



ELSEVIER

Contents lists available at ScienceDirect

## Neurocomputing

journal homepage: [www.elsevier.com/locate/neucom](http://www.elsevier.com/locate/neucom)

## A causal feature selection algorithm for stock prediction modeling

Xiangzhou Zhang<sup>a,1</sup>, Yong Hu<sup>b,\*</sup>, Kang Xie<sup>a</sup>, Shouyang Wang<sup>c,d</sup>, E.W.T. Ngai<sup>e</sup>, Mei Liu<sup>f</sup><sup>a</sup> School of Business, Sun Yat-sen University, No. 135, Xingang Xi Road, Guangzhou 510275, PR China<sup>b</sup> Business Intelligence and Knowledge Discovery, School of Management, Guangdong University of Foreign Studies School of Business, Sun Yat-sen University, Higher Education Mega Center, Guangzhou 510006, PR China<sup>c</sup> Institute of Systems Science, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, PR China<sup>d</sup> School of Management, Graduate University of Chinese Academy of Sciences, Beijing 100049, PR China<sup>e</sup> Department of Management and Marketing, The Hong Kong Polytechnic University, Kowloon, Hong Kong, PR China<sup>f</sup> Department of Computer Science, New Jersey Institute of Technology, University Heights Newark, NJ 07102, USA

## ARTICLE INFO

## Article history:

Received 16 November 2013

Received in revised form

26 January 2014

Accepted 26 January 2014

Available online 9 May 2014

## Keywords:

Stock prediction

Data mining

Feature selection

Causal discovery

V-structure

## ABSTRACT

A key issue of quantitative investment (QI) product design is how to select representative features for stock prediction. However, existing stock prediction models adopt feature selection algorithms that rely on correlation analysis. This paper is the first to apply observational data-based causal analysis to stock prediction. Causalities represent direct influences between various stock features (important for stock analysis), while correlations cannot distinguish direct influences from indirect ones. This study proposes the causal feature selection (CFS) algorithm to select more representative features for better stock prediction modeling. CFS first identifies causalities between variables and then, based on the results, generates a feature subset. Based on 13-year data from the Shanghai Stock Exchanges, comparative experiments were conducted between CFS and three well-known feature selection algorithms, namely, principal component analysis (PCA), decision trees (DT; CART), and the least absolute shrinkage and selection operator (LASSO). CFS performs best in terms of accuracy and precision in most cases when combined with each of the seven baseline models, and identifies 18 important consistent features. In conclusion, CFS has considerable potential to improve the development of QI product.

© 2014 Elsevier B.V. All rights reserved.

## 1. Introduction

Quantitative investment (QI) products (models/tools/systems) can provide accurate stock market prediction and help investors significantly alleviate risks of mispricing and irrational trading because of psychological factors, such as overconfidence, mental accounting, loss aversion, and so on [1,2]. One of the key issues of QI product design lies on how to select representative features for prediction. For example, Thawornwong and Enke [3] have adequately demonstrated the effectiveness of recent relevant variables (i.e., representative features) for improving stock direction prediction based on redeveloped probabilistic and feed-forward neural networks. Their work provides evidence of the importance of feature selection for QI product development.

Feature selection, a pre-processing step of data mining, can be used to filter redundant and/or irrelevant features [4]. Feature selection

results in simpler model, easier interpretation, and faster induction and structural knowledge [5]. Although many studies have claimed and/or verified that feature selection is the key process in stock prediction modeling [4], identifying more representative features and improving stock prediction are challenging issues that need to be considered. Common feature selection algorithms adopted in stock prediction models include stepwise regression analysis (SRA), principle component analysis (PCA), decision tree (DT), and information gain [3,4,6]. However, all these algorithms can only reveal underlying correlations/associations and cannot determine the direct (i.e., causal) influence of stock features (inputs) on stock return (output).

This paper aims to provide further insight on the application of observational data-based causal discovery approach on feature selection. Based on Pearl's theory [7] and Hu et al.'s study [8], this study proposes the causal feature selection (CFS) algorithm, with the goal of selecting the optimal feature subset for better stock prediction performance and identifying more representative features for better stock market analysis. This study is highlighted on the following two aspects:

First, causal analysis is applied to identify direct influences between variables. Correlation does not imply causation, while causation requires additional counterfactual dependence.

\* Corresponding author.

E-mail addresses: [zhxzhou@mail2.sysu.edu.cn](mailto:zhxzhou@mail2.sysu.edu.cn) (X. Zhang), [henryhu200211@163.com](mailto:henryhu200211@163.com) (Y. Hu), [mnsxk@mail.sysu.edu.cn](mailto:mnsxk@mail.sysu.edu.cn) (K. Xie), [swang@iss.ac.cn](mailto:swang@iss.ac.cn) (S. Wang), [eric.ngai@polyu.edu.hk](mailto:eric.ngai@polyu.edu.hk) (E.W.T. Ngai), [mei.liu@njit.edu](mailto:mei.liu@njit.edu) (M. Liu).

<sup>1</sup> Co-first authors.

Causal influences are more consistent over time, which is more attractive to stock investors.

