition*, 46(12):3483–3489, 2013. (Cited on page 112.)
- W. Deng, Z. Geng, and P. Luo. Identifiability of intermediate variables on causal paths. Frontiers of Mathematics in China, 8(3):517–539, 2013. (Cited on page 25.)
- M. Denil, D. Matheson, and N. De Freitas. Narrowing the gap: Random forests in theory and in practice. In *International conference on machine learning*, pages 665–673, 2014. (Cited on page 79.)
- P. A. Devijver and J. Kittler. *Pattern recognition: A statistical approach*. Prentice hall, 1982. (Cited on page 51.)
- R. Díaz-Uriarte and S. A. De Andres. Gene selection and classification of microarray data using random forest. *BMC bioinformatics*, 7(1):3, 2006. (Cited on pages 91, 99, and 112.)
- S. Diciotti, S. Ciulli, M. Mascalchi, M. Giannelli, and N. Toschi. The "peeking" effect in supervised feature selection on diffusion tensor imaging data. *American Journal of Neuroradiology*, 34(9):E107–E107, 2013. (Cited on page 54.)
- T. G. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine learning*, 40(2):139–157, 2000. (Cited on page 57.)
- T. G. Dietterich and E. B. Kong. Machine learning bias, statistical bias, and statistical variance of decision tree algorithms. Technical report, Technical report, Department of Computer Science, Oregon State University, 1995. (Cited on pages 71, 73, and 75.)
- A. Dobra and J. Gehrke. Bias correction in classification tree construction. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 90–97. Morgan Kaufmann Publishers Inc., 2001. (Cited on pages 103 and 106.)

- P. Domingos. Exploiting context in feature selection. In Workshop on Learning in Context-Sensitive Domains at the 13th International Conference on Machine Learning (ICML96), pages 15–20. Bari, Italy, 1996. (Cited on page 53.)
- F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608, 2017. (Cited on page 15.)
- M. Dramiński, A. Rada-Iglesias, S. Enroth, C. Wadelius, J. Koronacki, and J. Komorowski. Monte carlo feature selection for supervised classification. *Bioinformatics*, 24(1):110–117, 2008. (Cited on page 166.)
- M. Dramiński, M. J. Dabrowski, K. Diamanti, J. Koronacki, and J. Komorowski. Discovering networks of interdependent features in high-dimensional problems. In *Big Data Analysis: New Algorithms for a New Society*, pages 285–304. Springer, 2016. (Cited on page 166.)
- B. Efron and R. J. Tibshirani. *An introduction to the bootstrap.* CRC press, 1994. (Cited on page 75.)
- D. Ernst, P. Geurts, and L. Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005. (Cited on page 31.)
- R. M. Ewers and R. K. Didham. Confounding factors in the detection of species responses to habitat fragmentation. *Biological reviews*, 81(1):117–142, 2006. (Cited on page 27.)
- B. Frénay, G. Doquire, and M. Verleysen. Is mutual information adequate for feature selection in regression? *Neural Networks*, 48:1–7, 2013. (Cited on page 48.)
- J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001. (Cited on pages 88 and 94.)
- J. Gama. Functional trees. *Machine Learning*, 55(3):219–250, 2004. (Cited on page 60.)
- M. Ganz, D. N. Greve, B. Fischl, E. Konukoglu, A. D. N. Initiative, et al. Relevant feature set estimation with a knock-out strategy and random forests. *NeuroImage*, 122:131–148, 2015. (Cited on page 113.)
- H. J. Geissler, P. Hölzl, S. Marohl, F. Kuhn-Régnier, U. Mehlhorn, M. Südkamp, and E. R. de Vivie. Risk stratification in heart surgery: comparison of six score systems. *European Journal of Cardio-thoracic surgery*, 17(4):400–406, 2000. (Cited on page 146.)
- J. H. Gennari, P. Langley, and D. Fisher. Models of incremental concept formation. *Artificial intelligence*, 40(1-3):11–61, 1989. (Cited on page 33.)
- R. Genuer, J.-M. Poggi, and C. Tuleau-Malot. Variable selection using random forests. *Pattern Recognition Letters*, 31(14):2225–2236, 2010. (Cited on pages 44, 80, 90, 100, and 101.)
- P. Geurts. Contributions to decision tree induction: bias/variance tradeoff and time series classification. PhD thesis, University of Liège Belgium, 2002. (Cited on pages 58, 71, 73, 74, 77, and 79.)

- P. Geurts and Y. Saeys. Exploring signature multiplicity in microarray data using ensembles of randomized trees. In *5th International workshop on Machine Learning in Systems Biology (MLSB'11)*, pages 24–28. Technical University München, 2011. (Cited on page 36.)
- P. Geurts, D. Ernst, and L. Wehenkel. Extremely randomized trees. *Machine learning*, 63(1):3–42, 2006. (Cited on pages 50, 77, 79, 102, and 126.)
- P. Geurts, A. Irrthum, and L. Wehenkel. Supervised learning with decision tree-based methods in computational and systems biology. *Molecular Biosystems*, 5 (12):1593–1605, 2009. (Cited on page 99.)
- B. Ghimire, J. Rogan, and J. Miller. Contextual land-cover classification: incorporating spatial dependence in land-cover classification models using random forests and the getis statistic. *Remote Sensing Letters*, 1(1):45–54, 2010. (Cited on page 91.)
- C. Gini. Variabilità e mutabilità. Reprinted in Memorie di metodologica statistica (Ed. Pizetti E, Salvemini, T). Rome: Libreria Eredi Virgilio Veschi, 1912. (Cited on page 67.)
- B. Goebel, Z. Dawy, J. Hagenauer, and J. C. Mueller. An approximation to the distribution of finite sample size mutual information estimates. In *Communications*, 2005. ICC 2005. 2005 IEEE International Conference on, volume 2, pages 1102– 1106. IEEE, 2005. (Cited on pages 140 and 141.)
- T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, et al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *science*, 286(5439):531–537, 1999. (Cited on page 44.)
- B. Gregorutti, B. Michel, and P. Saint-Pierre. Correlation and variable importance in random forests. *Statistics and Computing*, 27(3):659–678, 2017. (Cited on pages 90, 91, 92, 93, 96, 97, 99, and 118.)
- U. Grömping. Variable importance assessment in regression: linear regression versus random forest. *The American Statistician*, 63(4):308–319, 2009. (Cited on page 99.)
- Guyon and A. Elisseeff. An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar):1157–1182, 2003. (Cited on pages 32, 44, 47, 48, 49, 50, 51, 52, and 53.)
- I. Guyon and A. Elisseeff. An introduction to feature extraction. In *Feature extraction*, pages 1–25. Springer, 2006. (Cited on pages 33, 35, 43, 44, 47, 50, 86, and 105.)
- A. Hapfelmeier and K. Ulm. A new variable selection approach using random forests. Computational Statistics & Data Analysis, 60:50–69, 2013. (Cited on pages 110 and 111.)
- D. Hardin, I. Tsamardinos, and C. F. Aliferis. A theoretical characterization of linear sym-based feature selection. In *Proceedings of the twenty-first international con*ference on Machine learning, page 48. ACM, 2004. (Cited on pages 36 and 37.)

