

STATISTICAL INFERENCE FOR MACHINE LEARNING: FEATURE IMPORTANCE, UNCERTAINTY QUANTIFICATION AND INTERPRETATION STABILITY

A Dissertation

Presented to the Faculty of the Graduate School
of Cornell University

in Partial Fulfillment of the Requirements for the Degree of
Doctor of Philosophy

by

Zhengze Zhou

May 2021

PREVIEW

© 2021 Zhengze Zhou
ALL RIGHTS RESERVED

STATISTICAL INFERENCE FOR MACHINE LEARNING: FEATURE
IMPORTANCE, UNCERTAINTY QUANTIFICATION AND
INTERPRETATION STABILITY

Zhengze Zhou, Ph.D.

Cornell University 2021

Machine learning has become ubiquitous in many areas, including high-stake applications such as autonomous driving, financial forecasting and clinical decisions. However, many models are complex in nature and act as “black boxes”, providing predictions but little insight as to how they were arrived at. In this thesis, we present our work from three different perspectives towards a better understanding of several types of machine learning models through the lens of statistical inference.

Tree-based methods, including decision trees, random forests and gradient boosting machines, are a popular class of nonparametric statistical model. They are widely used owing to their flexibility and superior performances. Many practitioners rely on some kind of feature importance measurements to examine model behavior. We propose a modification that corrects for split-improvement variable importance measures in random forests and other tree-based methods. These measurements have been shown to be biased towards increasing the importance of features with more potential splits. We show that by appropriately incorporating split-improvement as measured on out-of-sample data, this bias can be corrected yielding better summaries and screening tools.

Our next study focuses on understanding statistical properties and quantifying uncertainty for ensemble models. Tree-based ensembles like random

forests remain one such popular option for which several important theoretical advances have been made in recent years by drawing upon a connection between their natural subsampled structure and the classical theory of U -statistics. Unfortunately, the procedures for estimating predictive variance resulting from these studies are plagued by severe bias and extreme computational overhead. Here, we argue that the root of these problems lies in the structure of the resamples themselves. We develop a general framework for analyzing the asymptotic behavior of V -statistics, demonstrating asymptotic normality under precise regularity conditions and establishing previously unreported connections to U -statistics. Importantly, these findings allow us to produce a natural and efficient means of estimating the variance of a conditional expectation, a problem of wide interest across multiple scientific domains that also lies at the heart of uncertainty quantification for supervised learning ensembles. As an application, we apply this result to design a stopping rule for determining the ideal tree depth in model distillation.

Lastly, we investigate the stability for model explanation. Post hoc explanations based on perturbations are widely used approaches to interpret a machine learning model after it has been built. This class of methods has been shown to exhibit large instability, posing serious challenges to the effectiveness of the method itself and harming user trust. We propose a new algorithm called S-LIME, which utilizes a hypothesis testing framework based on central limit theorem for determining the number of perturbation points needed to guarantee stability of the resulting explanation. Experiments on both simulated and real world data sets are provided to demonstrate the effectiveness of our method.

BIOGRAPHICAL SKETCH

Zhengze Zhou was born in Zhoushan, an urbanized archipelago in eastern China, where he spent his childhood with his parents and grandmother. After one semester in the 7th grade in 2007, he transferred to Zhenhai Jiaochuan Middle School in Ningbo and then attended Zhenhai High School for three years until 2012. He got admitted to School of Mathematical Sciences at Peking University thereafter.

He was immersed in the beauty and difficulty of math for two years before he declared a major in Statistics and Probability at the end of sophomore year. He was fortunate to receive May 4th Scholarship and Community Contribution Award during undergraduate study. He graduated in 2016 as an Outstanding Graduate, and his thesis work was supervised by Prof. Jinzhu Jia.

Zhengze joined the PhD program at Cornell University in Fall 2016 and moved to Ithaca, a gorgeous yet remote city on Cayuga Lake in New York. His general interests lie at the intersection of machine learning and statistical inference. His research work was advised by Prof. Giles Hooker. During five years at Cornell, he enjoyed the never-ending snow season while also developed professional skills through research, teaching and internship. He defended his thesis in April 2021.

Dedicated to my parents, Huizhong Zhou and Jiufei Shi.

ACKNOWLEDGEMENTS

This thesis would not be possible without the help and support from many people along the journey.