Second, to verify the proposed algorithm more objectively, extensive experiments are conducted for multi-aspect performance comparison. These experiments involve seven baseline prediction models, three popular feature selection algorithms, and various performance measures [accuracy, precision, Sharpe ratio, Sortino ratio, information ratio, and maximum drawdown (MDD)].

To evaluate the proposed CFS, listed companies in the Shanghai Stock Exchanges of China are selected as back-testing subjects, and experimental data cover the period from 2000 to 2012. Experimental results show that CFS performs best in terms of accuracy and precision in most cases and identifies 18 representative features consistently over the entire testing period. Moreover, the constructed prediction models can obtain satisfying and stable investment returns. In conclusion, CFS has considerable potential to improve the performance of existing QI products.

The remainder of this paper is organized as follows: Section 2 reviews common filter-based feature selection algorithms and stock prediction models, and then compares related works in terms of datasets, prediction models, feature selection algorithms, and so on. Section 3 introduces the proposed feature selection algorithm based on the causal discovery algorithm. Section 4 presents the experiment setting of dataset, variables, sliding window test, and evaluation strategies, and then reports the empirical results. Section 5 provides a comprehensive conclusion.

## 2. Literature review

### 2.1. Stock prediction

Although the efficient market hypothesis [9] is against stock prediction based on past publicly available information, considerable studies suggest that some markets, especially the emerging markets, are not fully efficient, and prediction of future stock

prices/returns may produce better results than random selection [10,11].

Recent studies on stock prediction can be roughly grouped into two types: (a) time series forecasting [12–16] and (b) trend prediction [4,17–20]. A time series forecasting model is trained to fit the historical return/price series of individual stock and is used to predict the future return/price. A trend prediction model is trained to obtain the relationship between various fundamental and/or technical variables and the (rise and decline) movement of stock price (i.e., positive or negative return).

Many popular data mining algorithms have been widely used in stock trend prediction models, including logistic regression (LR) [3,6,21], neural network (NN) [4,21,22], support vector machine (SVM) [23], and decision tree (DT) [21,24]. However, the Bayesian network (BN) [including its variants and naïve Bayes (NB)] has seldom been used directly for stock forecasting/prediction; only Zuo and Kita [25] used BN according to our uncomprehensive search.

Table 1 lists related works in terms of their datasets, prediction models, and feature selection algorithms. Common feature selection algorithms used in stock prediction/forecasting models include SRA, PCA, genetic algorithm (GA), information gain, and so on. Numerous related studies consider both technical and fundamental variables (including economic variables). However, the number of input features used in these studies is different. Currently, no generally agreed-upon representative features for stock prediction are available, and no “best” feature selection algorithm exists. This information motivated us to explore a novel feature selection algorithm by introducing the technique of observational data-based causal discovery to identify a more representative and compact feature set for constructing a simple prediction model with excellent performance.

### 2.2. Feature selection

Common feature selection algorithms can be grouped into following two types: (1) filter and (2) wrapper approaches [28]. The filter approaches use the general characteristics of the training

**Table 1**  
Comparison of related work.

Work	Dataset	Prediction model	Input features	Feature selection
Chang et al. [26]	S&P500 index and 5 stocks in S&P500	Case based FDT <sup>a</sup>	8 Technical indices	SRA <sup>b</sup>
Tsai et al. [21]	Electronic industry of the Taiwan stock market	Classifier ensembles	19 Financial ratios and 11 economic indicators	–
Tsai and Hsiao [4]	Electronic corporations in Taiwan Stock Exchange	BPNN <sup>c</sup>	22 Fundamental indices and 63 macroeconomic indices	PCA, <sup>d</sup> CART, <sup>e</sup> GA <sup>f</sup>
Lai et al. [27]	3 Taiwan Stock Exchange Corporations	K-means + GAFDT <sup>g</sup>	7 Technical indices	SRA
Lee [23]	NASDAQ index	SVM <sup>h</sup>	17 Financial and economic variables	F_SFS <sup>i</sup>
Enke and Thawornwong [6]	S&P 500 index	BPNN, GRNN, <sup>j</sup> PNN, <sup>k</sup> LR <sup>l</sup>	31 Financial and economic variables	Info. gain
Thawornwong and Enke [3]	S&P 500 index	BPNN, GRNN, PNN, LR	31 Financial and economic variables	Info. gain
Lam [22]	364 S&P companies	BPNN	16 Financial variables and 11 macroeconomic variables	–

<sup>a</sup> FDT: fuzzy decision tree.

<sup>b</sup> SRA: step-wise regression analysis.

<sup>c</sup> BPNN: back propagation neural network.

<sup>d</sup> PCA: principle component analysis.

<sup>e</sup> CART: classification and regression tree.

<sup>f</sup> GA: genetic algorithm.

<sup>g</sup> GAFDT: genetic algorithm-based fuzzy decision tree.

<sup>h</sup> SVM: support vector machine.

<sup>i</sup> F\_SFS: F-score and supported sequential forward search.

<sup>j</sup> GRNN: generalized regression neural network.

<sup>k</sup> PNN: probabilistic neural network.

<sup>l</sup> LR: linear regression.

Download English Version:

<https://daneshyari.com/en/article/412259>

Download Persian Version:

<https://daneshyari.com/article/412259>

[Daneshyari.com](https://daneshyari.com)