- T. Hastie, R. Tibshirani, J. Friedman, and J. Franklin. The elements of statistical learning: data mining, inference and prediction. *The Mathematical Intelligencer*, 27(2):83–85, 2005. (Cited on pages 70, 73, 74, 79, and 81.)
- T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference, and prediction, 2nd Edition*, volume 1 of *Springer series in statistics*. Springer, 2009. (Cited on page 90.)
- Z. He and W. Yu. Stable feature selection for biomarker discovery. *Computational biology and chemistry*, 34(4):215–225, 2010. (Cited on pages 54 and 55.)
- D. Heath, S. Kasif, and S. Salzberg. Induction of oblique decision trees. In *IJCAI*, volume 1993, pages 1002–1007, 1993. (Cited on page 60.)
- D. Hernández-Lobato, G. MartíNez-MuñOz, and A. Suárez. How large should ensembles of classifiers be? *Pattern Recognition*, 46(5):1323–1336, 2013. (Cited on page 101.)
- T. K. Ho. The random subspace method for constructing decision forests. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(8):832–844, 1998. (Cited on pages 55, 57, 77, 112, and 166.)
- J. Hua, Z. Xiong, J. Lowey, E. Suh, and E. R. Dougherty. Optimal number of features as a function of sample size for various classification rules. *Bioinformatics*, 21(8): 1509–1515, 2004. (Cited on page 44.)
- J. Hua, W. D. Tembe, and E. R. Dougherty. Performance of feature-selection methods in the classification of high-dimension data. *Pattern Recognition*, 42(3):409–424, 2009. (Cited on page 54.)
- X. Huang, W. Pan, S. Grindle, X. Han, Y. Chen, S. J. Park, L. W. Miller, and J. Hall. A comparative study of discriminating human heart failure etiology using gene expression profiles. *BMC bioinformatics*, 6(1):205, 2005. (Cited on page 99.)
- V. A. Huynh-Thu. *Machine learning-based feature ranking: statistical interpretation and gene network inference*. PhD thesis, Université de Liège, 2012. (Cited on pages 110 and 186.)
- V. A. Huynh-Thu, L. Wehenkel, and P. Geurts. Exploiting tree-based variable importances to selectively identify relevant variables. In *JMLR: Workshop and Conference proceedings*, volume 4, pages 60–73. Microtome Publishing, 2008. (Cited on page 110.)
- V. A. Huynh-Thu, A. Irrthum, L. Wehenkel, and P. Geurts. Regulatory networks from expression data using tree-based methods. *PLoS ONE*, 5(9):e12776, 2010. (Cited on pages 54, 158, 185, 186, 187, 188, 189, and 194.)
- V. A. Huynh-Thu, Y. Saeys, L. Wehenkel, and P. Geurts. Statistical interpretation of machine learning-based feature importance scores for biomarker discovery. *Bioin-formatics*, 28(13):1766–1774, 2012. (Cited on pages 55, 101, 102, and 110.)
- T. Ideker and N. J. Krogan. Differential network biology. *Molecular systems biology*, 8(1), 2012. (Cited on page 146.)
- H. Ishwaran. Variable importance in binary regression trees and forests. *Electronic Journal of Statistics*, 1:519–537, 2007. (Cited on pages 91, 93, and 118.)

- H. Ishwaran and M. Lu. Standard errors and confidence intervals for variable importance in random forest regression, classification, and survival. *Statistics in medicine*, 2018. (Cited on page 112.)
- A. Jain and D. Zongker. Feature selection: Evaluation, application, and small sample performance. *IEEE transactions on pattern analysis and machine intelligence*, 19 (2):153–158, 1997. (Cited on pages 51 and 54.)
- A. K. Jain, R. P. Duin, and J. Mao. Statistical pattern recognition: A review. *IEEE Transactions on pattern analysis and machine intelligence*, 22(1):4–37, 2000. (Cited on page 51.)
- A. Jakulin. *Machine learning based on attribute interactions*. PhD thesis, Univerza v Ljubljani, 2005. (Cited on page 149.)
- A. Jakulin and I. Bratko. Analyzing attribute dependencies. Springer, 2003a. (Cited on page 149.)
- A. Jakulin and I. Bratko. Quantifying and visualizing attribute interactions. *arXiv* preprint cs/0308002, 2003b. (Cited on page 42.)
- A. Janecek, W. Gansterer, M. Demel, and G. Ecker. On the relationship between feature selection and classification accuracy. In *New Challenges for Feature Selection in Data Mining and Knowledge Discovery*, pages 90–105, 2008. (Cited on page 44.)
- C. Z. Janikow. Fuzzy decision trees: issues and methods. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 28(1):1–14, 1998. (Cited on page 60.)
- S. Janitza, C. Strobl, and A.-L. Boulesteix. An auc-based permutation variable importance measure for random forests. *BMC bioinformatics*, 14(1):119, 2013. (Cited on page 114.)
- S. Janitza, E. Celik, and A.-L. Boulesteix. A computationally fast variable importance test for random forests for high-dimensional data. *Advances in Data Analysis and Classification*, pages 1–31, 2015. (Cited on pages 108 and 114.)
- S.-y. Jiang and L.-x. Wang. Efficient feature selection based on correlation measure between continuous and discrete features. *Information Processing Letters*, 116 (2):203–215, 2016. (Cited on pages 54 and 106.)
- W. E. Johnson, C. Li, and A. Rabinovic. Adjusting batch effects in microarray expression data using empirical bayes methods. *Biostatistics*, 8(1):118–127, 2007. (Cited on page 158.)
- Jolliffe. Principal component analysis. Wiley Online Library, 2005. (Cited on page 200.)
- I. Jolliffe. Principal component analysis. In *International encyclopedia of statistical science*, pages 1094–1096. Springer, 2011. (Cited on page 32.)
- A. Joly. Exploiting random projections and sparsity with random forests and gradient boosting methods-Application to multi-label and multi-output learning, random forest model compression and leveraging input sparsity. PhD thesis, Université de Liège, Liège, Belgique, 2017. (Cited on pages 27, 65, and 74.)

- J. Kaiser and W. Reed. Data smoothing using low-pass digital filters. *Review of Scientific Instruments*, 48(11):1447–1457, 1977. (Cited on page 192.)
- A. Kalousis, J. Prados, and M. Hilario. Stability of feature selection algorithms: a study on high-dimensional spaces. *Knowledge and information systems*, 12(1): 95–116, 2007. (Cited on page 55.)
- F. Kamangar. Confounding variables in epidemiologic studies: basics and beyond. *Arch Iran Med*, 15(8):508–16, 2012. (Cited on page 26.)
- H. Kim and W.-Y. Loh. Classification trees with unbiased multiway splits. *Journal of the American Statistical Association*, 96(454):589–604, 2001. (Cited on page 106.)
- K. Kira and L. A. Rendell. The feature selection problem: Traditional methods and a new algorithm. In *Aaai*, volume 2, pages 129–134, 1992a. (Cited on pages 48, 49, 50, 51, and 52.)
- K. Kira and L. A. Rendell. A practical approach to feature selection. In *Machine Learning Proceedings* 1992, pages 249–256. Elsevier, 1992b. (Cited on page 49.)
- J. Kittler. Feature set search algorithms. *Pattern recognition and signal processing*, 1978. (Cited on page 51.)
- R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial intelligence*, 97(1-2):273–324, 1997. (Cited on pages 32, 33, 34, 35, 46, 47, 48, 49, 51, and 53.)
- D. Koller and M. Sahami. Toward optimal feature selection. Technical report, Stanford InfoLab, 1996. (Cited on pages 48 and 49.)
- E. Konukoglu and M. Ganz. Approximate false positive rate control in selection frequency for random forest. *arXiv preprint arXiv:1410.2838*, 2014. (Cited on pages 87, 89, 108, 110, 113, and 166.)
- L. I. Kuncheva. A stability index for feature selection. In *Artificial intelligence and applications*, pages 421–427, 2007. (Cited on pages 54 and 55.)
- L. I. Kuncheva and J. J. Rodríguez. On feature selection protocols for very low-sample-size data. *Pattern Recognition*, 81:660–673, 2018. (Cited on pages 22 and 54.)
- L. I. Kuncheva, J. J. Rodríguez, C. O. Plumpton, D. E. Linden, and S. J. Johnston. Random subspace ensembles for fmri classification. *Medical Imaging, IEEE Transactions on*, 29(2):531–542, 2010. (Cited on pages 166 and 167.)
- M. B. Kursa and W. R. Rudnicki. The all relevant feature selection using random forest. *arXiv preprint arXiv:1106.5112*, 2011. (Cited on pages 33, 34, 44, 46, 47, and 55.)
- S. W. Kwok and C. Carter. Multiple decision trees. In *Machine Intelligence and Pattern Recognition*, volume 9, pages 327–335. Elsevier, 1990. (Cited on pages 73 and 75.)

