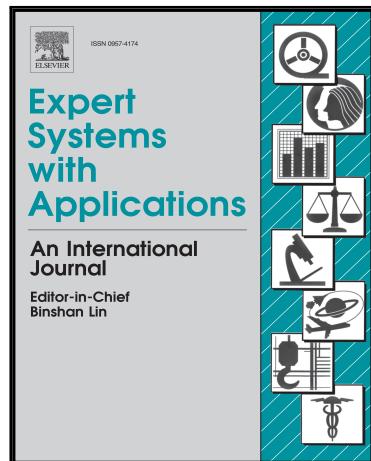


Accepted Manuscript

Literature Review: Machine Learning Techniques Applied to Financial Market Prediction

Bruno Miranda Henrique, Vinicius Amorim Sobreiro, Herbert Kimura

PII: S0957-4174(19)30017-X
DOI: <https://doi.org/10.1016/j.eswa.2019.01.012>
Reference: ESWA 12413



To appear in: *Expert Systems With Applications*

Received date: 15 May 2018
Revised date: 30 August 2018
Accepted date: 4 January 2019

Please cite this article as: Bruno Miranda Henrique, Vinicius Amorim Sobreiro, Herbert Kimura, Literature Review: Machine Learning Techniques Applied to Financial Market Prediction, *Expert Systems With Applications* (2019), doi: <https://doi.org/10.1016/j.eswa.2019.01.012>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Highlights

- The search for models to predict the prices is still a highly researched topic;
- The prices are financial time series that are difficult to predict; and
- The machine learning area applied to the prediction of financial market prices.

ACCEPTED MANUSCRIPT

Literature Review: Machine Learning Techniques Applied to Financial Market Prediction[☆]

Bruno Miranda Henrique^a, Vinicius Amorim Sobreiro^{a,*}, Herbert Kimura^a

^a*University of Brasília, Department of Management, Campus Darcy Ribeiro, Brasília, Federal District, 70910-900, Brazil.*

Abstract

The search for models to predict the prices of financial markets is still a highly researched topic, despite major related challenges. The prices of financial assets are non-linear, dynamic, and chaotic; thus, they are financial time series that are difficult to predict. Among the latest techniques, machine learning models are some of the most researched, given their capabilities for recognizing complex patterns in various applications. With the high productivity in the machine learning area applied to the prediction of financial market prices, objective methods are required for a consistent analysis of the most relevant bibliography on the subject. This article proposes the use of bibliographic survey techniques that highlight the most important texts for an area of research. Specifically, these techniques are applied to the literature about machine learning for predicting financial market values, resulting in a bibliographical review of the most important studies about this topic. Fifty-seven texts were reviewed, and a classification was proposed for markets, assets, methods, and variables. Among the main results, of particular note is the greater number of studies that use data from the North American market. The most commonly used models for prediction involve support vector machines (SVMs) and neural networks. It was concluded that the research theme is still relevant and that the use of data from developing markets is a research opportunity.

Keywords: Financial time series prediction, machine learning, literature review, main path analysis.

[☆]This document was a collaborative effort.

*Corresponding author

Email addresses: brunomhenrique@hotmail.com (Bruno Miranda Henrique), sobreiro@unb.br (Vinicius Amorim Sobreiro), herbert.kimura@gmail.com (Herbert Kimura)

1 1. Introduction

2 The prediction of stock markets is one of the most important and challenging problems involving time
 3 series (Chen et al., 2017, p. 340). Despite the establishment of the efficient market hypothesis (EMH) by
 4 Malkiel and Fama (1970), which was later revised in Fama (1991), according to which financial markets
 5 follow random pathways and therefore are unpredictable, the search for models and profitable systems is
 6 still attracting a lot of attention from academia (Weng et al., 2017, p. 153). In the specialized literature, there is
 7 evidence contrary to the efficiency of financial markets, as summarized by Kumar et al. (2016), Atsalakis and
 8 Valavanis (2009), Malkiel (2003), and Fama (1991). Additionally, a predictive model capable of consistently
 9 generating returns above the market indices over time would not only represent strong evidence contrary to
 10 the EMH but would enable large profits with financial operations.

11 The prediction of price time series in financial markets, which have a non-stationary nature, is very
 12 difficult (Zhang et al., 2017; Tay and Cao, 2001, p. 161, p. 309). They are dynamic, chaotic, noisy, non-linear
 13 series (Bezerra and Albuquerque, 2017; Göçken et al., 2016; Kumar and Thenmozhi, 2014, p. 180, p. 320,
 14 p. 285) that are influenced by the general economy, characteristics of the industries, politics, and even
 15 the psychology of investors (Zhong and Enke, 2017; Chen et al., 2017, p. 126, p. 340). Thus, the literature
 16 about financial market prediction is rich in methods and practical applications regarding historical data for
 17 evaluating the profitability of techniques.

18 Among the classic financial market prediction techniques, of particular note are the following: technical
 19 analysis (TA) with standards of support and resistance and indicators calculated from past prices (Chen et al.,
 20 2014, pp. 329–330) to indicate bearish or bullish trends (Lahmiri, 2014b, p. 1450013-3) and fundamentalist
 21 analysis, which seeks economic factors that have an influence on market trends (Cavalcante et al., 2016, p.
 22 194). Stock prices and market indices are also treated with time series analysis tools. The initial prediction
 23 techniques were moving averages, autoregressive models, discriminating analyses, and correlations (Kumar
 24 and Thenmozhi, 2014; Wang et al., 2012, p. 285; p. 758). More recently, a promising area of research in the
 25 prediction of time series is that of artificial intelligence (Yan et al., 2017; Wang et al., 2012, p. 2266; p. 758),
 26 given that the techniques are designed to address chaotic data, randomness, and non-linearity (Chen et al.,
 27 2017, pp. 340–341).

28 Technological advances have made it possible to analyse large historical price databases with computa-
 29 tional systems, as introduced by Chiang et al. (2016, p. 195). The intense computational use of intelligent
 30 predictive models is commonly studied under the title of machine learning. According to Hsu et al. (2016, p.
 31 215), it is common to test techniques for analysing time series using data from the financial market, given its
 32 difficult predictability. Thus, the literature regarding financial market prediction using machine learning
 33 is vast, which hinders revisions, systematizations of models and techniques, and searches for material to
 34 determine the state of the art. Tools are needed to objectively and quantitatively select the most relevant

³⁵ works for a literature review covering the most influential articles. Thus, the intention of this article is to
³⁶ present methods for the selection of the main advances in machine learning applied to financial market
³⁷ prediction, in order to present a review of the selected articles, clarifying the knowledge flow that the
³⁸ literature follows, and to propose a classification for the articles. Additionally, this paper brings a summary
³⁹ of the best procedures followed by the literature on applying machine learning to financial time series
⁴⁰ forecasting.

⁴¹ The selection of the most relevant literature for the proposed review was performed by searching the
⁴² theme in the *Scopus* database and validating the group of articles selected as a representative sample of the
⁴³ literature. For the review of the articles, objective parameters were proposed as a means of indicating those
⁴⁴ that are most relevant. Thus, the following were included in the review: the most-cited articles, the articles
⁴⁵ with the greatest bibliographic coupling, those with the greatest relationship in a co-citation network, those
⁴⁶ most recently published, and those that are part of the main path of the literature – a technique used to
⁴⁷ trace the flow of knowledge in a given scientific discipline ([Liu et al., 2013](#), p. 4). The articles were then
⁴⁸ objectively reviewed and subsequently classified in terms of the following: the markets used as data sources
⁴⁹ for tests, predictive variables, predicted variables, methods or models, and performance measures used in
⁵⁰ the comparisons. In all, 54 articles were reviewed and classified, covering the specialized literature from 1991
⁵¹ to 2017. Based on the searches in databases of related articles, no reviews were found with such objective
⁵² techniques with main path analysis on the theme proposed here.

⁵³ The literature review is a method for investigating the approaches of a studied topic, as stated by
⁵⁴ [Lage Junior and Godinho Filho \(2010](#), p. 14). The following section briefly presents a review of the main
⁵⁵ machine learning techniques covered in the articles selected for this study. Subsequently, in Sec. 3, the
⁵⁶ methods used in the selection of the literature most relevant to this work are described. This concerns
⁵⁷ seeking the state of the art of a science, systemizing the information, and indicating the challenges for future
⁵⁸ studies. The objective was to use quantitative methods in the selection of the most important articles about
⁵⁹ financial market prediction via machine learning, using information from citations and years of publication.
⁶⁰ In Sec. 4, the results of the bibliometric research are presented, revealing the most important articles in the
⁶¹ field under study. The main development path of the theme in the literature is also presented. The researched
⁶² articles are systematically reviewed and classified in Sec. 5. Finally, based on the reviewed works, the best
⁶³ practices for applying machine learning models to financial time series predictions are outlined in Sec. 6.

⁶⁴ 2. Brief review of machine learning techniques

⁶⁵ Machine learning techniques, which integrate artificial intelligence systems, seek to extract patterns
⁶⁶ learned from historical data – in a process known as training or learning to subsequently make predictions
⁶⁷ about new data ([Xiao et al., 2013](#), pp. 99–100). Empirical studies using machine learning commonly have

68 two main phases. The first one addresses the selection of relevant variables and models for the prediction,
 69 separating a portion of the data for the training and validation of the models, thus optimizing them. The
 70 second phase applies the optimized models to the data intended for testing, thus measuring predictive
 71 performance. The basic techniques used in the literature include the following: artificial neural networks
 72 (ANNs), support vector machines (SVMs), and random forests (RFs).

73 In general, neural networks model biological processes ([Adya and Collopy, 1998](#), p. 481) – specifically
 74 the human system of learning and identifying patterns ([Tsaih et al., 1998](#), p. 162). The basic unit of these
 75 networks, the neuron, emulates the human equivalent, with dendrites for receiving input variables to emit
 76 an output value ([Laboissiere et al., 2015](#), pp. 67–68), which can serve as input for other neurons. The layers
 77 of basic processing units of the neural networks are interconnected, attributing weights for each connection
 78 ([Lahmiri, 2014a](#), p. 1450008-5), which are adjusted in the learning process of the network ([Kumar and](#)
 79 [Thenmozhi, 2014](#), p. 291), in the first training phase mentioned in the previous paragraph. This phase
 80 optimizes not only the interconnections between the layers of neurons but also the parameters of the transfer
 81 functions between one layer and another, thus minimizing the errors. Finally, the last layer of the neural
 82 network is responsible for summing all the signals from the previous layer into just one output signal – the
 83 network's response to certain input data.

84 Whereas neural networks seek to minimize the errors of their empirical responses in the training stage,
 85 an SVM seeks to minimize the upper threshold of the error of its classifications ([Huang et al., 2005](#), p.
 86 2514). To do so, an SVM takes the training samples and transforms them from their original dimension
 87 space to another space, with a greater number of dimensions, in which a linear separation is approximated
 88 ([Kara et al., 2011](#), p. 5314) by a hyperplane. This algorithm, which is commonly used to classify data based
 89 on input variables in the model, seeks to minimize the margin of the classification hyperplane during
 90 the training stage of the model. The transformation of the space of original dimensions to the space in
 91 which the classifications are performed is done with the assistance of kernel functions, from estimated
 92 parameterization in the training of the model, as detailed by [Pai and Lin \(2005](#), pp. 498–499).