First and foremost, I would like to thank my advisor Giles Hooker. Our first connection dated back to February 2016 when he extended the PhD offer to me, and he was pretty helpful even before I came to Cornell. Since I started to work with him in the summer of 2017, he has been constantly available for advice and encouragement. I was impressed by his broad knowledge and deep insights as a statistician. To me, Giles has been and will always be a great mentor and friend.

I would also like to thank my committee members Madeleine Udell and Kilian Weinberger. Madeleine's group meeting has been a wonderful venue for me to gain new knowledge and seek valuable feedback. I also enjoyed a lot while taking Kilian's course on machine learning, and his high standard for presentation quality has inevitably made me better.

Furthermore, I want to express my gratitude to collaborators: Yichen Zhou, Lucas Mentch and Fei Wang; to my industry mentors: Anup Kotalwar, Kai Yuan and Dingxiong Deng; to an awesome statistics PhD cohort, and other faculty members and personnel in the department.

Finally, I am grateful for the support from my parents, family members and all my friends.

TABLE OF CONTENTS

| | |
|---|-----------|
| Biographical Sketch | iii |
| Dedication | iv |
| Acknowledgements | v |
| Table of Contents | vi |
| List of Tables | ix |
| List of Figures | x |
| 1 Introduction | 1 |
| 1.1 A Primer on Machine Learning and Statistical Inference | 1 |
| 1.2 Roadmap | 4 |
| 2 Unbiased Measurement of Feature Importance in Tree-Based Methods | 7 |
| 2.1 Introduction | 7 |
| 2.2 Tree-Based Methods | 10 |
| 2.2.1 Tree Building Process | 10 |
| 2.2.2 Random Forests and Gradient Boosting Trees | 12 |
| 2.3 Measurement of Feature Importance | 14 |
| 2.3.1 Permutation Importance | 14 |
| 2.3.2 Split-Improvement | 15 |
| 2.3.3 Bias in Split-Improvement | 16 |
| 2.3.4 Related Work | 19 |
| 2.4 Unbiased Split-Improvement | 21 |
| 2.4.1 Classification | 22 |
| 2.4.2 Regression | 25 |
| 2.5 Empirical Studies | 27 |
| 2.5.1 Simulated Data | 27 |
| 2.5.2 RNA Sequence Data | 29 |
| 2.5.3 Adult Data | 30 |
| 2.5.4 Boston Housing Data | 32 |
| 2.5.5 Summary | 34 |
| 2.6 Discussions | 35 |
| 3 V-statistics and Variance Estimation | 37 |
| 3.1 Introduction | 37 |
| 3.2 Related Work on U -statistics | 40 |
| 3.3 V -statistics | 44 |
| 3.3.1 Asymptotic Equivalence to U -statistics | 45 |
| 3.3.2 Representation As U -statistics | 47 |
| 3.4 Variance Estimation | 49 |
| 3.4.1 Internal Variance Estimation Method | 51 |
| 3.4.2 Balanced Variance Estimation Method | 52 |
| 3.4.3 Infinitesimal Jackknife | 53 |

| | | |
|----------|---|------------|
| 3.5 | Bias Corrections for Variance Estimates | 55 |
| 3.5.1 | Bias in Variance Estimation | 55 |
| 3.5.2 | A Bias-corrected Estimator | 58 |
| 3.6 | Randomized Ensembles | 60 |
| 3.7 | Empirical Studies | 62 |
| 3.7.1 | Predictive Performance | 63 |
| 3.7.2 | Asymptotic Normality and Variance Estimation | 66 |
| 3.8 | An Application: Determining Appropriate Tree Depth in Model Distillation | 70 |
| 3.8.1 | Motivation | 70 |
| 3.8.2 | Testing Procedures | 72 |
| 3.8.3 | Empirical Results | 74 |
| 3.9 | Conclusion | 77 |
| 4 | S-LIME: Stabilized-LIME for Model Explanation | 79 |
| 4.1 | Introduction | 79 |
| 4.2 | Background | 82 |
| 4.2.1 | LIME | 82 |
| 4.2.2 | LASSO and LARS | 84 |
| 4.2.3 | Instability with LIME | 85 |
| 4.3 | Asymptotic Properties of LARS Decisions | 87 |
| 4.4 | Stabilized-LIME | 89 |
| 4.5 | Empirical Studies | 92 |
| 4.5.1 | Breast Cancer Data | 92 |
| 4.5.2 | MARS Test Function | 95 |
| 4.5.3 | Early Prediction of Sepsis From Electronic Health Records | 96 |
| 4.6 | Discussions | 99 |
| 5 | Concluding Remarks | 102 |
| A | Appendix of Chapter 2 | 105 |
| A.1 | Proofs of Lemma 1 and 3 | 105 |
| A.2 | Additional Simulation Results | 108 |
| B | Appendix of Chapter 3 | 111 |
| B.1 | Proofs | 111 |
| B.1.1 | Proof of Theorem 6 | 111 |
| B.1.2 | Proof of Theorem 7 | 113 |
| B.1.3 | Proof of Theorem 8 | 116 |
| B.1.4 | Proof of Theorem 10 | 116 |
| B.1.5 | Proof of Theorem 12 | 118 |
| B.2 | Additional Simulation Results | 119 |
| B.3 | Derivations on Bias-corrected Estimator | 121 |
| B.4 | An Alternative Version of Bias Correction | 123 |