- C. Lai, M. J. Reinders, and L. Wessels. Random subspace method for multivariate feature selection. *Pattern recognition letters*, 27(10):1067–1076, 2006. (Cited on pages 112 and 166.)
- G. Langs, B. H. Menze, D. Lashkari, and P. Golland. Detecting stable distributed patterns of brain activation using gini contrast. *NeuroImage*, 56(2):497–507, 2011. (Cited on page 115.)
- P. Latinne, O. Debeir, and C. Decaestecker. Limiting the number of trees in random forests. In *International Workshop on Multiple Classifier Systems*, pages 178–187. Springer, 2001. (Cited on pages 80 and 101.)
- J. Lee Rodgers and W. A. Nicewander. Thirteen ways to look at the correlation coefficient. *The American Statistician*, 42(1):59–66, 1988. (Cited on page 43.)
- J. Li, S. Ma, T. Le, L. Liu, and J. Liu. Causal decision trees. *IEEE Transactions on Knowledge and Data Engineering*, 29(2):257–271, 2017. (Cited on page 207.)
- L. Li, B. Rakitsch, and K. Borgwardt. ccsvm: correcting support vector machines for confounding factors in biological data classification. *Bioinformatics*, 27(13): i342–i348, 2011. (Cited on pages 26 and 27.)
- A. Liaw, M. Wiener, et al. Classification and regression by randomforest. *R news*, 2 (3):18–22, 2002. (Cited on pages 81, 101, and 102.)
- M. Lichman. UCI machine learning repository, 2013. URL http://archive.ics. uci.edu/ml. (Cited on pages 157 and 182.)
- J. W. Lichtman and W. Denk. The big and the small: challenges of imaging the brains circuits. *Science*, 334(6056):618–623, 2011. (Cited on page 192.)
- Z. C. Lipton. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*, 2016. (Cited on page 15.)
- H. Liu and L. Yu. Toward integrating feature selection algorithms for classification and clustering. *IEEE Transactions on knowledge and data engineering*, 17(4): 491–502, 2005. (Cited on page 32.)
- Q. Liu and Y. Wu. Supervised learning. In *Encyclopedia of the Sciences of Learning*, pages 3243–3245. Springer, 2012. (Cited on page 27.)
- Y. Liu and H. Zhao. Variable importance-weighted random forests. *Quantitative Biology*, 5(4):338–351, 2017. (Cited on pages 112 and 113.)
- G. Louppe. *Understanding random forests: From theory to practice.* PhD thesis, Université de Liège, Liège, Belgique, 2014. (Cited on pages 39, 40, 58, 67, 73, 74, 75, 79, 81, 91, 94, 95, 97, 98, 106, 107, 108, 117, 118, 124, 125, 126, 131, 140, 185, 186, 187, and 188.)
- G. Louppe and P. Geurts. Ensembles on random patches. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 346–361. Springer, 2012. (Cited on pages 77, 108, 112, and 165.)
- G. Louppe, L. Wehenkel, A. Sutera, and P. Geurts. Understanding variable importances in forests of randomized trees. In *Advances in neural information processing systems*, pages 431–439, 2013. (Cited on pages 34, 91, 94, 95, 118, 120, 122, 128, 129, 150, 151, 156, 162, 168, and 169.)

- S. M. Lundberg and S.-I. Lee. Consistent feature attribution for tree ensembles. *arXiv preprint arXiv:1706.06060*, 2017. (Cited on pages 87, 99, and 114.)
- S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent individualized feature attribution for tree ensembles. *arXiv preprint arXiv:1802.03888*, 2018. (Cited on pages 87, 99, and 114.)
- K. L. Lunetta, L. B. Hayward, J. Segal, and P. Van Eerdewegh. Screening large-scale association study data: exploiting interactions using random forests. *BMC genetics*, 5(1):32, 2004. (Cited on page 99.)
- M. Luštrek, M. Gams, S. Martinčić-lpšić, et al. What makes classification trees comprehensible? *Expert Systems with Applications*, 62:333–346, 2016. (Cited on page 71.)
- D. Marbach, J. C. Costello, R. Küffner, N. Vega, R. J. Prill, D. M. Camacho, K. R. Allison, T. D. Consortium, M. Kellis, J. J. Collins, and G. Stolovitzky. Wisdom of crowds for robust network inference. *Nature methods*, 9(8):794–804, 2012. (Cited on pages 54, 158, 186, and 194.)
- D. Margaritis and S. Thrun. Bayesian network induction via local neighborhoods. In *Advances in neural information processing systems*, pages 505–511, 2000. (Cited on page 36.)
- T. Marill and D. Green. On the effectiveness of receptors in recognition systems. *IEEE transactions on Information Theory*, 9(1):11–17, 1963. (Cited on page 51.)
- L. Matthews, G. Gopinath, M. Gillespie, M. Caudy, D. Croft, B. de Bono, P. Garapati, J. Hemish, H. Hermjakob, B. Jassal, et al. Reactome knowledgebase of human biological pathways and processes. *Nucleic acids research*, 37(suppl 1):D619–D622, 2009. (Cited on page 158.)
- W. J. McGill. Multivariate information transmission. *Psychometrika*, 19(2):97–116, 1954. (Cited on pages 42 and 149.)
- N. Meinshausen and P. Bühlmann. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):417–473, 2010. (Cited on page 54.)
- P. E. Meyer and G. Bontempi. Information-theoretic gene selection in expression data. *Biological Knowledge Discovery Handbook: Preprocessing, Mining, and Postprocessing of Biological Data*, pages 399–420, 2013. (Cited on pages 43 and 48.)
- P. E. Meyer, C. Schretter, and G. Bontempi. Information-theoretic feature selection in microarray data using variable complementarity. *IEEE Journal of Selected Topics in Signal Processing*, 2(3):261–274, 2008. (Cited on pages 34, 35, 40, 42, and 48.)
- A. J. Miller. Subset selection in regression. number 40 in monographs on statistics and applied probability, 1990. (Cited on page 51.)
- R. Moddemeijer. On estimation of entropy and mutual information of continuous distributions. *Signal processing*, 16(3):233–248, 1989. (Cited on page 107.)

- K. Mohan, P. London, M. Fazel, D. Witten, and S.-I. Lee. Node-based learning of multiple gaussian graphical models. *The Journal of Machine Learning Research*, 15(1):445–488, 2014. (Cited on page 158.)
- A. M. Molinaro, R. Simon, and R. M. Pfeiffer. Prediction error estimation: a comparison of resampling methods. *Bioinformatics*, 21(15):3301–3307, 2005. (Cited on page 54.)
- P. Møller, L. E. Knudsen, S. Loft, and H. Wallin. The comet assay as a rapid test in biomonitoring occupational exposure to dna-damaging agents and effect of confounding factors. *Cancer Epidemiology and Prevention Biomarkers*, 9(10):1005–1015, 2000. (Cited on page 27.)
- K. V. S. Murthy and S. L. Salzberg. *On growing better decision trees from data*. PhD thesis, Citeseer, 1995a. (Cited on page 60.)
- S. Murthy and S. Salzberg. Lookahead and pathology in decision tree induction. In *IJCAI*, pages 1025–1033. Citeseer, 1995b. (Cited on page 67.)
- D. R. Nayak, R. Dash, and B. Majhi. Brain mr image classification using twodimensional discrete wavelet transform and adaboost with random forests. *Neu*rocomputing, 177:188–197, 2016. (Cited on page 91.)
- S. Nembrini, I. R. König, and M. N. Wright. The revival of the gini importance? *Bioinformatics*, 2018. (Cited on page 114.)
- T.-T. Nguyen, H. Zhao, J. Z. Huang, T. T. Nguyen, and M. J. Li. A new feature sampling method in random forests for predicting high-dimensional data. In *Advances in Knowledge Discovery and Data Mining*, pages 459–470. Springer, 2015. (Cited on page 166.)
- K. Nicodemus and J. Malley. Predictor correlation impacts machine learning algorithms: implications for genomic studies. *Bioinformatics*, 25(15):1884–1890, 2009. (Cited on pages 99, 104, and 105.)
- K. K. Nicodemus. Letter to the editor: On the stability and ranking of predictors from random forest variable importance measures. *Briefings in bioinformatics*, 12(4): 369–373, 2011. (Cited on page 106.)
- K. K. Nicodemus, J. D. Malley, C. Strobl, and A. Ziegler. The behaviour of random forest permutation-based variable importance measures under predictor correlation. *BMC bioinformatics*, 11(1):110, 2010. (Cited on pages 100 and 104.)
- R. Nilsson, J. M. Peña, J. Björkegren, and J. Tegnér. Consistent feature selection for pattern recognition in polynomial time. *The Journal of Machine Learning Re*search, 8:589–612, 2007. (Cited on pages 32, 34, 36, 37, 38, 44, 46, 47, 52, 53, and 173.)
- C. Olaru and L. Wehenkel. A complete fuzzy decision tree technique. *Fuzzy sets and systems*, 138(2):221–254, 2003. (Cited on page 60.)
- F. Olivier, A. Sutera, P. Geurts, R. Fonteneau, and D. Ernst. Phase identification of smart meters by clustering voltage measurements. In *Proceedings of the 20th Power Systems Computation Conference (PSCC 2018)*, 2018.