93 Just like ANNs and SVMs, decision trees are often used in the machine learning literature, as reviewed by
 94 [Barak et al. \(2017](#), p. 91). This method involves subdivision of the data into subsets separated by the values
 95 of the input variables until the basic classification unit is obtained, in accordance with the training samples.
 96 The consensual classifications of the most accurate trees are combined into a single one, comprising the RF
 97 algorithm proposed by [Breiman \(2001\)](#). The combination of decision trees in the RF technique can be used in
 98 regressions or classifications, leading to good results for financial market prediction, as demonstrated by
 99 [Krauss et al. \(2017\)](#), [Kumar et al. \(2016\)](#), [Ballings et al. \(2015\)](#), [Patel et al. \(2015\)](#), and [Kumar and Thenmozhi](#)
 100 ([2014](#)).

101 **3. Bibliometric analysis methods**

102 *Scopus* – which is a database of articles and citations from periodicals, of high relevance in the scientific
 103 community – was used to survey the most relevant literature about financial market prediction using
 104 machine learning. In the system of this database, it is possible to organize a database of citations containing
 105 information about the articles, such as title, author, periodical, year of publication, and references cited. The
 106 initial bibliometric analysis of the citation database revealed the most-cited articles and the distribution of the
 107 articles over the years, as presented in the study by [Seuring \(2013\)](#). The frequency of scientific publications
 108 for an area of knowledge obeys Lotka's law of 1926 ([Saam and Reiter, 1999](#), p. 135). Lotka discovered that
 109 the productivity of scientists in an area of knowledge follows a power law. Thus, the relative proportion of
 110 scientists with n publications is proportional to $1/n^2$ ([Saam and Reiter, 1999](#), p. 137), indicating that many
 111 scientists publish little material, while a few have a large number of publications. The law can be generalized
 112 to C/n^x , in which C is a constant of proportionality and x has a value of approximately 2. Lotka was unable
 113 to explain this law, but other studies interpreted it as being applicable to an area of science, not to individual
 114 scientists ([Saam and Reiter, 1999](#), p. 137). In this present study, Lotka's law is used to indicate the validity of
 115 the initial search in the *Scopus* database of articles. Thus, a bibliometric search in the area of financial market
 116 prediction using machine learning should be sufficiently comprehensive to obey Lotka's law.

117 One of the products of the bibliometric analysis is the listing of the most-influential authors in a research
 118 area. To measure this influence, the h and g indices for the individual performance of each author are used, in
 119 accordance with their published works, as performed by [Liu et al. \(2013\)](#). The h index incorporates citations
 120 and publications into a single number, following [Egghe \(2006](#), p. 132). This index is calculated by ordering
 121 the articles of a given author in terms of the respective number of citations and taking the highest h value of
 122 articles that have h or more citations ([Egghe, 2006](#), p. 132). Therefore, the articles below this ordination will
 123 have no more than h citations. However, [Egghe \(2006\)](#) argued that the h index is insensitive not only to the
 124 rarely cited articles of a given author but also to the articles with a very large number of citations. [Egghe](#)
 125 ([2006](#)) also noted that if the number of citations of an article increases over the years, the h index remains
 126 unchanged. Thus, Egghe introduced the g index – the most-cited g articles of an author that together have at
 127 least g^2 citations ([Egghe, 2006](#), p. 132).

128 One of the objectives of the bibliometric analysis performed in this study is to indicate the most important
 129 articles and journals about financial market prediction using machine learning. For this reason, the following
 130 are tabulated: the most-cited articles according to the *Scopus* database and the journals with the largest
 131 number of publications – as performed by [Mariano et al. \(2015](#), p. 38). In addition to the most-cited articles
 132 about the theme in the *Scopus* database, the articles most cited by the entire list searched are also recorded,
 133 that is, which articles are the most cited by those surveyed in the bibliometrics. Such a procedure reveals
 134 some articles that may not be part of the bibliometric survey but which are important references in the

¹³⁵ construction of the basic knowledge about financial market prediction.

¹³⁶ In addition to listing the most important bibliography on financial market prediction, this study intends
¹³⁷ to explain the relationship between the articles. For this reason, direct citations, bibliographic coupling, and
¹³⁸ a co-citation network are used. Citation analysis is a low-cost evaluation tool used to evaluate the acceptance
¹³⁹ of an academic article (Liu et al., 2013, p. 4). Bibliographic coupling networks – introduced by Kessler (1963)
¹⁴⁰ – relate to articles that use the same set of references. This involves a matrix A, whose a_{ij} element represents
¹⁴¹ how many references the articles i and j have in common. Such networks have the advantage of increasing
¹⁴² the chance of representativeness of more current articles, which do not yet have sufficient publication time
¹⁴³ to reach the level of a large number of citations. Small (1973) defined a co-citation network by the frequency
¹⁴⁴ with which authors are cited together. Both networks – bibliographic coupling and co-citation – are means of
¹⁴⁵ visualizing the structure of a scientific field. Direct citations, however, can be used to analyse the connectivity
¹⁴⁶ between the articles and the path followed by the knowledge, as performed by Hummon and Doreian
¹⁴⁷ (1989).

¹⁴⁸ Although widely used, the bibliographic coupling of Kessler (1963) is static, and the similarity between
¹⁴⁹ the articles is defined by the bibliography of the authors (Hummon and Doreian, 1989, pp. 57–58). The
¹⁵⁰ patterns of the co-citations are more dynamic, changing as the field of knowledge evolves (Small, 1973, p.
¹⁵¹ 265). Hummon and Doreian (1989) used a targeted network composed of the most important discoveries in
¹⁵² DNA theory, proposing a new way of analysing the path followed by a science. Representing the network
¹⁵³ of discoveries as a targeted graph, Hummon and Doreian (1989) proposed to discover the most important
¹⁵⁴ path via counting methods in the links between the events. This work constructed the network of direct
¹⁵⁵ citations – similar to the network discovered by Hummon and Doreian (1989) – in accordance with Alg. 1, as
¹⁵⁶ recommended by Henrique et al. (2018). The result of this algorithm is a graph with vertices representing
¹⁵⁷ articles and edges representing direct citations.

¹⁵⁸ Alg. 1 examines the entire list of articles surveyed in the search of the Scopus database, seeking the
¹⁵⁹ references of each article. When a reference of an article j is identified as an article k in the complete list
¹⁶⁰ of articles, there is a direct citation relationship. The edge of the graph originates in the cited work and
¹⁶¹ terminates in the article that references it. This procedure is adopted for the flow of presumed knowledge,
¹⁶² as performed by Hummon and Doreian (1989) and Liu et al. (2013); that is, when an article cites an earlier
¹⁶³ one, the scientific knowledge flows from the reference to the article making the citation (Liu et al., 2013, p. 4).
¹⁶⁴ Given that the publication of the cited article is inevitably before the one that cites it, the resulting graph of
¹⁶⁵ Alg. 1 is most likely acyclic.

¹⁶⁶ Please insert Alg. 1 here.

¹⁶⁷ The network of direct citations constructed in this work is used to calculate the paths of scientific

¹⁶⁸ knowledge in the area researched, as in the work of [Hummon and Doreian \(1989\)](#). Thus, the *source* and *sink* –
¹⁶⁹ nomenclature used by [Liu et al. \(2013](#), pp. 4–5) – vertices are defined. A source vertex is one from which
¹⁷⁰ only edges stem; that is, that vertex is only cited by others. Thus, the source vertex does not directly cite any
¹⁷¹ other vertex from the collection of articles studied – it is only cited. In turn, a sink vertex is not cited by any
¹⁷² other vertex. It only cites vertices from the collection of articles. A sink vertex therefore has edges only in
¹⁷³ its direction, and no edge stems from it. Paths are then formed in the literature researched – from source
¹⁷⁴ articles and destined to the sink articles, which of course are newer. Thus, one way to survey a bibliography
¹⁷⁵ using quantitative methods is to study the paths via which scientific literature evolved until arriving at the
¹⁷⁶ current state of the art.

¹⁷⁷ The method used in this work to count the weights of each path in the network of direct citations of
¹⁷⁸ the researched literature is known as search path count (SPC) – a method proposed by [Batagelj \(2003\)](#) to
¹⁷⁹ complement the techniques proposed by [Hummon and Doreian \(1989\)](#). SPC stems from the definition –
¹⁸⁰ according to criteria of the researcher – of the most suitable vertices as sources and sinks in a network of
¹⁸¹ targeted citations. From these definitions, all the paths between the sources and the sinks are computed,
¹⁸² recording how many times each of the edges of the graph is traversed. At the end of computing all possible
¹⁸³ paths, the edges receive weights according to the number of paths that pass through them. From these
¹⁸⁴ weights, the main pathways of the literature are calculated. The main path is defined as the one whose sum
¹⁸⁵ of the weights of the edges is the highest value among all the paths between source and sink vertices – an
¹⁸⁶ approach referred to as Global Search by [Liu and Lu \(2012\)](#). The SPC count method used in this work – with
¹⁸⁷ selection of the main path by the approach of that most used overall – is given by Alg. 2, adapted from
¹⁸⁸ [Henrique et al. \(2018\)](#).

¹⁸⁹ Please insert Alg. 2 here.

¹⁹⁰ 4. Results of the bibliometric analysis

¹⁹¹ The search of the *Scopus* database was performed on 13/4/2017. Combinations of the following terms
¹⁹² were used: stock market prediction/ forecasting, neural networks, data mining, stock price, classifiers,
¹⁹³ support vector machine, k-nearest neighbours, and random forest. The survey resulted in 1,478 documents,
¹⁹⁴ among which 629 were published articles and 23 were in press. To decrease the number of articles not related
¹⁹⁵ to the research theme, the results were restricted to the following areas: economics, econometrics and finance;
¹⁹⁶ business, administration, and accounting; social sciences; decision science; engineering; mathematics; and
¹⁹⁷ computer science. Thus, 547 articles were selected for this bibliometric work. A description of this database
¹⁹⁸ of articles is presented in Tab. 1.

199 The frequency of the publications per year is shown in Fig. 1, in dark bars, and it can be observed that
 200 there was growth, especially from 2007 onward. For comparison, a search in the *Scopus* database using
 201 the terms “credit risk” and “machine learning” – under the same conditions described in the preceding
 202 paragraph – returned 140 articles. Credit risk is a prominent area of finance research. The distribution of
 203 these publications is also shown in Fig. 1 in lighter bars. Examination of this figure suggests opportunities
 204 for future research about credit risk using new computational technologies. However, the predominance
 205 of the theme of financial market predictions compared to credit risk is clear in the number of publications,
 206 when both terms are associated with machine learning, with an annual growth rate of approximately 8.67%.
 207 Regarding this theme, the 20 most productive authors in the database of articles addressed are reported in
 208 Tab. 2, according to how many articles in this database are produced by each author, even in the case of
 209 co-authoring. The vast majority of the 1,151 authors searched (964 or approximately 84%) are authors or
 210 co-authors of only one article.

211 Please insert Tab. 1 here.

212 Please insert Fig. 1 here.

213 The frequency distribution of the number of publications per author among the articles researched is
 214 presented in Tab. 3. In accordance with Lotka’s law, it can be observed that most authors are an author or
 215 co-author of only one article, and few are responsible for an extensive production. The frequency distribution
 216 is shown in Fig. 2 by the circles marked as distribution observed. The frequency distribution theorized by
 217 Lotka, C/n^x , is illustrated by the continuous curve in Fig. 2. The circles represent the observed distribution,
 218 with coefficient x and constant C estimated to be 2.779123 and 0.567291, respectively. The R^2 of the estimate
 219 is 0.9252587. The Kolmogorov-Smirnov test of significance of the difference between the theoretical and
 220 observed distributions of Lotka returned a *p-value* of 0.1640792, thus indicating that there is no significant
 221 difference; that is, the bibliographic survey presented here follows Lotka’s law, as described previously.
 222 Consequently, the comprehensiveness of the results of the bibliometric search is validated as a significant
 223 sample of the totality of scientific publications about financial market prediction using machine learning.