| | | |
|----------|--|------------|
| B.5 | Bias Correction for U -statistics | 125 |
| B.6 | A Closer Look at Variance Components | 127 |
| B.7 | Additional Simulation Results on Normality for Ensembles | 130 |
| B.8 | Variance Estimation for V -statistics and Its Implications | 131 |
| B.9 | Datasets Information | 132 |
| B.10 | Approximation Tree for Breast Cancer Data | 133 |
| C | Appendix of Chapter 4 | 135 |
| C.1 | Instability with LASSO | 135 |
| C.2 | Additional Experiments | 136 |
| C.3 | Variables list for Sepsis detection | 139 |

PREVIEW

LIST OF TABLES

| | | |
|-----|---|-----|
| 2.1 | Average importance ranking of informative feature X_1 . R stands for regression and C for classification. The result averages over 100 repetitions. Lower values indicate better abilities in identifying informative features. In cforest, we set mincriterion to be 2.33 (0.99 percentile of normal distribution) for shallow trees and 1.28 (0.9 percentile) for deep trees. | 28 |
| 2.2 | Attribute description for adult data set. | 32 |
| 3.1 | Predictive performance on six datasets under different sampling strategies | 65 |
| 3.2 | Asymptotic normality and variance estimation results for MARS function | 68 |
| 3.3 | Asymptotic normality and variance estimation results for Protein Tertiary Structure across 20 test samples. | 69 |
| 4.1 | Average Jaccard index for 20 repetitions for LIME and S-LIME. The black box model is a random forests with 500 trees. | 94 |
| 4.2 | Average Jaccard index for 20 repetitions for LIME and S-LIME. The black box model is MARS. | 96 |
| 4.3 | Average Jaccard index for 20 repetitions for LIME and S-LIME on two randomly selected test samples. The black box model is a recurrent neural network. | 98 |
| B.1 | Normality Test for Ensembles. | 131 |
| B.2 | Feature names correspondence for Breast Cancer Data. | 133 |
| B.3 | Details of the approximation tree for Breast Cancer Data. Node column numbers the nodes from left to right in each layer in a complete binary tree. | 134 |
| C.1 | Average Jaccard index for 20 repetitions for LIME and S-LIME. The black box models are SVM and NN. | 138 |
| C.2 | Variables list for Sepsis detection. | 140 |

LIST OF FIGURES

| | | |
|-----|--|----|
| 2.1 | Split-improvement measures on five predictors. Box plot is based on 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. | 17 |
| 2.2 | Average feature importance ranking across different signal strengths over 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. | 18 |
| 2.3 | Unbiased split-improvement. Box plot is based on 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. Each tree is trained using bootstrap samples and <i>out-of-bag</i> samples are used as test set. | 22 |
| 2.4 | Unbiased feature importance ranking across different signal strengths averaged over 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. Each tree is trained using bootstrap samples and <i>out-of-bag</i> samples are used as test set. | 23 |
| 2.5 | Feature importance for RNA sequence data. 100 trees are built in the forest. Red error bars depict one standard deviation when the experiments are repeated 100 times. The x-axis denotes two classes of features: position relative to edited/non-edited site, other measures (cp, fe, dfe). | 31 |
| 2.6 | Feature importance for adult data. 20 trees are built in the forest. Red error bars depict one standard deviation when the experiments are repeated 100 times. The x-axis lists feature names for the model: Age, fnlgwt, Education, Sex, Capital Gain, Capital Loss, Hours/Week, Workclass, Education, Marital Status, Occupation, Relationship, Race, random. | 33 |
| 2.7 | Feature importance for Boston housing data. 100 trees are built in the forest. Red error bars depict one standard deviation when the experiments are repeated 100 times. The x-axis lists feature names for the model: CRIM, ZN, INDUS, CHAS, NOX, RM, AGE, DIS, RAD, TAX, PTRATIO, B, LSTAT, random. | 34 |
| 3.1 | Variance estimation by three different methods: Internal Variance Estimation Method (IM), Balanced Variance Estimation Method (BM) and Infinitesimal Jackknife (IJ). The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. For IM, we choose $n_{\text{OUT}} = (10, 20, 50)$ for $n_{\text{estimators}} = (100, 1000, 10000)$ respectively. | 57 |