- A. V. Oppenheim, A. S. Willsky, and S. H. Nawab. *Signals and systems*, volume 2. Prentice-Hall Englewood Cliffs, NJ, 1983. (Cited on page 192.)
- W. Paja. A decision rule based approach to generational feature selection. In *Industrial Conference on Data Mining*, pages 230–239. Springer, 2018. (Cited on page 32.)
- B. Pakkenberg, D. Pelvig, L. Marner, M. J. Bundgaard, H. J. G. Gundersen, J. R. Nyengaard, and L. Regeur. Aging and the human neocortex. *Experimental gerontology*, 38(1):95–99, 2003. (Cited on page 190.)
- G. Panagopoulos. A review of network inference techniques for neural activation time series. *arXiv* preprint *arXiv*:1806.08212, 2018. (Cited on page 186.)
- H. Pang, A. Lin, M. Holford, B. E. Enerson, B. Lu, M. P. Lawton, E. Floyd, and
 H. Zhao. Pathway analysis using random forests classification and regression.
 Bioinformatics, 22(16):2028–2036, 2006. (Cited on page 99.)
- L. Paninski. Estimation of entropy and mutual information. *Neural computation*, 15 (6):1191–1253, 2003. (Cited on page 107.)
- D. Patterson. Molecular genetic analysis of down syndrome. *Human Genetics*, 126 (1):195–214, Jul 2009. ISSN 1432-1203. doi: 10.1007/s00439-009-0696-8. URL https://doi.org/10.1007/s00439-009-0696-8. (Cited on page 26.)
- J. Paul, M. Verleysen, and P. Dupont. The stability of feature selection and class prediction from ensemble tree classifiers. In *ESANN*, 2012. (Cited on pages 101 and 102.)
- J. Paul, M. Verleysen, and P. Dupont. Identification of statistically significant features from random forests. In ECML workshop on Solving Complex Machine Learning Problems with Ensemble Methods, pages 69–80, 2013. (Cited on page 114.)
- J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988. (Cited on pages 35, 37, and 38.)
- J. Pearl. Direct and indirect effects. In *Proceedings of the seventeenth conference on uncertainty in artificial intelligence*, pages 411–420. Morgan Kaufmann Publishers Inc., 2001. (Cited on pages 25 and 26.)
- J. Pearl. Causality. Cambridge university press, 2009a. (Cited on pages 23 and 25.)
- J. Pearl. Simpson's Paradox, Confounding, and Collapsibility, pages 173–200. Cambridge University Press, 2009b. doi: 10.1017/CBO9780511803161.008. (Cited on pages 25, 26, and 28.)
- K. Pearson. Mathematical contributions to the theory of evolution. iii. regression, heredity, and panmixia. *Philosophical Transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 187: 253–318, 1896. (Cited on page 43.)
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(Oct):2825–2830, 2011. (Cited on page 79.)

- H. Peng, F. Long, and C. Ding. Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on pattern analysis and machine intelligence*, 27(8):1226–1238, 2005. (Cited on page 48.)
- F. Pereira, T. Mitchell, and M. Botvinick. Machine learning classifiers and fmri: a tutorial overview. *Neuroimage*, 45(1):S199–S209, 2009. (Cited on page 54.)
- P. Pudil, J. Novovičová, and J. Kittler. Floating search methods in feature selection. *Pattern recognition letters*, 15(11):1119–1125, 1994. (Cited on page 51.)
- Y. Qi, Z. Bar-Joseph, and J. Klein-Seetharaman. Evaluation of different biological data and computational classification methods for use in protein interaction prediction. *Proteins: Structure, Function, and Bioinformatics*, 63(3):490–500, 2006. (Cited on page 99.)
- J. R. Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986. (Cited on page 57.)
- J. R. Quinlan. *C4. 5: programs for machine learning*. Elsevier, 2014. (Cited on page 57.)
- S. Raschka. Model evaluation, model selection, and algorithm selection in machine learning: Part ii bootstrapping and uncertainties [blog post], 2016. URL https://sebastianraschka.com/blog/2016/model-evaluation-selection-part2.html. Accessed: 28 Oct. 2018. (Cited on page 75.)
- S. J. Raudys and A. K. Jain. Small sample size effects in statistical pattern recognition: Recommendations for practitioners. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 13(3):252–264, 1991. (Cited on page 54.)
- J. Reunanen. Overfitting in making comparisons between variable selection methods. *Journal of Machine Learning Research*, 3(Mar):1371–1382, 2003. (Cited on pages 51 and 54.)
- J. Richiardi, H. Eryilmaz, S. Schwartz, P. Vuilleumier, and D. Van De Ville. Brain decoding of fmri connectivity graphs using decision tree ensembles. In *Biomedical Imaging: From Nano to Macro, 2010 IEEE International Symposium on*, pages 1137–1140. IEEE, 2010. (Cited on page 115.)
- W. Rodenburg, A. G. Heidema, J. M. Boer, I. M. Bovee-Oudenhoven, E. J. Feskens, E. C. Mariman, and J. Keijer. A framework to identify physiological responses in microarray-based gene expression studies: selection and interpretation of biologically relevant genes. *Physiological genomics*, 33(1):78–90, 2008. (Cited on pages 99, 110, and 113.)
- J. J. Rodriguez, L. I. Kuncheva, and C. J. Alonso. Rotation forest: A new classifier ensemble method. *IEEE transactions on pattern analysis and machine intelligence*, 28(10):1619–1630, 2006. (Cited on page 79.)
- M. Rohrbach, S. Ebert, and B. Schiele. Transfer learning in a transductive setting. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 46–54. Curran Associates, Inc., 2013. URL http://papers.nips.cc/