224 Please insert Tab. 2 here.

225 In accordance with the previous description, the citation performance of each author can be measured by
 226 the *h* and *g* indices. The 20 authors with the highest *g* indices are listed in Tab. 5. However, ordering the
 227 authors according to the *h* index, the authors reported in the list of Tab. 5 are added to the list of Tab. 6. Thus,

the two indices are used to evaluate the contribution and influence of the authors, as performed by Liu et al. (2013, pp. 6–7). As indicated by these authors, the indices are highly correlated, which can be observed in Tabs. 5 and 6. In these tables, a correlation can also be observed between the h and g indices with the number of articles researched in the database surveyed, but not with the number of citations. For example, Lin, C. – author with h and g values of 4 – is cited 362 times in the database surveyed, whereas Wang, J. – author with h and g values of 9 and 16, respectively – is cited 296 times in the database searched. Additionally, studying the 27 authors of these two tables, it can be observed that many of them also appear as authors with the greatest number of articles produced in Tab. 2. Only O, J.; Wang, Y.; Kim, S.; Liu, M.; Chen, T.; Fan, C.Y.; Lu, C.J.; Quek, C.; and Lin, C. do not appear as authors with the most number of articles in Tab. 2. The number of publications per country, in turn, is summarized in Tab. 4. The number of articles and the frequency in the database searched are recorded for the first 20 countries in publications. Together, these countries account for approximately 85% of the database's articles. Similarly, the number of citations per country is listed in Tab. 8. China, Taiwan, and the United States are at the top of the production of articles and citations.

241 Please insert Tab. 3 here.

242 Please insert Fig. 2 here.

243 Please insert Tab. 4 here.

244 Please insert Tab. 5 here.

245 Please insert Tab. 6 here.

246 The 20 journals with the highest number of publications in the database of articles searched are listed in
247 Tab. 9. These journals account for approximately 44% of the articles surveyed and therefore are important
248 sources of research in financial market prediction using machine learning. It is worth highlighting that
249 the journal *Expert Systems with Applications* accounts for 13% of the total publications, which indicates not
250 only the number of useful references for the area published in this journal but also that it is a potential
251 target journal for future studies. Listed in Tab. 9, the journals *Neurocomputing*, *Applied Soft Computing Journal*,
252 *Decision Support Systems*, *Neural Computing and Applications*, *Neural Network World*, and *Journal of Forecasting*
253 together account for 17% of the articles surveyed, thus being alternatives for submissions of new studies. For

²⁵⁴ searches in this area, the 20 most commonly used keywords in the database searched are listed in Tab. 10.
²⁵⁵ Almost 90% of the articles searched use one or more of these keywords. The systems for search and analysis
²⁵⁶ of keywords distinguish terms in singular and plural form, considering them distinct for the purposes of
²⁵⁷ searches.

²⁵⁸ Please insert Tab. 7 here.

²⁵⁹ Please insert Tab. 8 here.

²⁶⁰ Upon recording the previous bibliographical statistics, it is then necessary to survey the most relevant
²⁶¹ literature on financial market prediction using machine learning. An important relationship is that of the
²⁶² most-cited articles – the top 20 are listed in Tab. 11. It is emphasized that the number of citations from each
²⁶³ article is related to the entire *Scopus* database of articles. It is also worth listing the references most cited by
²⁶⁴ the articles analysed in the bibliometric research. Such references do not necessarily appear in the database
²⁶⁵ of *Scopus* articles, but they are important sources for the area of financial market prediction. These most-cited
²⁶⁶ articles are listed in Tab. 12. Also listed for review are the 10 most recent articles among those surveyed in
²⁶⁷ the initial bibliometric search, all of which were published in 2017. These articles are listed in Tab. 13.

²⁶⁸ Please insert Tab. 9 here.

²⁶⁹ Please insert Tab. 10 here.

²⁷⁰ Please insert Tab. 11 here.

²⁷¹ Please insert Tab. 12 here.

²⁷² Please insert Tab. 13 here.

²⁷³ The network resulting from the bibliographic coupling of Kessler (1963) is shown in Fig. 3 only for the 20
²⁷⁴ articles with the highest number of relationships between each other. Each node in the figure represents
²⁷⁵ an article, and the arcs, or links, are the relationships between them. As described in Sec. 3, a relationship

²⁷⁶ represents similarities in the references of the articles. Thus, the articles in Fig. 3 have similar references
²⁷⁷ regarding financial market prediction using machine learning. The review of these articles is therefore a
²⁷⁸ means of summarizing the evolution of the field researched, considering the literature researched by the
²⁷⁹ authors. The data from the articles of Fig. 3 are explained in Tab. 14. This table has more recent articles than
²⁸⁰ the most-cited articles listed in Tab. 11 and 12. The co-citation network for this bibliographic survey – in the
²⁸¹ model proposed by Small (1973) – is illustrated by Fig. 4 for the 10 greatest relationships. The co-citation
²⁸² network of this figure reveals the articles listed in Tab. 15, also selected for review.

²⁸³ According to the previous description, a network of direct citations is constructed between the 547
²⁸⁴ articles of the bibliometric research. Alg. 1, which returns an acyclic graph, is used. Isolated vertices are
²⁸⁵ removed; that is, those that do not have direct citations to the others or are not cited. Such vertices will not
²⁸⁶ be considered in the main path calculations because they do not relate to any other. With the network of
²⁸⁷ direct citations at hand, the source vertices and the sink vertices must be selected, in accordance with Liu
²⁸⁸ et al. (2013). Such a choice may be subjective, at the discretion of the researcher of the network; however, this
²⁸⁹ work opted to test all the paths between all possible sources and all possible sinks. The main path selected
²⁹⁰ was the one with the highest sum of weights, as described earlier. Thus, in the network of direct citations
²⁹¹ among the 547 articles constructed by Alg. 1, 60 articles were identified that are only cited, not referencing
²⁹² any other article from the database. Such articles constitute the sources, from which only edges originate.

293

 Please insert Fig. 3 here.

294

 Please insert Tab. 14 here.

295

 Please insert Fig. 4 here.

296

 Please insert Tab. 15 here.

²⁹⁷ The pairing combination then occurs for all 60 source vertices and all 165 sink vertices, calculating all the
²⁹⁸ possible paths between each one. According to Alg. 2, each edge receives a weight according to how many
²⁹⁹ paths pass through it. The main path of the literature – in accordance with the method of Alg. 2 – is shown
³⁰⁰ in Fig. 5, calculated with a total weight of 4686. The articles that constitute the vertices are listed in Tab. 16. It
³⁰¹ can be observed that half of the articles of the main path calculated were published by the journal *Expert*
³⁰² *Systems with Applications*, which is consistent with the results presented in Tab. 9.

303

Please insert Fig. 5 here.

304

Please insert Tab. 16 here.

305 **5. Review of the selected literature**

306 Brief comments on the literature selected via the quantitative methods described in Secs. 3 and 4 follow.
 307 Commented upon are the most-cited articles – in accordance with the *Scopus* database – listed in Tab. 11, in
 308 addition to those articles most cited by the compiled database of the 547 articles in this work, listed in Tab.
 309 12. The most recent articles selected for the review are shown in Tab. 13. Also selected and reviewed are
 310 the articles with the greatest bibliographic coupling – see Tab. 14. Aiming to address the evolution of the
 311 state of the art of predictive models of machine learning applied to the financial market, the articles that are
 312 part of the main path of the literature are described – see Tab. 16. Finally, the articles reviewed were then
 313 classified according to markets, assets, and methods and variables, seeking to highlight some of the main
 314 characteristics of the literature.

315 *5.1. Most-cited articles*

316 Among the articles most cited by the bibliographical survey of Sec. 4, the classic work of [Malkiel and](#)
 317 [Fama \(1970\)](#) deserves attention, given that it established the EMH. According to this theory, financial
 318 markets immediately adjust to the information available, and it is impossible to predict their movements.
 319 The weak form of the EMH considers available information to be only the past prices of the asset ([Malkiel](#)
 320 [and Fama, 1970](#), p. 388). By adding other publicly available information, such as annual reports and the
 321 issuing of new shares, [Malkiel and Fama \(1970\)](#) addressed the semi-strong form of the EMH. Finally, the
 322 strong form of the EMH corresponds to when there is internal information monopolized by some investors.
 323 The theory proposed by [Malkiel and Fama \(1970\)](#) is critical for the prediction of financial markets because
 324 the construction of consistently profitable systems may mean the existence of evidence contrary to the EMH
 325 ([Timmermann and Granger, 2004](#), p. 16).

326 Also present in Tab. 12, the articles of [Engle \(1982\)](#) and [Bollerslev \(1986\)](#) introduced important econo-
 327 metric models used in financial market prediction. [Engle \(1982\)](#) modelled a time series using a process
 328 called autoregressive conditional heteroskedasticity (ARCH). In this model, the conditional variance present
 329 depends on the terms of previous errors, keeping the unconditional variance constant. [Bollerslev \(1986\)](#),
 330 in turn, generalized the ARCH model, considering its own variance to be an autoregressive process, intro-
 331 ducing the generalized autoregressive conditional heteroskedasticity (GARCH) model. Although widely
 332 applied in the prediction of time series, ARCH and GARCH assume a linear process of generating the values

333 of the time series (Cavalcante et al., 2016, p. 197). However, markets are characterized by nonlinearities,
 334 interacting with political and economic conditions and the expectations of their operators (Göçken et al.,
 335 2016, p. 320), making GARCH assumptions inadequate for many financial time series applications (Lahmiri
 336 and Boukadoum, 2015, p. 1550001-2). Thus, other methods are used, such as the one proposed by Elman
 337 (1990), which introduced a prediction network with precursor memory of some models of neural networks.
 338 Campbell (1987), in turn, sought to document variables that predict stock returns in two distinct periods.
 339 However, Campbell (1987, p. 393) concluded that no simple model can anticipate all the variations in return
 340 on stock prices.

341 The most-cited article in the *Scopus* database among those listed in Sec. 4 is that by Kim (2003). As
 342 recorded in Tab. 1, this article relies on the total of 546 citations from others in the *Scopus* database and is
 343 referenced 39 times on average per year. Kim (2003) addressed the application of SVMs to classify the daily
 344 direction of the Korean stock market index (KOSPI), using technical analysis (TA) indicators as predictive
 345 variables. The results were compared to those obtained with neural networks and case-based reasoning
 346 (CBR), and the support vector machine (SVM) achieved better performance, as measured by the accuracy of
 347 the predictions. However, the work of Kim (2003) may be used as a reference with care, for its results contain
 348 data-snooping bias: the author selects the best SVM parameters based on test accuracy. Machine learning
 349 models should choose the best parameters based on training the model using historical data, given that test
 350 data are not known *a priori*. That is precisely what the model seeks to predict. Following the same line of
 351 reasoning as Kim (2003), the work of Huang et al. (2005) – present in Tab. 12 – also uses an SVM to classify
 352 the direction of the market, comparing its performance with linear discriminant analysis (LDA), quadratic
 353 discriminant analysis (QDA), and the Elman backpropagation neural network (EBNN). Those authors,
 354 however, make use of training data for models' parameters selection. The results indicated greater accuracy
 355 in predictions using the SVM alone and combined with the other methods. The tests were performed on the
 356 Japanese market index (NIKKEI 225) using weekly quotes. On the other hand, among the articles listed Tab.
 357 12, Pai and Lin (2005) also used an SVM as a prediction method, not for the direction of prices, like Kim
 358 (2003) and Huang et al. (2005), but to predict stock values. Additionally, Pai and Lin (2005) combined an
 359 SVM with autoregressive integrated moving average (ARIMA) in a hybrid system capable of fewer errors
 360 than those obtained using the models separately.