| | | |
|-----|--|-----|
| 3.2 | Variance estimation by three different methods: corrected-V, BM and IJ. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. | 59 |
| 3.3 | Approximation tree built with cp (complexity parameter) = 0.001. Data set is generated by $X_1 \sim U[0, 1]$ and $X_2 \sim U[0, 1]$, independent of each other. $Y + 1 \sim N(0, 1)$ when $X_1 < 0.5$ and $Y - 1 \sim N(0, 1)$ when $X_1 \geq 0.5$. The number of trees $B = 10,000$ | 75 |
| 3.4 | Histograms for p -values based of 500 simulated data set with a three-node tree as an underlying generator. Left plot gives p -values at the root node, with power 0.994 to detect a real split. Right gives p -values at the children where no further splitting is warranted. The non-uniform distribution of the right-hand plot is due to averaging p -values over 10 sets of 10 query points. | 76 |
| 3.5 | A mock adaptive screening tool extracted from the dominant tree structure for the Breast Cancer Wisconsin Data Set, truncated at most 3 questions. | 77 |
| 4.1 | Empirical selection probability for features in Breast Cancer Data. The black box model is a random forests classifier with 500 trees. LIME is run 100 times on a randomly selected test point and top 5 features are selected via LASSO. | 86 |
| 4.2 | Four iterations of LIME on Breast Cancer Data. The black box model is a random forests classifier with 500 trees. LIME explanations are generated with 1000 synthetic perturbations. | 93 |
| 4.3 | Two iterations of S-LIME on Breast Cancer Data. The black box model is a random forests classifier with 500 trees. | 95 |
| 4.4 | Output of S-LIME for two randomly selected test samples. The black box model is a recurrent neural network. | 99 |
| A.1 | Split-improvement measures on five predictors, where we treat categorical features as ordered discrete values. Box plot is based on 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. | 109 |
| A.2 | Average feature importance ranking across different signal strengths over 100 repetitions, where we treat categorical features as ordered discrete values. 100 trees are built in the forest and maximum depth of each tree is set to 5. | 109 |
| A.3 | Unbiased split-improvement, where we treat categorical features as ordered discrete values. Box plot is based on 100 repetitions. 100 trees are built in the forest and maximum depth of each tree is set to 5. Each tree is trained using bootstrap samples and <i>out-of-bag</i> samples are used as test set. | 110 |

| | | |
|-----|--|-----|
| A.4 | Unbiased feature importance ranking across different signal strengths averaged over 100 repetitions, where we treat categorical features as ordered discrete values. 100 trees are built in the forest and maximum depth of each tree is set to 5. Each tree is trained using bootstrap samples and <i>out-of-bag</i> samples are used as test set. | 110 |
| B.1 | Variance estimation by three different methods: IM, BM and IJ. The kernel size $k_n = 250$. The variance shown is for prediction at test point $x = 10$. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. | 119 |
| B.2 | Variance estimation by three different methods: IM, BM and IJ. The kernel size $k_n = 400$. The variance shown is for prediction at test point $x = 10$. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. | 120 |
| B.3 | ζ_{1,k_n} estimated by three different methods: IM, BM and IJ. The kernel size $k_n = 100$. The variance shown is for prediction at test point $x = 10$ | 121 |
| B.4 | ζ_{k_n,k_n} estimated by three different methods: IM, BM and IJ. The kernel size $k_n = 100$. The variance shown is for prediction at test point $x = 10$ | 122 |
| B.5 | Variance Estimation by three different methods: corrected-IJ, BM and IJ. The kernel size $k_n = 100, 250, 400$. The variance shown is for prediction at test point $x = 10$. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. . . . | 126 |
| B.6 | Variance estimation by three different methods: corrected-V, BM and IJ. The variance shown is for prediction at test point $x = 10$. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. | 127 |
| B.7 | Variance Estimation by three different methods: corrected-U, BM and IJ. The kernel size $k_n = 100, 250, 400$. The variance shown is for prediction at test point $x = 10$. The red line denotes true (log) variance obtained by generating data, training the ensemble 100 times and calculating the empirical variance of predictions. | 128 |
| B.8 | Variance components for different B_n . The number of training observations $n = 1000$ and kernel size $k_n = 10$. The variance shown is for prediction at test point $x = 10$. Four lines shown are: empirical variance, two variance components ($\frac{k_n^2}{n}\zeta_{1,k_n}$ and $\frac{1}{B_n}\zeta_{k_n,k_n}$) and their sum as estimated variance. | 130 |