- paper/5209-transfer-learning-in-a-transductive-setting.pdf. (Cited on page 31.)
- L. Rokach. Data mining with decision trees: theory and applications. series in machine perception and artificial intelligence: Volume 69. vol. 69, 2008. (Cited on page 60.)
- W. R. Rudnicki, M. Kierczak, J. Koronacki, and J. Komorowski. A statistical method for determining importance of variables in an information system. In *International Conference on Rough Sets and Current Trends in Computing*, pages 557–566. Springer, 2006. (Cited on page 55.)
- Y. Saeys, I. Inza, and P. Larrañaga. A review of feature selection techniques in bioinformatics. *bioinformatics*, 23(19):2507–2517, 2007. (Cited on pages 32, 47, 48, 49, 50, and 54.)
- Y. Saeys, T. Abeel, and Y. de Peer. Towards robust feature selection techniques. In *Proceedings of Benelearn*, pages 45–46. Citeseer, 2008a. (Cited on page 55.)
- Y. Saeys, T. Abeel, and Y. Van de Peer. Robust feature selection using ensemble feature selection techniques. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 313–325. Springer, 2008b. (Cited on pages 54, 55, 101, and 102.)
- M. Sandri and P. Zuccolotto. A bias correction algorithm for the gini variable importance measure in classification trees. *Journal of Computational and Graphical Statistics*, 17(3):611–628, 2008. (Cited on page 114.)
- G. Saporta. *Probabilités, analyse des données et statistique*. Editions Technip, 2006. (Cited on page 141.)
- J. Schäfer and K. Strimmer. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical applications in genetics and molecular biology*, 4(32):1175, 2005. (Cited on page 193.)
- M. Schrynemackers. Supervised inference of biological networks with trees: Application to genetic interactions in yeast. PhD thesis, Université de Liège, 2015. (Cited on pages 27 and 71.)
- M. Schrynemackers, R. Küffner, and P. Geurts. On protocols and measures for the validation of supervised methods for the inference of biological networks. *Fron*tiers in genetics, 4, 2013. (Cited on pages 185, 187, and 194.)
- M. Schrynemackers, L. Wehenkel, M. M. Babu, and P. Geurts. Classifying pairs with trees for supervised biological network inference. *Molecular BioSystems*, 11(8): 2116–2125, 2015. (Cited on page 187.)
- T. Schürmann. Bias analysis in entropy estimation. *Journal of Physics A: Mathematical and General*, 37(27):L295, 2004. (Cited on page 107.)
- E. Scornet. Random forests and kernel methods. *IEEE Transactions on Information Theory*, 62(3):1485–1500, 2016. (Cited on page 81.)
- E. Scornet, G. Biau, J.-P. Vert, et al. Consistency of random forests. *The Annals of Statistics*, 43(4):1716–1741, 2015. (Cited on page 79.)

- C. E. Shannon and W. Weaver. *The Mathematical Theory of Communication*. Urbana, 1949. (Cited on page 67.)
- C. Sima and E. R. Dougherty. What should be expected from feature selection in small-sample settings. *Bioinformatics*, 22(19):2430–2436, 2006. (Cited on page 54.)
- T. J. Simons. Calcium and neuronal function. *Neurosurgical review*, 11(2):119–129, 1988. (Cited on page 186.)
- E. H. Simpson. The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 238–241, 1951. (Cited on page 28.)
- P. Smialowski, D. Frishman, and S. Kramer. Pitfalls of supervised feature selection. *Bioinformatics*, 26(3):440–443, 2009. (Cited on page 54.)
- P. Somol, P. Pudil, J. Novovičová, and P. Paclık. Adaptive floating search methods in feature selection. *Pattern recognition letters*, 20(11-13):1157–1163, 1999. (Cited on pages 51 and 52.)
- O. Sporns. Brain connectivity. *Scholarpedia*, 2(10):4695, 2007. doi: 10.4249/scholarpedia.4695. revision #91084. (Cited on page 186.)
- A. Statnikov and C. F. Aliferis. Analysis and computational dissection of molecular signature multiplicity. *PLoS computational biology*, 6(5):e1000790, 2010. (Cited on page 36.)
- A. Statnikov, L. Wang, and C. F. Aliferis. A comprehensive comparison of random forests and support vector machines for microarray-based cancer classification. *BMC bioinformatics*, 9(1):319, 2008. (Cited on page 91.)
- A. Statnikov, N. I. Lytkin, J. Lemeire, and C. F. Aliferis. Algorithms for discovery of multiple markov boundaries. *Journal of Machine Learning Research*, 14(Feb): 499–566, 2013. (Cited on pages 35, 36, 37, 38, and 47.)
- S. Stearns. On selecting features for pattern classifiers. In *Proceedings of the 3rd International Conference on Pattern Recognition (ICPR 1976)*, pages 71–75, 1976. (Cited on page 51.)
- O. Stetter, D. Battaglia, J. Soriano, and T. Geisel. Model-free reconstruction of excitatory neuronal connectivity from calcium imaging signals. *PLoS computational biology*, 8(8):e1002653, 2012. (Cited on pages 192, 193, and 194.)
- H. Stoppiglia, G. Dreyfus, R. Dubois, and Y. Oussar. Ranking a random feature for variable and feature selection. *Journal of Machine Learning Research*, 3:1399– 1414, 2003a. (Cited on page 176.)
- H. Stoppiglia, G. Dreyfus, R. Dubois, and Y. Oussar. Ranking a random feature for variable and feature selection. *Journal of machine learning research*, 3(Mar): 1399–1414, 2003b. (Cited on pages 48, 55, and 108.)
- C. Strobl and A. Zeileis. Danger: High power!—exploring the statistical properties of a test for random forest variable importance. Technical report, Department of Statistics, University of Munich, 2008. (Cited on page 109.)

- C. Strobl, A.-L. Boulesteix, and T. Augustin. Unbiased split selection for classification trees based on the gini index. *Computational Statistics & Data Analysis*, 52(1): 483–501, 2007a. (Cited on page 106.)
- C. Strobl, A.-L. Boulesteix, A. Zeileis, and T. Hothorn. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics*, 8(1):25, 2007b. (Cited on pages 87, 88, 106, 108, and 140.)
- C. Strobl, A.-L. Boulesteix, T. Kneib, T. Augustin, and A. Zeileis. Conditional variable importance for random forests. *BMC bioinformatics*, 9(1):307, 2008. (Cited on pages 99, 100, 101, 102, 104, 105, 110, 111, and 204.)
- M. Studeny. *Probabilistic conditional independence structures*. Springer Science & Business Media, 2006. (Cited on page 38.)
- A. Sutera, A. Joly, V. François-Lavet, A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In *Neural Connectomics Workshop*, pages 23–35, 2015. (Cited on page 18.)
- A. Sutera, G. Louppe, V. A. Huynh-Thu, L. Wehenkel, and P. Geurts. Context-dependent feature analysis with random forests. In *Uncertainty In Artificial Intelligence: Proceedings of the Thirty-Second Conference*, 2016. (Cited on pages 86 and 114.)
- A. Sutera, A. Joly, V. François-Lavet, Z. A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In *Neural Connectomics Challenge*, pages 23–36. Springer, 2017.
- A. Sutera, C. Châtel, G. Louppe, L. Wehenkel, and P. Geurts. Random subspace with trees for feature selection under memory constraints. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 929–937, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. URL http://proceedings.mlr.press/v84/sutera18a.html. (Cited on pages 37, 53, 86, 91, 95, 100, 113, 118, 119, and 127.)
- V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan, and B. P. Feuston. Random forest: a classification and regression tool for compound classification and qsar modeling. *Journal of chemical information and computer sciences*, 43 (6):1947–1958, 2003. (Cited on page 91.)
- R. Tang, J. P. Sinnwell, J. Li, D. N. Rider, M. de Andrade, and J. M. Biernacka. Identification of genes and haplotypes that predict rheumatoid arthritis using random forests. In *BMC proceedings*, volume 3, page S68. BioMed Central, 2009. (Cited on pages 109 and 111.)
- D. Taralla, Z. Qiu, A. Sutera, R. Fonteneau, and D. Ernst. Decision making from confidence measurement on the reward growth using supervised learning: A study intended for large-scale video games. In *Proceedings of the 8th International Conference on Agents and Artificial Intelligence (ICAART 2016)-Volume 2*, pages 264–271, 2016.
- L. Tian, S. A. Hires, T. Mao, D. Huber, M. E. Chiappe, S. H. Chalasani, L. Petreanu, J. Akerboom, S. A. McKinney, E. R. Schreiter, et al. Imaging neural activity in