361 The SVM classification model can be adapted as a regression to predict values in financial time series –
 362 in this case it is known as Support Vector Regression (SVR). This model has been used, for example, in the
 363 work of Huang and Tsai (2009). The authors combined SVR with a Self-Organizing Feature Map (SOFM)
 364 in two stages to predict the value of an index of the Taiwan market (FITX) in daily quotes. The SOFM is a
 365 method for spatial mapping of the training samples according to their similarities (Huang and Tsai, 2009,
 366 p. 1531). The inputs are indicators of TA, grouped by the SOFM according to their similarities in clusters,
 367 which are in turn fed into SVR models. The results achieved by the hybrid model of Huang and Tsai (2009)

368 are superior to those obtained with only the use of SVR as a predictive model. In their model, Yu et al.
 369 (2009) – present in Tab. 12 – used a variation of the SVM known as Least Squares SVM (LSSVM), which
 370 has a lower computational cost than the original SVM and good generalization capacity (Yu et al., 2009, p.
 371 88). The authors proposed an evolutionary LSSVM with the use of a Genetic Algorithm (GA) – a process of
 372 selecting values optimized for each generation (Yu et al., 2009, p. 88). Thus, in the proposal of Yu et al. (2009),
 373 a GA is used at the time of the selection of variables, among the TA values and fundamentalist indices,
 374 and in the optimal parameterization of the LSSVM. The results obtained are superior to conventional SVM
 375 models, ARIMA, LDA, and neural networks of the Backpropagation Neural Network (BPNN) type. It is
 376 worth noting Yu et al. (2009) apply McNemar's tests to formally compare models' performance. Not all
 377 works go beyond looking at hit rates.¹

378 As observed in Tabs. 11 and 12, the majority of the most-cited articles in the market prediction literature
 379 researched apply variations of neural networks. However, another algorithm that is quite present in the
 380 texts on prediction is the SVM, which was used by Kim (2003) – the most-cited article from Tab. 11. The work
 381 of Kara et al. (2011) – present in Tab. 11 – compares the two basic models of SVMs and ANNs regarding their
 382 predictive capabilities for the daily direction in the Turkish market. In addition to using the index values
 383 from an emerging market as data, Kara et al. (2011) are also notable because they consider all 10 years of
 384 daily prices in the parameterization of the models, ensuring they remain as general as possible. Samples
 385 from all 10 years are also considered when training the models. The authors conclude on the predictive
 386 superiority of ANNs under their proposed conditions. However, Kara et al. (2011) made predictions for the
 387 current direction of daily prices at market close, using TA indicators calculated by those same closing prices
 388 as inputs to the models. Obviously, closing prices are available only after the market closes, and the market's
 389 direction is already defined. Therefore, predicting the current prices' direction upon market closure using
 390 machine learning and those current closing prices, in the form of TA values, as inputs may not make sense
 391 for practical use, but this approach enables straightforward model comparisons.

392 Among the articles compiled in Tab. 12, the earliest reference that effectively addresses the prediction of
 393 stock prices is the work of Yoon et al. (1993). The authors applied neural networks to the prediction of the
 394 financial performance of stocks relative to the market, and they based their study on the good predictive
 395 performances reported in previous works (Yoon et al., 1993, p. 51), demonstrating that neural networks
 396 achieve better results than those obtained by discriminant analysis. Neural networks – a model based on the
 397 human nervous system (Göçken et al., 2016, p. 322) – have been widely used as prediction methods (Göçken
 398 et al., 2016, p. 320). Among the most-cited references about this subject are the works of Hornik et al. (1989)
 399 and Hornik (1991), present in Tab. 12. Laying the bases for general applications, such articles rigorously
 400 demonstrate the approximation abilities of neural networks for mathematical functions with a certain

¹ McNemar's tests are built pairing the distributions given by the results from each model, applying a χ^2 test on the null hypothesis that the models have the same performance. For more details, refer to Dietterich (1998).

accuracy. This approximation ability of the neural networks is explored in the article by [Abu-Mostafa and Atiya \(1996\)](#), which provides an initial approach regarding financial market predictions before proposing the system based on neural networks and hints. These hints are a learning process joining training data and previous knowledge ([Abu-Mostafa and Atiya, 1996](#), p. 209), such as a known property of any financial asset. However [Abu-Mostafa and Atiya \(1996\)](#) are ultimately vague on formal hints definition and validation.

Neural networks continued to be explored in the articles listed in Tabs. [11](#) and [12](#). Of particular note is the review presented by [Adya and Collopy \(1998\)](#), which aimed to summarize the criteria for evaluating predictive work on neural networks. Among the criteria suggested by the authors were validation in test data (out-of-sample) and generalization capacity and stability of the proposed model. Another review study regarding the use of neural networks in financial market predictions – listed in Tab. [12](#) – is the article by [Zhang et al. \(1998\)](#). The authors presented ANNs as predictive models and commented on previous results from the literature, concluding with the adaptation of neural networks to predictions in the financial market due to their adaptability and ability to handle nonlinearities present in time series ([Zhang et al., 1998](#), p. 55), among other factors.

As observed from the articles listed in Tabs. [11](#) and [12](#), the literature about financial market prediction involves the use of many models based on neural networks. For example, the work of [Kamstra and Donaldson \(1996\)](#) used ANNs to combine predictions about the indices of developed markets; for example, the S&P500, NIKKEI, TSEC, and FTSE. The American S&P 500 index was also used to test a hybrid predictive model proposed by [Tsaih et al. \(1998\)](#). These authors constructed their model from variables of the TA and rules given by experts and scholars, using them as inputs for the ANNs in the prediction of the direction of the S&P500. They rely on cases termed “obvious” and “non-obvious”, but with neat definitions for each one. The direction of returns was also the dependent variable sought by [Fernandez-Rodriguez et al. \(2000\)](#) – see Tab. [11](#). The authors applied ANNs to Madrid’s market index, using the returns from nine days earlier as independent variables. The results showed the superiority of the neural networks over buy-and-hold strategies for almost all of the tested periods ([Fernandez-Rodriguez et al., 2000](#), p. 93). Buy-and-hold means the acquisition and holding of an asset for a given period of time ([Chiang et al., 2016](#), p. 201), thus exposing it to the variations in its market price. A distinguishing characteristic of [Fernandez-Rodriguez et al. \(2000\)](#) is the use of lagged returns as inputs to the neural networks instead of the conventional TA indicators. The authors, however, do not provide comparisons between both approaches.

ANNs were also used in comparison with other models in the work of [Leung et al. \(2000\)](#) – an important article present in both Tabs. [11](#) and [12](#). The authors contrasted predictions of values and monthly direction of the S&P500, FTSE, and NIKKEI indices using neural networks, LDA, and regressions. Among the input variables used in the models are interest rates, indices of industrial production and consumer prices, and previous returns. The main conclusion of the study is that the models for direction prediction have superior performance to the models for predicting value ([Leung et al., 2000](#), p. 188), measured not only by the

436 successful predictions but also by the return obtained in operations strategies. [Chen et al. \(2003\)](#) – another
 437 important study present in Tabs. 11 and 12 – applied ANNs in the prediction of returns. The authors
 438 conducted their predictions in the emerging market of Taiwan, using for this purpose a Probabilistic Neural
 439 Network (PNN). Such networks use Bayesian probability, deal better with the effects of outliers, and are
 440 faster in the learning process ([Chen et al., 2003](#), p. 906). [Kim and Han \(2000\)](#), in turn, addressed the reduction
 441 in dimensionality of variables to obtain the input variables of the ANNs through a GA. This algorithm
 442 is used to discretize the continuous values of the TA indicators and optimize the weights in the network
 443 connections, which leads to better predictive performance. [Leigh et al. \(2002\)](#) – present in Tab. 11 – also
 444 discussed the use of GAs for optimizing neural networks, comparing their predictive performance with a
 445 visual pattern of the TA called “bull flag”. Although that pattern finds definition in the text of [Leigh et al.](#)
 446 ([2002](#)), results may vary as how other researchers build their template for a pattern recognition tool. Such
 447 is a big challenge of TA visual patterns recognition for they are generally based on loose or ambiguous
 448 definitions.

449 [Thawornwong and Enke \(2004\)](#) and [Enke and Thawornwong \(2005\)](#) – presented in Tab. 12 – investigated
 450 how the predictions of neural networks vary according to the selection of the input variables. To do so, they
 451 proposed a measure of the relevance of each input variable according to how much information is added
 452 to the model with their use. The authors concluded that the models that dynamically select variables with
 453 the period are more profitable and less risky ([Thawornwong and Enke, 2004](#), pp. 226–227). [Armano et al.](#)
 454 ([2005](#)), in turn, proposed the use of a layer – prior to the neural network predictors – responsible for the
 455 selection of the predictors, taking indicators of TA as inputs. Only the best predictors are selected by means
 456 of a GA. Despite using a limited number of TA indicators, results outperform buy-and-hold strategies,
 457 even considering trading costs, which is a differential feature of [Armano et al. \(2005\)](#). [Hassan et al. \(2007\)](#)
 458 also used a GA; however, it was to optimize the parameters of a hidden Markov model (HMM), using an
 459 ANN to transform the input variables of the general model. The HMM is based on transition matrices and
 460 probabilities used in general predictions, such as DNA sequencing and voice recognition ([Hassan et al.](#),
 461 [2007](#), p. 171). The combined use of an ANN, a GA, and an HMM proposed by [Hassan et al. \(2007\)](#) was tested
 462 on the next day's closing price forecasting. However, the proposed model has only been applied to three
 463 stocks in the computing area, with the reported results being very close to an ARIMA model. The historical
 464 data was also very limited, revealing roughly 2 years of daily prices.

465 A good review of studies involving neural networks in financial market prediction is the article by
 466 [Atsalakis and Valavanis \(2009\)](#) – see Tab. 12. It is worth noting that this article is also one of the most cited
 467 among the articles compiled in the bibliographical survey of Sec. 4, listed in Tab. 11. In addition, the article
 468 by [Atsalakis and Valavanis \(2009\)](#) is in the list of the articles with the most bibliographic coupling – as
 469 indicated in Tab. 14 – and is part of the main path of the literature about financial market prediction, in
 470 accordance with the method presented in Sec. 3. This important work analysed 100 articles, classifying them

471 regarding markets analysed, input variables of each model, sample size, and performance and comparative
 472 measurements. Among the findings of the review, of note are the following: the average use of 4 to 10
 473 variables; an estimated total of 30% of the articles apply indicators on closing prices; and 20% use TA
 474 indicators as input variables ([Atsalakis and Valavanis, 2009](#), pp. 5933–5936). However, most of the articles
 475 resort to a combination of technical and fundamentalist indicators as inputs for their models ([Atsalakis and](#)
 476 [Valavanis, 2009](#), p. 5936).