| | | |
|-----|--|-----|
| C.1 | Two cases of variable ordering in LASSO path. | 136 |
| C.2 | S-LIME on Breast Cancer Data with SVM and NN as black box models. | 139 |

PREVIEW

CHAPTER 1

INTRODUCTION

1.1 A Primer on Machine Learning and Statistical Inference

Machine learning has made its appearance in everyday life: the spam filter in your email box, the search results returned based on your Google query, a virtual assistant like Siri or Alexa, just to name a few. In essence, machine learning is the practice of using algorithms to learn through experiences. Rather than explicitly tell a program what to do, you teach the computer to develop an algorithm based on previous experience (i.e., data) to improve itself and complete the task.

There are three major types of machine learning:

- **Supervised learning** is to learn a function which maps from input to output from labeled data. The two main types of supervised learning are classification and regression. Spam filter is an example of the classification task, where an email needs to be classified as spam or not-spam. Regression involves continuous outputs such as predicting housing prices based on a few attributes.
- **Unsupervised learning** is a type of algorithm which learns patterns from unlabeled data. Two of the main methods are clustering and principal component analysis. The goal of clustering algorithms is to group or segment datasets based on similar properties. Principal component analysis is the process of computing the principal components and thus is usually used as a means for dimensional reduction.

- **Reinforcement learning** is concerned with how an agent interacts with the environment, collects feedback, and improves its own decision making process. Unlike supervised and unsupervised learning, reinforcement learning does not aim to find a “correct” output for a given instance. Instead the focus is to improve how the agent interacts with environment such that the rewards it receives in the long run can be maximized.

In this thesis we focus on supervised learning. Mathematically, suppose one has a set of n training examples denoted by $\mathcal{D} = \{z_1 = (\mathbf{x}_1, y_1), z_2 = (\mathbf{x}_2, y_2), \dots, z_n = (\mathbf{x}_n, y_n)\}$ where \mathbf{x}_i is the feature vector of i -th sample and y_i is the corresponding label, a supervised learning algorithm seeks a function $g : X \rightarrow Y$ which maps from the input space X to the output space Y . The function g is chosen from some space of possible functions denoted by \mathcal{G} , usually called hypothesis space.

For example in the spam filtering case, \mathbf{x}_i denotes some features related to an email such as bag-of-words, and $y_i \in \{0, 1\}$ denoting whether this email is spam or not. To predict housing prices, \mathbf{x}_i is a vector collecting information like square feet and number of bedrooms, $y_i \in \mathcal{R}$ being the price of the house. If we choose \mathcal{G} to be the space of linear functions, then we’ll learn a linear relationship between \mathbf{x}_i and y_i ; on the other hand, if we use more complex models such as random forests, then the hypothesis space \mathcal{G} contains any function representable by an ensemble of decision trees. We’ll provide more details in each chapter when it deals with certain model classes.

Now let’s briefly introduce the concept of statistical inference. We’ll build on the notations we just introduced. Assume all data in \mathcal{D} comes from some *unknown* distribution \mathcal{F}_Z : $z_1, z_2, \dots, z_n \stackrel{i.i.d.}{\sim} \mathcal{F}_Z$. A typical question we might ask

is:

Given a sample z_1, z_2, \dots, z_n drawn from an unknown distribution \mathcal{F}_Z , how to estimate \mathcal{F}_Z , or some quantities related to \mathcal{F}_Z ?

In a more general sense, z_i can represent a pair of values (x_i, y_i) as in the supervised learning case, or just a scalar value.

There are three types of common statistical inference problems.

