

DISS. ETH No. 19158

Effective Causal Analysis

Methods for Structure Learning and Explanations

A dissertation submitted to

ETH ZURICH

for the degree of
Doctor of Sciences

presented by

JEAN-PHILIPPE PELLET

Ing. info. dipl. EPF,
École polytechnique fédérale de Lausanne

born November 12th, 1982

citizen of St-Livres, VD

accepted on the recommendation of

Prof. Dr. Joachim M. Buhmann

Prof. Dr. Peter Widmayer

Dr. André Elisseeff

2010

Effective Causal Analysis Methods for Structure Learning and Explanations

Doctoral Thesis

Jean-Philippe Pellet

March 12th, 2010 (*submitted*)

July 2nd, 2010 (*defended*)

EIDGENÖSSISCHE TECHNISCHE HOCHSCHULE ZÜRICH
Pattern Analysis and Machine Learning Group

IBM ZURICH RESEARCH LABORATORY
Data Analytics and Business Optimization Groups

Supervised by

Dr. André Elisseeff

Prof. Dr. Joachim M. Buhmann

Co-referent

Prof. Dr. Peter Widmayer

ABSTRACT

This study is concerned with causal analysis, an approach to multivariate-data analysis that aims to identify the structure of the data-generating process to understand how external interventions influence the system. In particular, causal knowledge can help understand the data better by preventing erroneous interpretations (e.g., Simpson’s paradox), providing causal diagnostics to explain observed values, and evaluating alternative scenarios and policies that affect the data distribution.

We study two key complementary aspects of causal analysis in depth: causal-structure learning and explanatory causal reasoning. The former deals with the computationally hard task of building a causal model; the latter uses the causal model to extract targeted contextual causal information for analysts.

Our contributions to the structure-learning task are novel algorithms which improve upon current state of the art in accuracy and time complexity. We deal with the two main causal contexts: when all potential confounders are observed and when hidden confounders may exist. We argue that the presented algorithms, thanks to their reduced complexity, can tackle a new range of large causal problems.

We then examine the causal-diagnostics task and present a new approach to evidence explanation: causal-explanation trees. Given some observed variables and a fully specified causal model, we show how to extract qualitative and quantitative information about how and why the system ended up in such a state.

In addition to continuous empirical evaluation of the presented algorithms on standard benchmarks, we present in a final chapter two real-world applications of (subsets of) the introduced approaches. We discuss how validation of causal methods is difficult without artificial data, and show pragmatic alternatives to causal-model specification in adverse conditions.

Keywords: causality, causal analysis, causal networks, structure learning, hidden variables, causal explanations.

RÉSUMÉ

Cette étude traite de l'analyse causale, une approche d'analyse de données multivariées dont le but est l'identification de la structure du processus générateur de données pour comprendre comment des interventions influencent le système. Plus spécifiquement, les informations causales aident à mieux comprendre les données en empêchant des interprétations erronées (par exemple, le paradoxe de Simpson), en établissant des diagnostics causaux pour expliquer des valeurs observées, et en étant capable d'évaluer des politiques et scénarios alternatifs qui modifient la distribution des données.

Nous étudions en profondeur deux aspects complémentaires de l'analyse causale: l'apprentissage de la structure causale, et le raisonnement causal explicatif. Le premier traite de la tâche NP-difficile de construction d'un modèle causal; le second utilise le modèle causal pour extraire de manière ciblée des informations contextuelles dans le cadre d'une analyse causale.

Nos contributions en matière d'apprentissage de structure sont des algorithmes novateurs qui dépassent l'état de l'art en précision et en complexité. Nous examinons deux principaux contextes causaux: lorsque toutes les variables causalement pertinentes sont observées, et lorsque des variables cachées parasites peuvent exister. Nous argumentons que les algorithmes présentés peuvent s'attaquer à une nouvelle dimension de problèmes causaux grâce à leur complexité réduite.

Ensuite, nous nous intéressons aux diagnostics causaux et présentons une nouvelle approche pour l'explication d'observations: les arbres d'explications causales. Étant donné certaines variables observées et un modèle causal complet, nous montrons comment extraire des informations qualitatives et quantitatives pour savoir pourquoi et comment le système a généré ces observations.

En plus de l'évaluation empirique continue des algorithmes sur des tests de performance standards, nous présentons dans un chapitre final deux applications réelles de l'analyse causale de données. Nous discutons de la difficulté de la validation des méthodes causales sans données artificielles et proposons des solutions pragmatiques pour la caractérisation de modèles causaux dans ces circonstances délicates.

Mots clés: causalité, analyse causale, réseaux causaux, apprentissage de structure, variables cachées, explications causales.