477 In addition to ANNs and SVMs, the most-cited articles among those surveyed in Sec. 4 also address other
 478 prediction methods. [Chiu \(1994\)](#) – listed in Tab. 12 – used fuzzy logic for grouping data and classification.
 479 The purpose of grouping into clusters is to create larger groups, thus representing the behaviour of the
 480 system ([Chiu, 1994](#), p. 267). The article by [Chiu \(1994\)](#) did not address classification in the financial market,
 481 but the method can be used in the prediction of the direction of stock prices and indices. Fuzzy logic was
 482 also used by [Wang \(2002\)](#) and [Wang \(2003\)](#) in the construction of predictor systems for Taiwan's market.
 483 Another method for predicting stock values is the combination of representation algorithms and textual
 484 relevance of news to machine learning, as in [Schumaker and Chen \(2009\)](#), p. 20). These authors reported
 485 prediction results superior to those obtained by the SVM and ARIMA hybrid used by [Pai and Lin \(2005\)](#).

486 5.2. Articles with the greatest bibliographic coupling

487 As discussed earlier, the method of bibliographic survey by coupling enables the listing of more recent
 488 literature, even without a large number of citations. The articles with greater bibliographic coupling – as
 489 described in Sec. 3 – are listed in Tab. 14 and can be observed in Fig. 3. Some of these works have already
 490 been commented in Sec. 5.1, which addresses the most cited articles within the bibliographical survey
 491 performed: the article by [Atsalakis and Valavanis \(2009\)](#), which is a review of works about predictions with
 492 neural networks, and the articles by [Huang and Tsai \(2009\)](#), [Yu et al. \(2009\)](#), and [Kara et al. \(2011\)](#), which
 493 effectively applied neural networks and SVMs in predictions in the financial market. This current section is
 494 dedicated to the other articles listed in Tab. 14.

495 As noted in the articles studied in Sec. 5.1, neural networks are mainly applied to the prediction of prices
 496 and movements in the financial market. For example, [Thawornwong et al. \(2003\)](#) used only neural network
 497 algorithms to predict the direction of three American stocks in daily quotes, using the indicators from the TA
 498 as inputs. The results indicated that TA may be inconsistent in the prediction of short-term trends and that
 499 the use of neural networks and their indicators may increase predictive performance ([Thawornwong et al.,](#)
 500 [2003](#), p. 323). Additionally, TA and neural networks result in better strategies than the simple buy-and-hold
 501 ([Thawornwong et al., 2003](#), pp. 320–321), which is evaluated not only by accuracy of predictions but by
 502 profitability in simulations. As another example, the work of [Rodríguez-González et al. \(2011\)](#) applied neural
 503 networks to the relative strength indicator (RSI), achieving better predictive performance than the simple
 504 use of such an indicator. In turn, exploring the predictive effects of neural networks in developing markets,

505 Cao et al. (2005) applied univariate and multivariate ANNs to Chinese stocks, measuring performance
 506 through the mean absolute percentage error (MAPE). The authors concluded that the neural networks
 507 surpassed the predictive performance of linear models, such as the capital asset pricing model (CAPM). This
 508 model assumes that the return of an asset is a linear function of its risk in relation to the market (Cao et al.,
 509 2005, p. 2501). In their conclusions, the authors pointed out that neural networks may be more effective in
 510 predictions of developing markets' prices.

511 Neural networks have also been applied to hybrid models. Wang et al. (2012) combined predictions using
 512 the Exponential Smoothing Model (ESM) and ARIMA, which are capable of capturing linear characteristics
 513 of time series, with BPNN – a neural network for addressing the non-linear characteristics of these series.
 514 Using for tests the monthly closing prices of the Chinese SZII index, and the monthly opening prices of the
 515 American DJIA index, the authors concluded that the predictive performance of the hybrid model is superior
 516 to that obtained using the models individually. Wang et al. (2012) could, however, explicitly state which
 517 variables are used as input and briefly justify choosing opening prices for an index and closing prices for the
 518 other. Kumar and Thenmozhi (2014) also worked with hybrid models in a developing market, combining
 519 ARIMA, ANN, SVM, and RF to predict the daily returns of an index of the Indian market. The authors
 520 argue that financial time series are not absolutely linear or non-linear (Kumar and Thenmozhi, 2014, p.
 521 288), which justifies the combination of the two types of predictive models. The article indicates superiority
 522 in prediction and profitability with the use of the ARIMA and SVM hybrid. Other than traditional error
 523 measures, like mean absolute error and root mean square error, the hybrid models are also compared in
 524 terms of return, volatility, maximum drawdown and percentage of winning up and down periods, which
 525 could be incorporated as performance measurements in future articles.

526 The article by Tsai and Hsiao (2010) uses methods of selecting variables before processing the data
 527 by ANN to predict the direction of prices. Principal Component Analysis (PCA), GA, decision trees and
 528 their combinations are used, demonstrating that the previous selection of variables increases the predictive
 529 performance of neural networks. PCA is a multivariate statistical method that extracts a reduced number
 530 of factors, or components, of highly correlated elements, from the original variables (Tsai and Hsiao, 2010,
 531 p. 260). This method is presented in variants, some of which are examined by Zhong and Enke (2017) in
 532 the selection of variables before applying them to an ANN. Although it is concluded that pre-processing in
 533 the selection of variables increases the predictive performance of neural networks, the work of Zhong and
 534 Enke (2017, p. 137) indicates traditional PCA as a simpler and more efficient method than its variants in
 535 use combined with ANN. Both Tsai and Hsiao (2010) and Zhong and Enke (2017) make respective results
 536 more robust by using t-tests to assure differences in performance measures from each model is statistically
 537 significant. Similarly, Chiang et al. (2016) selected variables based on the information gain inherent to each
 538 one, before applying them to the ANN. In their approach, Chiang et al. (2016) apply variable selection in a
 539 set smaller than Tsai and Hsiao (2010) and Zhong and Enke (2017), but increment results by smoothing data

540 with wavelet transformations.

541 This paragraph was improved. Pre-processing data with wavelets prior to the application of neural
 542 networks was also explored by [Li and Kuo \(2008\)](#), who used a technique – known as the Discrete Wavelet
 543 Transform (DWT) – that involves decomposition of digital signals into their components. The components
 544 are then processed by a special class of neural network known as the Self-Organizing Map (SOM) to
 545 generate short- and long-term purchase and sale signals. The work done by [Li and Kuo \(2008\)](#) is remarkable
 546 for its pattern labeling scheme of long- and short-term trading signals as well as a defined, reproducible
 547 and profitable trading strategy. [Chang and Fan \(2008\)](#) also used decomposition into wavelets as data pre-
 548 processing, but they grouped them by homogeneous characteristics into clusters that are mapped into fuzzy
 549 logic rules. Subsequently, [Chang and Fan \(2008\)](#) applied a system proposed for the interpretation of these
 550 rules in the generation of predictions, using also the k-Nearest Neighbours (kNN) algorithm to reduce errors.
 551 The authors highlighted results superior to other models such as BPNN, although they used a small 2 years
 552 dataset of daily prices from a Taiwanese index. Transformations involving wavelets have also been used to
 553 mitigate short-term noise in index prices; for example, in [Chiang et al. \(2016\)](#). The authors showed larger
 554 returns when the data are smoothed through transformations with wavelets before applying them to neural
 555 networks ([Chiang et al., 2016](#), p. 205), as stated before.

556 Widely applied to time series, the fuzzy logic theory was developed in human linguistic terms ([Chen
 557 et al., 2014](#), p. 330). For example, the work of [Chen et al. \(2014\)](#) uses fuzzy logic in time series, aiming to
 558 overcome limitations regarding linearity assumptions in other models, such as ARIMA and GARCH. In a
 559 previous study, [Ang and Quek \(2006\)](#) combined neural networks with fuzzy logic to obtain daily predictions
 560 of stock prices better than other traditional neural network models. The system proposed by [Ang and Quek
 561 \(2006\)](#) also provides interpretability of rules – a property that is often absent in traditional ANN and SVM
 562 approaches ([Al Nasser et al., 2015; Yu et al., 2009](#)).

563 As noted in the articles listed in Sec. 5.1, many works compare predictions obtained with various methods.
 564 In this context, the article by [Ballings et al. \(2015\)](#) – listed in Tab. 14 – compares the combination of predictions
 565 of multiple classifiers, among them ANN, SVM, kNN, and Random Forest (RF). The results indicated better
 566 accuracy in the prediction of the direction of stock prices in a year using RF ([Ballings et al., 2015](#), p. 7051).
 567 Distinguishing features of [Ballings et al. \(2015\)](#) are the wide list of classifiers selected by the authors and
 568 their numerous input variables, drawn from fundamentals of 5.767 European companies. Therefore, the
 569 authors are limited to yearly predictions, for that is a very low frequency type of data, typically released
 570 monthly or yearly. [Gorenc Novak and Velušček \(2016](#), p. 793) worked with the prediction of the stock price
 571 direction for the following day, but applying only the SVM. The work of these authors stands out due to
 572 using the daily maximums of the assets, as opposed to the traditional closing price. [Gorenc Novak and
 573 Velušček \(2016](#), p. 793) observed that the volatility of the maximums is less than that of the prices at the end
 574 of the negotiation session, at closing. Therefore, the maximums would be easier to predict, as stated by the

575 authors.

576 Finally, the articles (in Tab. 14) that address financial market predictions using textual analysis deserve to
 577 be highlighted. Hájek et al. (2013) analysed sentiment about the annual reports of companies, processing
 578 terms that exert positive or negative influence on asset prices. Corporate reports are tools for communication
 579 with investors, which contain terms loaded with qualitative data (Hájek et al., 2013, p. 294). The authors
 580 processed these terms through previously constructed dictionaries, and the resulting categorizations serve
 581 as inputs for both neural networks and SVR models. The method proposed by Hájek et al. (2013) proved to
 582 be capable of predicting returns for one year in advance of the data used in the tests. Al Nasseri et al. (2015)
 583 also analysed sentiments, but in publications of a blog specializing in stock markets. The authors concluded
 584 that variations in the terms used in the texts of the blog predict trends of the American DJIA index.

585 *5.3. Articles with greater co-citation relationships*

586 According to Small (1973), a co-citation occurs when two articles are quoted by a third party in the same
 587 work. The more frequent these joint citations, the stronger the relationship between the two articles. Sec. 4
 588 lists the 10 articles with the greatest co-citation relationships in Tab. 15. Some of these articles were reviewed
 589 earlier, in Sec. 5.1, and therefore are not commented on in this section. The articles are as follows: Atsalakis
 590 and Valavanis (2009), Chen et al. (2003), Hornik et al. (1989), Huang et al. (2005), Kim and Han (2000), and
 591 Zhang et al. (1998). Thus, the next few paragraphs briefly comment only on the other articles of greater
 592 co-citation frequency from Tab. 15.

593 The study by Kimoto et al. (1990) used combined neural networks to form a single prediction of weekly
 594 buying or selling for a Japanese stock market index. It applied the basic ANN model to six economic
 595 indicators as input variables, with more lucrative results than the basic buy-and-hold strategy. The major
 596 innovation by Kimoto et al. (1990) is the notion that prediction rules must change with time, according to
 597 market conditions. Therefore, the neural networks are re-trained as time passes, in a fixed moving window
 598 fashion. A prediction period is defined and as new data arrives, the training period is forward shifted.
 599 Chang et al. (2009) also worked with the classic neural network model; however, they increased the returns
 600 obtained in simulations combining the predictions of the networks with CBR. The authors also used a stock
 601 selection model based on financial health indicators of the respective companies. The returns obtained with
 602 the combined ANN and CBR model of Chang et al. (2009) were greater than the individual returns of each
 603 model for the nine stocks selected in that work. In that paper, the innovation rests on previously selecting
 604 potentially profitable stocks and reusing only price pattern cases in which the system had previous success.