- worms, flies and mice with improved gcamp calcium indicators. *Nature methods*, 6(12):875, 2009. (Cited on page 186.)
- P. Tieri, L. Farina, M. Petti, L. Astolfi, P. Paci, and F. Castiglione. Network inference and reconstruction in bioinformatics. *Network Inference and Reconstruction in Bioinformatics.*, 2016. (Cited on page 186.)
- L. Toloşi and T. Lengauer. Classification with correlated features: unreliability of feature ranking and solutions. *Bioinformatics*, 27(14):1986–1994, 2011. (Cited on pages 99, 102, 104, 105, and 113.)
- Tsamardinos and C. Aliferis. Towards principled feature selection: Relevancy, filters and wrappers. In in Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics. Citeseer, 2003. (Cited on pages 35, 36, 37, 38, 46, and 47.)
- I. Tsamardinos, C. F. Aliferis, and A. Statnikov. Time and sample efficient discovery of markov blankets and direct causal relations. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 673–678. ACM, 2003a. (Cited on pages 37 and 47.)
- I. Tsamardinos, C. F. Aliferis, A. R. Statnikov, and E. Statnikov. Algorithms for large scale markov blanket discovery. In *FLAIRS conference*, volume 2, pages 376– 380, 2003b. (Cited on page 47.)
- A. M. Turing. Computing machinery and intelligence, 1950. URL http://cogprints.org/499/. One of the most influential papers in the history of the cognitive sciences: http://cogsci.umn.edu/millennium/final.html. (Cited on page 21.)
- P. Turney. The identification of context-sensitive features: A formal definition of context for concept learning. In 13th International Conference on Machine Learning (ICML96), Workshop on Learning in Context-Sensitive Domains, pages 60–66, 1996. (Cited on page 149.)
- E. Tuv, A. Borisov, and K. Torkkola. Feature selection using ensemble based ranking against artificial contrasts. In *Neural Networks, 2006. IJCNN'06. International Joint Conference on*, pages 2181–2186. IEEE, 2006. (Cited on pages 55 and 109.)
- T. Van de Cruys. Two multivariate generalizations of pointwise mutual information. In *Proceedings of the Workshop on Distributional Semantics and Compositionality*, pages 16–20. Association for Computational Linguistics, 2011. (Cited on page 149.)
- M. J. Van der Laan. Statistical inference for variable importance. *The International Journal of Biostatistics*, 2(1), 2006. (Cited on pages 102 and 109.)
- R. G. Verhaak, K. A. Hoadley, E. Purdom, V. Wang, Y. Qi, M. D. Wilkerson, C. R. Miller, L. Ding, T. Golub, J. P. Mesirov, et al. Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in pdgfra, idh1, egfr, and nf1. *Cancer cell*, 17(1):98–110, 2010. (Cited on page 157.)
- F. Wang and J. Liang. An efficient feature selection algorithm for hybrid data. *Neurocomputing*, 193:33–41, 2016. (Cited on page 54.)

- S. Watanabe. Information theoretical analysis of multivariate correlation. *IBM Journal of research and development*, 4(1):66–82, 1960. (Cited on page 42.)
- M. Wehenkel. Characterization of neurodegenerative diseases with tree ensemble methods: the case of Alzheimer's disease. PhD thesis, Université de Liège, Liège, Belgique, 2018. (Cited on pages 32, 80, 99, 101, 110, 112, 113, 115, and 140.)
- M. Wehenkel, C. Bastin, C. Phillips, and P. Geurts. Tree ensemble methods and parcelling to identify brain areas related to alzheimerś disease. In *Pattern Recognition in Neuroimaging (PRNI), 2017 International Workshop on*, pages 1–4. IEEE, 2017. (Cited on pages 85 and 110.)
- M. Wehenkel, A. Sutera, C. Bastin, P. Geurts, and C. Phillips. Random forests based group importance scores and their statistical interpretation: application for alzheimer's disease. *Frontiers in Neuroscience Brain Imaging Methods*, 2018. (Cited on pages 39, 101, 104, and 115.)
- H. White, K. Chalak, and X. Lu. Linking granger causality and the pearl causal model with settable systems. In *NIPS Mini-Symposium on Causality in Time Series*, pages 1–29, 2011. (Cited on page 23.)
- A. W. Whitney. A direct method of nonparametric measurement selection. *IEEE Transactions on Computers*, 100(9):1100–1103, 1971. (Cited on page 51.)
- W. Wienholt and B. Sendhoff. How to determine the redundancy of noisy chaotic time series. *International Journal of Bifurcation and Chaos*, 6(01):101–117, 1996. (Cited on page 42.)
- S. Wold, A. Ruhe, H. Wold, and W. Dunn, III. The collinearity problem in linear regression. the partial least squares (pls) approach to generalized inverses. *SIAM Journal on Scientific and Statistical Computing*, 5(3):735–743, 1984. (Cited on page 32.)
- D. H. Wolpert and W. G. Macready. An efficient method to estimate bagging's generalization error. *Machine Learning*, 35(1):41–55, 1999. (Cited on page 81.)
- Z. Wu, H. Wang, M. Cao, Y. Chen, and E. P. Xing. Fair deep learning prediction for healthcare applications with confounder filtering. *arXiv* preprint arXiv:1803.07276, 2018. (Cited on page 27.)
- L. Yu and H. Liu. Efficient feature selection via analysis of relevance and redundancy. *Journal of machine learning research*, 5(Oct):1205–1224, 2004. (Cited on pages 11, 39, 40, 43, 44, and 47.)
- F. Zaklouta, B. Stanciulescu, and O. Hamdoun. Traffic sign classification using kd trees and random forests. In *Neural Networks (IJCNN), The 2011 International Joint Conference on*, pages 2151–2155. IEEE, 2011. (Cited on page 91.)
- N. L. Zhang and D. L. Poole. On the role of context-specific independence in probabilistic inference. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence, IJCAI 99, Stockholm, Sweden, July 31 August 6, 1999. 2 Volumes, 1450 pages*, pages 1288–1293, 1999. (Cited on page 150.)
- G. Zhao. *A new perspective on classification*. PhD thesis, Utah State University, Department of Mathematics and Statistics, 2000. (Cited on page 79.)

R. Zhu, D. Zeng, and M. R. Kosorok. Reinforcement learning trees. *Journal of the American Statistical Association*, 110(512):1770–1784, 2015. (Cited on pages 92, 93, and 118.)

Part IV APPENDICES



NOTATIONS AND SYMBOLS

We collect below the most important and most frequently used notation. All symbols and notations are nevertheless defined precisely in the first place when they are introduced in the main text.

In Section 2.4 and all subsequent sections, we use uppercase letters to denote both individual random variables and sets of random variables, and we reserve lower case letters to denote values of variables or configurations of subsets of variables (unless explicitly specified differently). In order to lighten the presentation, we assume that all considered random variables are discrete unless explicitly specified differently.

LIST OF NOTATIONS

\mathcal{A}	a learning algorithm
α	the percentage of memory devoted to previously found features in the sequential random subspace algorithm
α_{j}	the coefficient of \boldsymbol{X}_j in a linear combination of variables
В	a subset $B \subseteq V$ of variables
$C_p^k = \frac{p!}{k!(p-k)!}$	the number of combinations of \boldsymbol{k} elements from a set of \boldsymbol{p} elements
cov(A, B)	the covariance of A and B
c_{j}	a class
С	the number of classes
d	the size of the selected feature subset in heuristic search methods (d \leqslant p)
	the maximal depth parameter in tree-based methods
D	the maximal depth parameter in tree-based methods
$\mathbf{D} = \{\mathbf{o}^i\}_{i=1}^N$	a dataset of N observations
$\mathfrak D$	a sample of input-output pairs (x,y)
$\tilde{\mathbb{D}}_{\mathfrak{m}}$	a modified sample obtained from $\ensuremath{\mathfrak{D}}$ by permuting the values of the variable $X_{\mathfrak{m}}$ randomly
$deg(X_m)$	degree of variable $X_{\mathfrak{m}}$
$\Delta i(s,t)$	the impurity reduction of the split s at node t
$\Delta i(s^*,t)$	the impurity reduction of the best split s^{\ast} at node t
Δi_{min}	the minimal impurity reduction
Err(f)	the generalisation error of f
$Err(f_B)$	the residual error, i.e., the generalisation error of the Bayes model