- **Point estimation** aims to find a particular quantity of interest. We can write the estimate as $\hat{\theta}_n = g(z_1, z_2, \dots, z_n)$ for some function g . Here the true unknown quantity θ can be population mean, or parameters related to the underlying distribution \mathcal{F}_Z .
- **Confidence interval**. Instead of producing a point estimate, one may need a range of plausible values which contains the actual value with high probability. Given a confidence level α , a $1 - \alpha$ confidence interval for θ is an interval (l_n, u_n) such that $P(\theta \in (l_n, u_n)) \geq 1 - \alpha$.
- **Hypothesis testing** is the process of distinguishing between the null hypothesis H_0 and the alternative hypothesis H_1 . Practically, a test statistic T is calculated based on the samples z_1, z_2, \dots, z_n and if T exceeds some threshold we will conclude to reject the null hypothesis. Statistical hypothesis testing is widely used in survey or experiment studies.

For a supervised learning problem, usually we can view it as a point estimation task. If the hypothesis space \mathcal{G} can be parameterized by θ (such as linear models or neural networks), fitting a model is equivalent to finding a point estimate $\hat{\theta}$. For nonparametric case (such as tree ensembles), the model prediction

of a given instance x is the point of interest, where the true unknown quantity is $\theta = E_F[y|x]$.

At this point, it is worth mentioning one important philosophical difference between machine learning and statistical inference. Machine learning tries to make the most accurate predictions. Statistical inference is more about finding the relationship between variables, which also takes predictions into account by studying the statistical properties therein.

1.2 Roadmap

The topics in this thesis lie at the intersection of machine learning and statistical inference. Starting from some practical machine learning problems, we try to understand its behavior from the perspective of statistical inference. This section provides a roadmap where we introduce the motivation behind our work and present a high level summary of each chapter that follows.

Tree-based methods, such as random forests or gradient boosting decision trees, are arguably the most popular machine learning algorithms deployed in industry. However, large ensembles of trees act as “black boxes”, providing predictions but little insight as to how they were arrived at. One of the approaches taken by practitioners is to look at variable importance scores. Previous studies have shown that split-improvement variable importance measures are biased towards increasing the importance of features with more potential splits. In Chapter 2, we show that by appropriately incorporating split-improvement as measured on out-of-sample data, this bias can be corrected yielding better summaries and screening tools. We also prove the proposed method is statistically

unbiased for features uncorrelated with the response.

In Chapter 3, we shift our attention to the predictions generated by ensembles, where the quantity of interest is $\theta = E_F[y|x]$. Machine learning aims to build algorithms which can get better estimate $\hat{\theta}$ of θ , while we focus on asymptotic properties of $\hat{\theta}$. Several important theoretical advances have been made in recent years by drawing upon a connection between the subsampled structure of ensembles and the classical theory of U -statistics. But the theory is unsatisfactory from two perspectives: contrary to common practice, it requires sampling without replacement; the procedures for estimating predictive variance resulting from these studies are plagued by severe bias and extreme computational overhead. In this work, we develop a general framework for analyzing the asymptotic behavior of V -statistics, demonstrating asymptotic normality under precise regularity conditions and establishing previously unreported connections to U -statistics. In addition, we propose an efficient variance estimation algorithm which also synthesizes most of the existing methods. As an application, we apply this result to design a stopping rule for determining the appropriate tree depth in model distillation.

Chapter 4 focuses on the stability of model explanation. Post hoc explanations based on perturbations are widely used approaches to interpret a machine learning model after it has been built. This class of method usually involves generating synthetic samples, and as a result exhibits large instability due to the randomness introduced. An unstable interpretation method can hardly be trusted by practitioners. We propose S-LIME, which utilizes a hypothesis testing framework based on central limit theorem for determining the number of perturbation points needed to guarantee stability of the resulting explanation.

We demonstrate its effectiveness by applying S-LIME on several real world data sets to explain different model types.

At the core of this thesis is a central theme: to better understand machine learning with the tool of statistical inference: we develop new feature importance measures and prove its statistical unbiasedness (Chapter 2); we analyze the asymptotics of ensemble predictions and propose efficient algorithms for variance estimation (Chapter 3); using hypothesis testing, we develop a stabilized post hoc model explanation method (Chapter 4). In the final chapter, we discuss some possible future directions.