605 Contrary to the models used by Kimoto et al. (1990) and Chang et al. (2009), the article by Tay and Cao
 606 (2001) obtained superior predictions with the use of SVM. The authors compared their results to those
 607 obtained by neural networks, and they concluded that the best performance of the SVM is due to: the
 608 minimization of the structural risk, to a fewer number of parameters to be optimized by the SVM; and

609 the possibility of the neural networks converging for local solutions (Tay and Cao, 2001, p. 316). Finally,
 610 also present in Tab. 15, the book of Box et al. (2015) is an ample introduction to the subject of prediction
 611 of time series for general applications that is not restricted to only financial series. The book covers linear
 612 techniques, correlations, moving averages, and autoregressive models. Being a general compendium of time
 613 series theory, Box et al. (2015) did not specifically address machine learning techniques, but even so, it is still
 614 listed as one of the works with the highest number of co-citations among the references surveyed in this
 615 present study.

616 *5.4. Main path*

617 Below is a brief review of the main path of the literature on financial market prediction using machine
 618 learning. As stated earlier, it is a chronological survey of the main articles published on the subject, which
 619 are described in detail in Sec. 3. The main path of the literature addressed in this work – illustrated in Fig.
 620 5 – suggests the most important works for the review of methods, experiments, findings, and scientific
 621 conclusions regarding the proposed theme. Thus, this section is dedicated to detailing the main aspects of
 622 the articles listed in Tab. 16, exploring the state of the art of this literature.

623 The main path of Fig. 5 begins with the article by Kamstra and Donaldson (1996). This article is also
 624 listed in Tab. 11 and has already been commented on in Sec. 5.1. It addresses the use of neural networks
 625 in the prediction of daily volatility of the S&P500, NIKKEI, TSEC, and FTSE indices, compared with the
 626 popular linear model GARCH, based on out-of-sample test data. As the ANN is a collection of non-linear
 627 transfers that relate the output variables to the inputs (Kamstra and Donaldson, 1996, p. 51), it is a more
 628 suitable proposal for potentially non-linear data. In fact, the empirical tests of Kamstra and Donaldson
 629 (1996) indicated that the predictions of the volatility of market indices using GARCH have deviations from
 630 the actual values that are higher than those obtained with ANN. The results are confirmed by the second
 631 article of the main path, which is by the same authors as the first article. Thus, Donaldson and Kamstra
 632 (1999) concluded that combination of predictions with ANNs can provide significant improvements when
 633 compared to the linear combinations approach. It is important to notice Donaldson and Kamstra (1999)
 634 aim to predict returns, differently from Kamstra and Donaldson (1996), which work to predict volatility.
 635 Moreover, both papers confront linear and non-linear models, although some authors seem to agree that
 636 non-linear models are more suitable for financial time series predictions.

637 The authors of the first two articles of the main path did not yet consider the SVM classification model,
 638 which the initial publication is credited to Vapnik (1995). Approximately five years after Donaldson and
 639 Kamstra (1999), the work of Pai and Lin (2005) considered the combination of SVM with a linear model in
 640 the prediction of stock prices. At the time, one of the linear models most commonly used in predictions
 641 was the ARIMA (Pai and Lin, 2005, p. 498), which has limitations for capturing non-linear characteristics of
 642 time series. Thus, Pai and Lin (2005) stand out due to combining this model with an SVM, which is based

643 on minimizing structural risk through limitations in the error thresholds. Despite the limited amount of
 644 data used by the authors – just over a year of daily closing data ([Pai and Lin, 2005](#), pp. 499–500), the work
 645 concluded that the proposed hybrid model can overcome the individual use of its components but suggests
 646 optimization of parameters to achieve the best results. As a very well written paper and a great reference on
 647 early models hybridization, [Pai and Lin \(2005\)](#) could state more clearly the input variables to their system. It
 648 should be noted that the manuscript by [Pai and Lin \(2005\)](#) is also one of the most-cited articles according to
 649 the survey of Sec. 4 – see Tabs. 11 and 12.

650 Reviewing 100 articles about predictive models in the financial market using computational techniques,
 651 [Atsalakis and Valavanis \(2009\)](#) provided a classification of the studies regarding markets, variables, predic-
 652 tion methods, and performance measures. The importance of this article for the main path is demonstrated
 653 by its presence in all the results of the bibliographical survey of Sec. 4, listed, therefore, in Tabs. 11, 12,
 654 and 14. In general, the studies used four to ten predictive variables ([Atsalakis and Valavanis, 2009](#), pp.
 655 5933–5936), with the most common being the opening and closing prices of market indices. Also according
 656 to the authors, 30% of the proposed models in the articles analysed use closing prices, and 20% use TA
 657 variables, most of which combine them with statistical and fundamentalist data. [Atsalakis and Valavanis](#)
 658 ([2009](#)) highlighted ANN and SVM, with data pre-processing using normalization or PCA, among other
 659 methods. The most common performance measures used by the authors are also listed; for example, Root
 660 Mean Square Error (RMSE), Mean Absolute Error (MAE), profitability, and annual return. The authors
 661 concluded that neural and neuro-fuzzy networks are suitable predictive algorithms for the stock market,
 662 but there is still no definition for the structures of such networks – they are determined by trial and error
 663 ([Atsalakis and Valavanis, 2009](#), p. 5938).

664 A few years after the review of [Atsalakis and Valavanis \(2009\)](#), the main path continued with the work
 665 of [Kara et al. \(2011\)](#) – an article from the most-cited group, as observed in Sec. 5.1, and listed as having
 666 the greatest bibliographic coupling, as addressed in Sec. 5.2. [Kara et al. \(2011\)](#) is an excellent reference
 667 for comparing predictions using ANN and SVM in their basic forms. To compare the models, [Kara et al.](#)
 668 ([2011](#)) used 10 years of daily market prices from the Istanbul index, practically balanced on days of high
 669 days and lows in the prices. It is worth highlighting the procedure for selecting the samples for the sets of
 670 parameters, training, and tests of the authors, which ensures the presence of samples of each year in each set.
 671 Ten indicators were calculated from TA – they were pre-processed only with normalization of values before
 672 their use in the models. ANN had slightly better performance than the SVM, with statistical significance
 673 calculated via *t*-tests. Following the same methods of [Kara et al. \(2011\)](#), [Patel et al. \(2015\)](#) included in the
 674 comparisons – in addition to ANN and SVM – the RF and Naïve Bayes (NB) models. These authors also
 675 used TA indicators as predictive variables, but they innovated by considering – in addition to the continuous
 676 values of the indicators – the price trend indicated by each variable. The overall result indicated that this
 677 approach of the indicators – referred to as discrete – improves the predictions ([Patel et al., 2015](#), p. 268).

678 In an article contemporaneous to that of [Patel et al. \(2015\)](#), [Laboissiere et al. \(2015\)](#) sought the prediction
 679 of stock prices in the Brazilian market using the basic ANN form. These authors differed from previous
 680 approaches by focusing on a specific sector of the market – electric energy – and using as input variables not
 681 only the prices but also the index of the São Paulo stock exchange (BOVESPA), a specific index for the electric
 682 energy market (IEE), and the US dollar. In addition, [Laboissiere et al. \(2015\)](#) focused on the prediction of
 683 daily maximums and minimums, seeking to define thresholds for operations with the stocks studied. The
 684 authors opted for pre-processing of the prices with the Weighted Moving Average (WMA), filtering noisy
 685 fluctuations and highlighting trends ([Laboissiere et al., 2015](#), p. 68). Additionally, they also used correlation
 686 analysis between stock prices and the indices to select the most important ones as inputs for the ANN model.
 687 Finally, it is important to record that the articles of [Patel et al. \(2015\)](#) and [Laboissiere et al. \(2015\)](#) delimit the
 688 most recent publications of the main path of the literature. Thus, such articles and their successors in the
 689 main path of the literature do not have sufficient publication time to reach the number of citations of the
 690 previous articles, present in Tabs. 11 and 12.

691 After the article by [Laboissiere et al. \(2015\)](#), the main path of the literature about financial market
 692 prediction continued with an example of pattern recognition. This study by [Chen and Chen \(2016\)](#) specified
 693 an algorithm for operations based on a visual signal of a high in the prices of a stock or index. As a means
 694 of reducing dimensionality, the authors used a method of weighting the most important points of a time
 695 series ([Chen and Chen, 2016](#), p. 262), based on a visual pattern used in the TA known as a bull-flag. To avoid
 696 incurring subjectivities, a specific definition of the bull-flag pattern is given and parameterized ([Chen and](#)
 697 [Chen, 2016](#), p. 264). Based on the pattern sought, a model was calculated dynamically, seeking to assess
 698 how well the pattern adapted to the daily prices of the NASDAQ and TAIEX indices. Thus, the pattern was
 699 recognized computationally, and an operation initiated, varying as a parameter the time until its liquidity.
 700 [Chen and Chen \(2016\)](#) evaluated their algorithm through the return generated, comparing it with more
 701 advanced models, such as Genetic Algorithms (GAs).

702 At the end of the main path of the literature is the article by [Zhong and Enke \(2017\)](#), which was already
 703 commented on in Sec. 5.2. [Zhong and Enke \(2017\)](#) sought to predict the daily direction of a fund based on
 704 the American index S&P500 using the basic ANN algorithm but differing in the selection of the predictive
 705 variables. For this purpose, the authors used PCA and its variations to select the most significant variables
 706 for predictions among economic indicators, such as American Treasury rates, foreign exchange, indices
 707 of international markets, and returns for companies with large capitalization. The performance of each
 708 method of variable selection – known as dimensionality reduction – allied to ANN was measured by the
 709 MSE, confusion matrices, and profitability in a simple strategy of following purchase signals and investing
 710 in American Treasury securities in the case of selling signals. [Zhong and Enke \(2017\)](#) concluded that there is
 711 no statistically significant difference in performances using the different variations of PCA. However, the
 712 profitability measured is slightly higher using the traditional PCA for dimensionality reduction ([Zhong and](#)

⁷¹³ Enke, 2017, p. 135).

⁷¹⁴ Finally, the article by Chen et al. (2017) represents the sink node in the main path of the literature shown in
⁷¹⁵ Fig. 5; that is, where the state of the art of financial market prediction applying machine learning converges.
⁷¹⁶ Chen et al. (2017, p. 341) argued that most of the previous research considers that the predictive variables
⁷¹⁷ make equal contributions to the value obtained by the predictive market model. Chen et al. (2017) innovated
⁷¹⁸ by applying a measure of information gain in the weighting of the variables, taken from indicators of the TA.
⁷¹⁹ The variables subsequently weighted are used as inputs for the SVM and kNN models, comparing the results
⁷²⁰ to the use of the classifiers without the weighting of variables. Subsequently, a hybrid model – combining
⁷²¹ SVM and kNN – was proposed in a manner similar to that of Nayak et al. (2015). However, the model
⁷²² described by Chen et al. (2017) considered the weighting of the variables, obtaining better performance
⁷²³ measures than those reported by Nayak et al. (2015).