$\widehat{Err}(\widehat{f}_{LS},LS)$	the training error or the empirical risk
$\widehat{\operatorname{Err}}(\widehat{f}_{LS},LS')$	the training error or the empirical risk
err (ILS, LS)	the average prediction error
	the out-of-bag error estimate
$\mathbb{E}_{X}\{X\}$	expectation value of X
$\mathbb{E}_{X}\{f(X)\}$	expectation of a function of a random variable X
$\mathbb{E}_{X,Y}\{f(X,Y)\}$	expectation of a function of random variables X and Y
$\mathbb{E}_{X Y}\{f(X,Y)\}$	expectation of a function of a random variables X given Y
$\hat{f}_{\mathbf{LS}}$	a model learnt from a learning set LS
$f(\mathbf{x})$	prediction of a model f for an input vector x
f_B	Bayes model
F	the subset of selected features
G_{T}	a tree structure
$G(\cdot)$	the impurity decrease for a generic impurity measure
H(X)	entropy of X
H(X,Y)	joint entropy of X and Y
H(X Y)	conditional entropy of X given Y
$Imp(X_m)$	the mean decrease of impurity importance of $\boldsymbol{X}_{\mathfrak{m}}$
Imp^{x_c} , $\operatorname{Imp}^{x_c}_s$, $\operatorname{Imp}^{ x_c }$	contextual mean decrease of impurity importances in the context $\boldsymbol{\kappa}_c$
$Imp^{X_c}(X_m)$	contextual mean decrease of impurity importance of $X_{\mathfrak{m}}$ given the contextual variable X_c in asymptotic conditions
$Imp^{freq}(X_m)$	the feature selection frequency importance of $X_{\mathfrak{m}}$
$Imp^{mdi}(X_m)$	the mean decrease of impurity importance of $X_{\mathfrak{m}}$
$Imp_{\infty}^{\mathfrak{m}di}(X_{\mathfrak{m}})$	the mean decrease of impurity importance of $X_{\mathfrak{m}}$ in asymptotic conditions
$Imp_{\infty}^{\mathfrak{mdi},K}(X_{\mathfrak{m}})$	the mean decrease of impurity importance of $X_{\rm m}$ in asymptotic conditions as computed by an ensemble of trees with randomisation parameter ${\sf K}$
${ m Imp}_{{ m N},{ m N}_{ m T}}^{{ m K},{ m D}}$	the mean decrease of impurity importance of $X_{\rm m}$ as computed by an ensemble of $N_{\rm T}$ trees with parameters K and D from a learning set of N samples.
$Imp_{q,\infty}^{K,\alpha}(X_m)$	the mean decrease of impurity importance of X_m as computed in asymptotic conditions in the context of the sequential random subspace algorithm with parameters α , q
$\operatorname{Imp}_{f}^{\operatorname{mda}}(X_{\operatorname{\mathfrak{m}}},f,\mathbb{D},\tilde{\mathbb{D}}_{\operatorname{\mathfrak{m}}})$	the mean decrease of accuracy estimate of $X_{\mathfrak{m}}$ in f over ${\mathfrak{D}}$ for a particular permutation ${\mathfrak{D}}_{\mathfrak{m}}$
$Imp_f^{mda}(X_m, f, D)$	the mean decrease of accuracy estimate of $X_{\mathfrak{m}}$ in f over \mathfrak{D}
$\text{Imp}_{\text{Algo}}^{\text{mda}}(X_{\text{m}}, \text{Algo}, \textbf{LS})$	the mean decrease of accuracy importance of $\ensuremath{X_{\mathfrak{m}}}$

the mean decrease of accuracy importance of $X_{\rm m}$ in $Imp_{\infty}^{mda}(X_m)$

asymptotic conditions

 $Imp_z^{mda}(X_m, Algo, LS)$ the z-score of X_m

 $Imp_f^{infl}(X_i, f)$ the relative influence of X_i in f

i(t)the impurity of node t

 $i_h(t)$ the Shannon impurity of node t

 $i_q(t)$ the Gini impurity of node t

 $i_{\nu}(t)$ the variance estimate impurity of node t

imin the minimal impurity

I(X;Y)mutual information of X and Y

 $I(X_1,\ldots,X_p;Y)$ mutual information of X_1, \ldots, X_p and Y

I(X;Y|Z)conditional mutual information of X and Y given Z

I(X;Y;Z)multivariate mutual information of X,Y,Z

I(X; Y; Z|B)multivariate mutual information of X,Y,Z given B $X \perp\!\!\!\perp Y$ X is independent of Y (the same as $X \perp \{Y\}$) $X \perp \!\!\! \perp Y|Z$ X is conditionally independent of Y given Z

 $X \not\perp\!\!\!\perp Y$ X is dependent on Y

 $X \perp \!\!\! \perp Y | Z$ X is dependent on Y given Z

the indicator function which equals 1 when its argu- $\mathbb{1}(\cdot)$

ment is true. 0 otherwise

K the number of folds

the number of input variables drawn at each node for

finding a split

the number of input variables drawn for each tree in random subspace and random patches methods

LS a learning set (of size $N \times p$)

 LS_{train} a training set LS_{test} a test set

the learning set associated to node t, i.e., the set of all LS_t

learning samples reaching node t

 LS^B a bootstrap sample set

LS_i a bootstrap sample set for T_i LSoob an out-of-bag sample set

 $LS_{i}^{oob} = LS \setminus LS_{i}^{B}$ the out-of-bag sample set for Ti

a set of elements

 $\mathcal{L}^{\mathfrak{i}}$ a subset $\mathcal{L}^{i} \subseteq \mathcal{L}$ of elements

 $L(f(\mathbf{x}), \mathbf{y})$ a loss function

the number of samples drawn for each tree in random L

subspace and random patches methods

 Σ_{ij}

L^{0-1}	the zero-one $(0-1)$ loss function
Lae	the absolute error loss function
Lse	the squared error loss function
M	a Markov boundary or Markov blanket
m	cardinality of a (m-ary) variable
μ_A	the mean of A (the same as $\mathbb{E}_A\{A\}$)
N	the number of samples or observations
N_t	the number of samples reaching node t
N_{T}	the number of trees in an ensemble of trees or forest
n_{min}	the minimal number of samples required to split a node
n_{leaf}	the minimal number of samples required in child nodes after the split
N_{nodes}	the maximal number of nodes
N_{leaf}	the maximal number of leaves
$\mathbf{o} = (o_1, \dots, o_p)$	an observation, sample, example
$\mathbf{o}^{\mathfrak{i}}=(\mathfrak{o}_{1}^{\mathfrak{i}},\ldots,\mathfrak{o}_{\mathfrak{p}}^{\mathfrak{i}})$	the i^{th} observation in ${f D}$
p	the number of features
$\mathfrak{p}_{\mathfrak{j},\mathfrak{i}}$	the confidence level for the putative edge from \boldsymbol{X}_j to \boldsymbol{X}_i
	partial correlation between X_j and X_i
$p(t) = \frac{N_t}{N}$	the ratio of samples reaching node t
$p(c_j t)$	the proportion of samples in $\boldsymbol{L}\boldsymbol{S}_t$ such that $\boldsymbol{y}=\boldsymbol{c}_j$
$P_{X,Y,Z}$	joint probability density of variables X, Y, Z
$P_{X,Y,Z}(x,y,z)$	the value of the joint probability density $P_{X,Y,Z}$ for a combination of values of variables X,Y,Z
$P_{X,Y Z}$	conditional joint probability density of \boldsymbol{X} and \boldsymbol{Y} given \boldsymbol{Z}
$P_{X,Y Z}(x,y z)$	the value of the conditional joint density $P_{X,Y\mid Z}$ for a combination of values of variables X,Y,Z
φ	a partitioning of LS provided by TLS
ϕ^*	the optimal partitioning of LS
$\mathcal{P}_k(V^{-m})$	the set of all subsets of cardinality k of $V^{-\mathfrak{m}}$
q	the number of features in a subspace
r	the number of relevant features
ρ	the Pearson correlation coefficient
$\rho(A, B)$	Pearson correlation between A and B
σ_{A}	the standard deviation of A
Σ	the covariance matrix
Σ^{-1}	the precision or concentration matrix
7	

the element $(\mathfrak{i},\mathfrak{j})$ of the covariance matrix Σ

S	a split
s _t	the split associated
s*	the best split st in S

the cardinality of a split, i.e., the number of created $|s_t|$

subsets or the number of children of node t

to node t

the set of all candidate splitting function for node t (on S_{t}

any feature)

the set of all candidate splitting function for node t on $S_{t,m}$

feature X_m

t a node in a decision tree

 t_0 the root node

the left child of a node t in a binary decision tree t_L t_R the right child of a node t in a binary decision tree the successor node of t corresponding to value x_m of t_{x_m}