CHAPTER 2

UNBIASED MEASUREMENT OF FEATURE IMPORTANCE IN TREE-BASED METHODS

2.1 Introduction

This chapter examines split-improvement feature importance scores for tree-based methods. Starting with Classification and Regression Trees (CART; Breiman et al., 1984) and C4.5 (Quinlan, 2014), decision trees have been a workhorse of general machine learning, particularly within ensemble methods such as Random Forests (RF; Breiman, 2001) and Gradient Boosting Trees (Friedman, 2001). They enjoy the benefits of computational speed, few tuning parameters and natural ways of handling missing values. Recent statistical theory for ensemble methods (e.g. Denil et al., 2014; Scornet et al., 2015; Mentch and Hooker, 2016; Wager and Athey, 2018; Zhou and Hooker, 2018) has provided theoretical guarantees and allowed formal statistical inference. Variants of these models have also been proposed such as Bernoulli Random Forests (Yisen et al., 2016; Wang et al., 2017) and Random Survival Forests (Ishwaran et al., 2008). For all these reasons, tree-based methods have seen broad applications including in protein interaction models (Meyer et al., 2017), in product suggestions on Amazon (Sorokina and Cantú-Paz, 2016) and in financial risk management (Khaidem et al., 2016).

However, in common with other machine learning models, large ensembles of trees act as “black boxes”, providing predictions but little insight as to how they were arrived at. There has thus been considerable interest in providing tools either to explain the broad patterns that are modeled by these methods,

or to provide justifications for particular predictions. This chapter examines variable or feature¹ importance scores that provide global summaries of how influential a particular input dimension is in the models' predictions. These have been among the earliest diagnostic tools for machine learning and have been put to practical use as screening tools, see for example Díaz-Uriarte and De Andres (2006) and Menze et al. (2009). Thus, it is crucial that these feature importance measures reliably produce well-understood summaries.

Feature importance scores for tree-based models can be broadly split into two categories. Permutation methods rely on measuring the change in value or accuracy when the values of one feature are replaced by uninformative noise, often generated by a permutation. These have the advantage of being applicable to any function, but have been critiqued by Hooker (2007); Strobl et al. (2008); Hooker and Mentch (2019) for forcing the model to extrapolate. By contrast, in this chapter we study the alternative split-improvement scores (also known as Gini importance, or mean decrease impurity) that are specific to tree-based methods. These naturally aggregate the improvement associated with each node split and can be readily recorded within the tree building process (Breiman et al., 1984; Friedman, 2001). In Python, split-improvement is the default implementation for almost every tree-based model, including **RandomForestClassifier**, **RandomForestRegressor**, **GradientBoostingClassifier** and **GradientBoostingRegressor** from **scikit-learn** (Pedregosa et al., 2011a).

Despite their common use, split-improvement measures are biased towards features that exhibit more potential splits and in particular towards continuous features or features with large numbers of categories. This weakness was

¹We use “feature”, “variable” and “covariate” interchangeably here to indicate individual measurements that act as inputs to a machine learning model from which a prediction is made.

already noticed in Breiman et al. (1984) and Strobl et al. (2007) conducted thorough experiments followed by more discussions in Boulesteix et al. (2011) and Nicodemus (2011)². While this may not be concerning when all covariates are similarly configured, in practice it is common to have a combination of categorical and continuous variables in which emphasizing more complex features may mislead any subsequent analysis. For example, gender will be a very important binary predictor in applications related to medical treatment; whether the user is a paid subscriber is also central to some tasks such as in Amazon and Netflix. But each of these may be rated as less relevant to age which is a more complex feature in either case. In the task of ranking single nucleotide polymorphisms with respect to their ability to predict a target phenotype, researchers may overlook rare variants as common ones are systematically favoured by the split-improvement measurement. (Boulesteix et al., 2011).

We offer an intuitive rationale for this phenomenon and design a simple fix to solve the bias problem. The observed bias is similar to overfitting in training machine learning models, where we should not build the model and evaluate relevant performance using the same set of data. To fix this, split-improvement calculated from a separate test set is taken into consideration. We further demonstrate that this new measurement is unbiased in the sense that features with no predictive power for the target variable will receive an importance score of zero in expectation. These measures can be very readily implemented in tree-based software packages. We believe the proposed measurement provides a more sensible means for evaluating feature importance in practice.

In the following, we introduce some background and notation for tree-based methods in Section 2.2. In Section 2.3, split-improvement is described in detail

²See <https://explained.ai/rf-importance/> for a popular demonstration of this.