⁷²⁴ As observed in the preceding paragraphs, the initial studies of the main literature about financial market
⁷²⁵ prediction using machine learning contrast linear models of prediction with non-linear ones. The literature
⁷²⁶ seems to agree that the former are surpassed by the latter and are therefore only used as a benchmark in the
⁷²⁷ most recent works. Thus, the main predictive models explored in the main path of Fig. 5 are ANNs and
⁷²⁸ SVMs. The applications of both prediction approaches have variations in pre-processing, the selection of
⁷²⁹ variables, and, more recently, the hybridization of the models. Most of the studies – especially the most
⁷³⁰ recent ones of the main path – present a comparison of models and combinations between them, thus
⁷³¹ representing the state-of-the-art aspects of the theme.

⁷³² 5.5. Most recent articles

⁷³³ The following paragraphs are dedicated to reviewing the most recent publications among the articles
⁷³⁴ researched in the bibliometrics of this study. The 10 most recent articles are listed in Tab. 12. These articles
⁷³⁵ summarize the most modern approaches in the use of machine learning for financial market prediction. For
⁷³⁶ example, Weng et al. (2017) used algorithms that have been widely studied in previous literature; however,
⁷³⁷ they used input data originating from public sources of knowledge via the Internet. Specifically, Weng
⁷³⁸ et al. (2017) derived news indicators from Google News and from visits to the Apple® page on Wikipedia,
⁷³⁹ in addition to the products of this company, combining traditional indicators from the TA and predicting
⁷⁴⁰ the direction of the following day's prices using ANNs, SVMs and decision trees. This approach attained
⁷⁴¹ accuracies higher than 80%, allowing the authors to claim a higher performance than the use of SVMs in
⁷⁴² the manner proposed by Kim (2003). It should be noted that Weng et al. (2017) explores crowd-sourced
⁷⁴³ information bases freely available on the Internet, leveraging their prediction system with theoretically other
⁷⁴⁴ potential traders' opinions. Naturally, how to obtain, interpret and apply all of the relevant data are rich,
⁷⁴⁵ open research questions.

⁷⁴⁶ As stated in Sec. 5.4, the most recent approaches to predicting stocks and indices of financial markets

747 use data pre-processing and hybridization of classifiers. For example, in the work of [Zhang et al. \(2017\)](#),
 748 time series of four market indices were decomposed into individual components before the application of a
 749 multidimensional variant of kNN for predictions of closing prices and maximums. The results obtained by
 750 [Zhang et al. \(2017\)](#) were better than those obtained using traditional ARIMA and kNN. However, [Zhang](#)
 751 [et al. \(2017\)](#) do not compare their unique method to the more popular machine learning approaches of SVM
 752 or ANN. The approach of [Barak et al. \(2017\)](#) involved the diversification of classifiers regarding the data
 753 used for training, using methods of multiple partitioning of these data, as well as sampling techniques.
 754 [Barak et al. \(2017\)](#) used the selection of variables among fundamentalist indices and finally merged the
 755 classifications of the methods of greater accuracy into a single prediction of returns and risk. [Barak et al.](#)
 756 [\(2017\)](#) concluded with the superiority of combined predictions, reporting accuracies exceeding 80%. A
 757 fundamental base for the work of [Barak et al. \(2017\)](#) rests on the assumption that diversity of results
 758 by different classifiers may be combined into a superior prediction. Supposedly, diversification schemes
 759 could balance each machine learning classifiers' weakness, while strengthening overall accuracy. Therefore,
 760 combining results of individual machine learning models may be a research area as prolific as the fusion of
 761 those models themselves.

762 Still concerning the hybridization of classification models for predictions in financial markets, [Krauss](#)
 763 [et al. \(2017\)](#) worked with the combination of variants of neural networks and random trees and random
 764 forests by comparing stock portfolios created with these techniques. The authors verified that the models
 765 used in combinations have better performance than when they are used individually ([Krauss et al., 2017](#),
 766 p. 694) – the individual model used as the basis, which had the best predictive performance regarding
 767 the data used by the authors, was the RF. Working with a long period of S&P500 daily prices history, the
 768 paper from [Krauss et al. \(2017\)](#) accounts for survivorship bias in their portfolio building proposal, that is,
 769 not all stocks have always been constituent of the index. Portfolio building schemes must take that into
 770 consideration when testing for return using historical prices. [Pan et al. \(2017\)](#), in turn, applied an SVM
 771 to multi-frequency independent variables to obtain weekly price predictions of the S&P 500 index – they
 772 reported results better than those of the models using single frequency variables. [Pan et al. \(2017\)](#) stand
 773 out from other works that use SVM to predict financial time series because the independent variables are
 774 not necessarily sampled at the same rate as the model's output. It would be interesting for future works to
 775 explore economic implications from this approach. [Bezerra and Albuquerque \(2017\)](#) also applied an SVM,
 776 but in its regressor mode, combined with GARCH, to predict the volatility of the returns on financial assets.
 777 The greatest innovation of the SVR – GARCH model proposed by [Bezerra and Albuquerque \(2017\)](#) was the
 778 combination of Gaussian functions as the kernel function of the SVR. The predictions were compared to
 779 those obtained using traditional GARCH models, among others, and, for most of the tests, they achieved
 780 results with fewer MAE and RMSE errors.

781 Following the line of research of [Schumaker and Chen \(2009\)](#) and [Al Nasseri et al. \(2015\)](#), which have

already been reviewed in previous paragraphs, the work of Oliveira et al. (2017) proposed the prediction of return, volume, and volatility of multiple portfolios using textual analysis and the sentiments of specialized blogs. Listed among the most recent of Tab. 13, Oliveira et al. (2017) created indicators of sentiments with textual data, and they applied SVMs, RF and neural networks, among other techniques, to evaluate the contribution that the information from blogs has on financial market predictions. Among the various results, of particular note were the superior accuracies when taking into account the data of the blogs and SVM classifiers in the prediction of returns, especially for companies with lower capitalization, in addition to the technology, energy, and telecommunications sectors. However, the textual data did not significantly increase the accuracy of the predictions of volatility and volume. Researchers exploring sentiment analysis combined with machine learning for stock prices prediction will find a rich review of previous literature and a well developed framework example in Oliveira et al. (2017).

Neural networks continue to be researched for use in predicting financial market prices. Yan et al. (2017) used neural networks combined with Bayesian probability theory to obtain predictions better than those obtained via SVMs and traditional neural networks. The errors are further reduced by the authors with the application of Particle Swarm Optimization (PSO), which is an optimization technique based on grouping and migration of artificial life forms (Yan et al., 2017, p. 2278), in which each state, or particle, is a candidate for an optimal solution, adjusting itself according to the others. Pei et al. (2017) modified traditional neural networks by applying Legendre² polynomials in the internal layers of the networks, in addition to a special time function in the method for updating the weights of the connections between the layers. One difference in the work of Pei et al. (2017) is that the authors sought to predict the moving averages of different periods for the prices, not the prices directly or their direction. According to the authors, this approach removes the influence of accidental factors in identifying the direction of the trend (Pei et al., 2017, p. 1694). Finally, Mo and Wang (2017) also proposed neural networks modified by time functions; however, they applied them in the prediction cross-correlations between Chinese and American market indices. The work of Mo and Wang (2017) can therefore be used in the optimization of asset portfolios.

5.6. Classification of the articles

This section is dedicated to a classification of the articles reviewed in Secs. 5.1, 5.2, 5.3, 5.4, and 5.5. The articles are classified according to the markets addressed, the assets used in the empirical analyses, the types of predictive variables used as inputs, the dependent variables of the predictions, the main predictive methods used in the models, and the performance measures considered in each author's evaluations. The

²Denoted by $L_p(X)$, the Legendre polynomials are a set of orthogonal polynomials of order p that are solutions for the differential equation $\frac{d}{dx} \left[(1 - x^2) \frac{dy}{dx} \right] + p(p + 1)y = 0$. Legendre polynomials can be used to expand the coefficients of the hidden layers of neural networks. For more details, see Dash (2017).

classification subsequently proposed is in Tab. 17, with the following being sought: the most commonly used methods, methods of measuring performance considered, and markets addressed.

Tab. 17 lists 57 articles reviewed in the previous sections. The review works and those that did not directly consider financial market prediction were excluded from the classifications, as well as those used as basic references to the methods. Thus, the articles of Adya and Collopy (1998), Zhang et al. (1998), and Atsalakis and Valavanis (2009) were not classified because they are review papers. The article by Malkiel and Fama (1970), which addresses EMH, was also excluded from the classification, in addition to Engle (1982) and Bollerslev (1986), who introduced ARCH and GARCH, respectively. Although Campbell (1987) addressed the prediction of stock prices, this author's methods are not directly related to machine learning, and therefore, it was excluded from the classification. The article by Elman (1990) and the book of Box et al. (2015) are used as bases for the construction of models and are therefore outside the scope of the classification proposed in this section. Finally, the following articles were not classified: Chiu (1994), which introduced fuzzy logic in predictions but did not address financial markets; and Hornik et al. (1989) and Hornik (1991), which generically demonstrated the predictive capabilities of neural networks but did not directly address financial markets.

Analysing the articles listed in Tab. 17 regarding the markets addressed, it can be observed that almost half of the studies used North American data (approximately 47%), and one sixth of them (approximately 17%) refer to data from Taiwan. This is expected, given the economic hegemony of the USA and the vast academic production of Taiwan, as quantified in Tab. 4. Another interesting fact is that despite the large Chinese academic productivity, only six articles from Tab. 16 use data from China in their predictions (approximately 10% of the articles). Likewise, studies are recorded using data from Brazil, Russia, India, China, and South Africa (BRICS) – 10 articles or approximately 17%. As for the assets for which the predictions are calculated, most articles in Tab. 17 focus on stock market indices (more than 60%). In addition, only two studies applied prediction models simultaneously to indices and stocks.

Regarding the variables used as inputs in the models of financial market prediction, the TA indicators are the most popular in Tab. 17 – they were used in approximately 37% of the studies, followed by fundamentalist information, which was used in 26% of the studies. Only two studies explicitly applied both types of variables in their predictive models. Also of note are some studies that applied the prices of the assets themselves – or lagged prices – as inputs in their models. Regarding the prediction sought by the articles of Tab. 17, the models are basically divided between those that seek future prices or returns, and those that seek only the future direction or trend of the market analysed. In this respect, the articles whose dependent variable is the direction of the markets are predominant – approximately 42% of the articles aimed to predict the direction of selected indices or stocks, while 31% of the articles predicted prices.

845

Please insert Tab. 17 here.

846 Among the prediction methods used by the articles in Tab. 17, it can be observed that approximately 70%
 847 of the studies used at least some type of neural network; therefore, it was the classification method most used
 848 among those of machine learning. The second most used model was SVMs/SVR – a more recent approach
 849 than neural networks, which was used in approximately 37 % of the articles reviewed. The hegemony of the
 850 ANN and SVM models has already been observed in the main path of literature, in accordance with the
 851 reviews of Sec. 5.4. Few studies use other models, such as kNN, RF, or NB. Thus, the articles that only used
 852 techniques of classification or regression different from ANNs and SVMs/SVR accounted for approximately
 853 14% of the total researched.