 $X_{\mathfrak{m}}$

Τ a decision tree model

the number of iterations in the sequential random sub-

space algorithm

TLS a decision tree model learnt on LS

 T^* the best subtree $T^* \subseteq T$

a random forest model made of a set of N_T different $T = \{T_i | i = 1, ..., N_T\}$

trees T_i

cut-point, split value, or threshold value of a split

 $V = \{X_1, \dots, X_p\}$ set of all input features

 $V^{-i} = V \setminus \{X_i\}$ set of all input features V without X_i

cardinality of a set of variables, i.e., the number of vari-|V|

ables in V

 $v(s_t)$ a split variable, i.e., the variable used for the split st

 $var\{Y|B=b\}$ empirical variance of Y given B = b

Χ an input feature or variable

 χ_{i} the ith input feature or variable (of V)

 X_{c} a context variable

cardinality of a variable, i.e., the number of possible |X|

values for X

 χ the input space

 χ_{i} an input subspace (i.e., $\mathcal{X}_i \subseteq \mathcal{X}$)

 χ_{t} the input subspace associated to node t

 χ^{s} the part of the input subspace that satisfies the test s

the part of the input subspace that does not satisfy the $\mathfrak{X}^{\bar{\mathbf{s}}} = \mathfrak{X} \setminus \mathfrak{X}^{\mathbf{s}}$

test s

$\mathbf{x} = \{x_1, \dots, x_p\}$	a value of the vector of input variables			
x^{i}	\mathfrak{i}^{th} sample of a learning set LS			
Υ	output feature, target variable			
y	output space			
y	a value of the output variable Y			
y ⁱ	the value of variable Y for the \mathfrak{i}^{th} sample			
ŷ	approximated value of y			
ŷt	the value associated to node t			

LIST OF SYMBOLS

union \bigcup

intersection

difference

logical not

logical exclusive-or (xor) \oplus

estimation, approximation of a quantity

independence \perp

dependence 1

1 indicator function

NOTATIONS, AND DEFINITIONS OF ENTROPIES AND MUTUAL INFORMATION

To be self-contained, we first recall several definitions from information theory (see Cover and Thomas [2012], for further properties).

We suppose that we are given a probability space $(\Omega, \mathcal{E}, \mathbb{P})$ and consider random variables defined on it taking a finite number of possible values. We use upper case letters to denote such random variables (e.g. $X, Y, Z, W \ldots$) and calligraphic letters (e.g. $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{W} \ldots$) to denote their image sets (of finite cardinality), and lower case letters (e.g. $x, y, z, w \ldots$) to denote one of their possible values. For a (finite) set of (finite) random variables $X = \{X_1, \ldots, X_i\}$, we denote by $P_X(x) = P_X(x_1, \ldots, x_i)$ the probability $\mathbb{P}(\{\omega \in \Omega \mid \forall \ell : 1, \ldots, i : X_\ell(\omega) = x_\ell\})$, and by $\mathcal{X} = \mathcal{X}_1 \times \cdots \times \mathcal{X}_i$ the set of joint configurations of these random variables. Given two sets of random variables, $X = \{X_1, \ldots, X_i\}$ and $Y = \{Y_1, \ldots, Y_j\}$, we denote by $P_{X|Y}(x \mid y) = P_{X,Y}(x,y)/P_Y(y)$ the conditional density of X with respect to Y.

With these notations, the joint (Shannon) entropy of a set of random variables $X = \{X_1, \dots, X_i\}$ is thus defined by

$$H(X) = -\sum_{x \in \mathcal{X}} P_X(x) \log_2 P_X(x),$$

while the mean conditional entropy of a set of random variables $X = \{X_1, ..., X_i\}$, given the values of another set of random variables $Y = \{Y_1, ..., Y_i\}$ is defined by

$$H(X \mid Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{X,Y}(x,y) \log_2 P_{X|Y}(x \mid y).$$

The mutual information among the set of random variables $X = \{X_1, ..., X_i\}$ and the set of random variables $Y = \{Y_1, ..., Y_i\}$ is defined by

$$I(X;Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{X,Y}(x,y) \log_2 \frac{P_X(x)P_Y(y)}{P_{X,Y}(x,y)}$$

= $H(X) - H(X \mid Y)$
= $H(Y) - H(Y \mid X)$.

The mean conditional mutual information among the set of random variables $X = \{X_1, \ldots, X_k\}$ and the set of random variables $Y = \{Y_1, \ldots, Y_j\}$, given the values of a third set of random variables $Z = \{Z_1, \ldots, Z_i\}$, is defined by

$$\begin{split} \mathrm{I}(X;Y \mid Z) &= \mathrm{H}(X \mid Z) - \mathrm{H}(X \mid Y, Z) \\ &= \mathrm{H}(Y \mid Z) - \mathrm{H}(Y \mid X, Z) \\ &= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} \mathrm{P}_{X,Y,Z}(x,y,z) \log_2 \frac{\mathrm{P}_{X\mid Z}(x \mid z) \mathrm{P}_{Y\mid Z}(y \mid z)}{\mathrm{P}_{X,Y\mid Z}(x,y \mid z)}. \end{split}$$

We also recall the chaining rule

$$I(X, Z; Y | W) = I(X; Y | W) + I(Z; Y | W, X),$$

¹To avoid problems, we suppose that all probabilities are strictly positive, without fundamental limitation.

and the symmetry of the (conditional) mutual information among sets of random variables

$$I(X;Y\mid Z) = I(Y;X\mid Z).$$

DIGIT RECOGNITION PROBLEM

The problem of digit recognition was introduced in [Breiman et al., 1984] and is used in several occasions in this thesis for illustrating variable importances computed from tree-based methods.

It models a seven-segment display displaying numerals using horizontal and vertical lights in on-off combinations, as illustrated in Figure C.0.1.

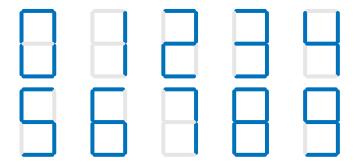


Figure C.0.1: Numerals as represented by a 7-segment display.

Variables of this problem are defined as follows: Let Y be a random variable taking its value in $\{0,1,\ldots,9\}$ with equal probability and let X_1,\ldots,X_7 be binary variables, each representing the on-off state of one segment as shown in Figure C.0.2, whose values are each determined univocally given the corresponding value of Y in Table C.0.1.

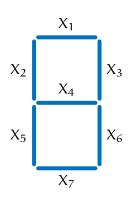


Figure C.0.2: Correspondence between segments and input variables.

y	x_1	x_2	χ_3	χ_4	χ_5	χ_6	x ₇
0	1	1	1	0	1	1	1
1	0	0	1	0	0	1	0
2	1	0	1	1	1	0	1
3	1	0	1	1	0	1	1
4	0	1	1	1	0	1	0
5	1	1	0	1	0	1	1
6	1	1	0	1	1	1	1
7	1	0	1	0	0	1	0
8	1	1	1	1	1	1	1
9	1	1	1	1	0	1	1

Table C.0.1: Values of $Y, X_1, ..., X_7$.