854 Finally, it should be noted that the method for measuring performance varies with the type of prediction
 855 – direction or price – sought by the articles. Thus, articles that seek to predict the direction of the market tend
 856 to measure the performance of their models by means of accuracy. Similarly, articles that seek to predict
 857 prices, verify their performance by calculating prediction errors. Specifically, MAE and RMSE measure the
 858 average magnitude of the error, as given in Bezerra and Albuquerque (2017, p. 188), being the RMSE only
 859 a square root applied to the MSE. The MAE is also known as Mean Absolute Deviation (MAD) by some
 860 authors, such as Cao et al. (2005, p. 2506). The MAPE measure, in turn, is a percentage measure of the error.
 861 MAE, MSE, and MAPE are given by Eqs. 1, 2, and 3, respectively (Bezerra and Albuquerque, 2017; Cao et al.,
 862 2005, p. 188, p. 2506), in which N is the number of observations, y is the class or real value of an observation,
 863 and \hat{y} is the class or value estimated by the model.

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (1)$$

$$MSE = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_i}{y_i} \right)^2 \quad (2)$$

$$MAPE = \frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (3)$$

864 Some articles in Tab. 17 that involve classification of the direction of financial markets measure the
 865 performance of their respective models via a measure denominated Area Under the Curve (AUC). This is
 866 the area under a curve known as Receiver Operating Characteristics (ROC), which is constructed from the
 867 variation of the classification thresholds of a model, denoting the rates of false positives on the axis of the
 868 abscissa and the ratios of the true positives on the ordinate axis (Zweig and Campbell, 1993, pp. 564–565).
 869 There are also many articles that measure the financial returns made possible by the use of suggested

870 strategies based on their respective predictive models. One of these measures is the Sharpe rate, which
 871 is calculated as the rate between average return from an operations strategy and its standard deviation
 872 ([Fernandez-Rodriguez et al., 2000](#), p. 92).

873 For future reference, the main forecasting and optimization methods considered in machine learning
 874 applied to financial markets forecasting, according to the research conducted in this paper, are listed in Tab.
 875 [18](#) and Tab. [19](#), respectively. References are grouped by each applied method, making future citations of
 876 those papers easier. It should be noted, however, that much of the machine learning literature consists of
 877 customization and tuning of base models. Therefore, models are listed as their basic form in Tab. [18](#) and Tab.
 878 [19](#), even though some authors call them by different names once they alter or hybridize them according to
 879 their applications.

880 Please insert Tab. [18](#) here.

881 Please insert Tab. [19](#) here.

882 6. Robust research structure

883 Based on the reviewed works on the previous section, the next paragraphs summarize the research
 884 structure followed by the selected literature on financial markets prediction using machine learning. This
 885 section describes in detail the steps to a robust research on predicting markets variables applying machine
 886 learning, the desired results and their validation. It aims to be a future reference for the best practices
 887 established by the most prominent papers in the field.

888 Most financial markets prediction papers first acknowledge the difficult task at hand. As commented
 889 on Sec. [1](#), financial markets prices are influenced by a myriad of factors and there has been a large number
 890 of proposals for their prediction. Among many markets values, researchers begin by setting their target
 891 for prediction, for example, returns and prices, as [Zhong and Enke \(2017\)](#) and [Chen et al. \(2017\)](#), direction,
 892 as [Ballings et al. \(2015\)](#) and [Patel et al. \(2015\)](#) or volatility, as [Oliveira et al. \(2017\)](#). The selected dependent
 893 variable to be predicted is generally some future, unknown value of a time series, supposedly useful for
 894 trading or hedging strategies. As examples, [Pan et al. \(2017\)](#) aim to predict USA stock index for weeks ahead
 895 and [Gorenc Novak and Velušček \(2016\)](#) set the next day prices of USA stocks as the dependent variables.

896 Although not common, present market values may also be chosen as variables of interest to forecast. For
 897 instance, [Patel et al. \(2015\)](#) and [Kara et al. \(2011\)](#) evaluate prediction models' performance on closing prices
 898 directions forecasting using those very same closing prices for calculating their independent variables. As
 899 argued in Sec. [5](#), for practical implementations, no models would be necessary to forecast market closing

900 direction if the closing prices are already available at the time of the forecast. Arguably, this approach is
 901 more suited to explore models' capabilities than to guide the building of profitable strategies. However,
 902 as financial markets future values prediction is the main interest of the academy and practitioners, best
 903 practices would dictate the choosing of a future market variable as the forecasting target.

904 According to the dependent variable chosen for prediction, a performance measure is more adequate.
 905 For instance, accuracy definitions, like F-measure, apply to works predicting prices direction classification
 906 (up or down), as [Weng et al. \(2017\)](#) and [Patel et al. \(2015\)](#), while error measurements, like RMSE or MAE,
 907 are more suited for prices and returns predictions, as applied by [Yan et al. \(2017\)](#) and [Göcken et al. \(2016\)](#).
 908 Overall, financial return may or may not be used as performance measure, because not all research focus on
 909 building trading strategies. In fact, some works classified on Tab. 17 propose and test models for predicting
 910 financial market variables, regardless as how those predictions should be used for profit. Some models
 911 obviously translate to trading strategies. For example, direction prediction can surely be interpreted as a
 912 buy/sell signals strategy, depending on the direction of the forecast. However, some authors refrain from
 913 advising a specific practical usage of their models, concentrating on the report of results using measures
 914 other than returns. In case financial returns are used to evaluate models and strategies, it is advisable to
 915 include other practical measures, such as drawdown and volatility as risk parameters.

916 After specifying the prediction target and the measurement of how close the results are from the real
 917 outcome, the researcher moves to the methods of prediction, arguably the richest and most creative part of
 918 the work. At this point, models are conceived, created, modified, tuned or even hybridized, as commented
 919 in Sec. 5. Tab. 18 and Tab. 19 bring the main machine learning algorithms and optimizations compiled by
 920 the papers considered in this review. They are only a small number of possibilities, although promising,
 921 for market forecasting. There is a large number of new methods being developed and studied, some tuned
 922 to other areas of forecasting, but plainly customizable for financial market time-series analysis. Visual
 923 recognition models, for example, could greatly enhance predictions based on given patterns, as in the
 924 work of [Chen and Chen \(2016\)](#). As another example, semantics and sentiment recognition innovations are
 925 potentially applicable to forecasting, as in [Oliveira et al. \(2017\)](#).

926 Machine learning model building does not necessarily involve the development of an algorithm entirely
 927 new. Customizing and tweaking well known models can lead to improved prediction results. Even pre-
 928 processing data before running the model is subject to research innovation, as exemplified by the PCA
 929 variations applied by [Zhong and Enke \(2017\)](#). Other examples of optimization techniques are found in Tab.
 930 19. As for model customizations examples, basic neural networks present many possibilities throughout
 931 specialized literature. For instance, [Yan et al. \(2017\)](#) combine traditional ANN with Bayesian probability
 932 theory and [Pei et al. \(2017\)](#) modify internal layers of neural networks with Legendre polynomials. SVM
 933 also presents opportunities for innovations through modifications of the traditional approach and the vast
 934 possibilities for kernel functions. Combining methods can also result in new and improved algorithms for

935 predictions, as exemplified by [Krauss et al. \(2017\)](#), [Zhang et al. \(2017\)](#) and [Chen et al. \(2017\)](#).

936 The next task at researching financial time series predictions using machine learning is obtaining data.
 937 However independent variables are chosen by the researcher, historical data must be gathered and somehow
 938 processed. All of the above reviewed papers apply their models to real, past data, being them stock prices
 939 ([Gorenc Novak and Velušček, 2016](#)), market indexes ([Pei et al., 2017](#)), fundamental company data ([Barak
 940 et al., 2017](#)) or text ([Oliveira et al., 2017](#)). Few authors rely exclusively on simulated data and, therefore,
 941 robust research in the area should always include real data. In fact, independent variables planning can be
 942 conducted simultaneously with data gathering, mainly because of availability concerns. Economic data,
 943 for example, are not frequently released, which may hinder higher frequency works, and involve reading
 944 or processing many company reports. High frequency prices, or even tick data, are commonly not freely
 945 available and its processing may require high computing power.

946 Data gathering may incur in high expenses and should be considered in the budget of the research. Some
 947 stock markets provide prices and transaction information free of costs. Brazilian BMF&Bovespa, for example,
 948 hosts a database of tick-by-tick level information, with some time constraints, free of costs at the Internet.
 949 Yahoo!Finance, a free online quotes database, is a common choice for data, as in [Weng et al. \(2017\)](#), [Bezerra
 950 and Albuquerque \(2017\)](#), [Chiang et al. \(2016\)](#) and [Gorenc Novak and Velušček \(2016\)](#). Another aspect that
 951 should be considered in selecting an information database is data quality. Few authors dedicate time to the
 952 treatment of data problems, such as outliers and missing data, as done by [Zhong and Enke \(2017\)](#). Data
 953 problems can be harmful to research results, specially for data-intense algorithms of machine learning, and
 954 should be dealt with by a well specified treatment process.

955 The variables to be used by machine learning models are intimately related to the available data. As
 956 exemplified by the papers listed in Tab. 17, predictive or independent variables may be chosen from TA,
 957 fundamentalist data, text from news, Internet blogs or prices themselves. Predictive variables choosing may
 958 in fact be the research question pursued by some works, as in [Chen et al. \(2017\)](#) and [Göçken et al. \(2016\)](#).
 959 Regardless of the type or number of independent variables selected, researchers must be cautious about
 960 only using data that would be available at the moment of the forecast. Making use of information prior to its
 961 release in the prediction procedure imply in a problem known as data snooping, which harms the validity of
 962 results. For example, using closing prices of a given period for calculating TA indicators and using them to
 963 predict the closing prices for that same period surely raises questions about validity, which can be observed
 964 in the work of [Patel et al. \(2015\)](#) and [Kara et al. \(2011\)](#). Results would be more robust if authors explicitly
 965 reported avoidance of snooping problems, as done by [Chen et al. \(2003, pp. 905–906\)](#), [Tsai and Hsiao \(2010,
 966 p. 264\)](#) and [Gerlein et al. \(2016, p. 198\)](#).

967 In machine learning applications, the whole dataset is divided into subsets for training and testing the
 968 chosen models. Training data, also called in-sample ([Krauss et al., 2017](#), p. 692), is commonly used in a
 969 supervised manner, that is, the real past outcomes are known by the models so they can be optimized to the

970 training data, expecting the testing data will hold the same underlying characteristics. A good practice is to
 971 avoid reusing data, which can lead to data snooping bias, by further dividing the training set into subgroups
 972 for model parametrization, as done by [Bezerra and Albuquerque \(2017\)](#). For testing purposes and measuring
 973 performance, test subset, called out-of-sample ([Huang et al., 2005](#), p. 2518) or holdout dataset ([Kim and](#)
 974 [Han, 2000](#), p. 192), is used. No observation from this set should be used in the training or optimization of
 975 the models, under the risk of snooping. Rigidly separating the subsets ensures validity of results, as the
 976 testing of the models are conducted on new data, unknown by the algorithms, simulating a real situation of
 977 forecast.

978 According to the reviewed literature of Tab. 17, training machine learning models is necessary to select
 979 the best parameters for future predictions. Values commented on Sec. 2, such as interconnections' weights
 980 in neural networks layers, SVM kernel function parameters and an optimal number of decision trees in a
 981 RF model, are obtained in the training phase. Assuming the test dataset behaves similarly to the training
 982 dataset, those optimized parameters are used on the models to predict the target variable using test data
 983 samples. The best parameters' values must always be chosen by in-sample training, without knowledge of
 984 test data, for validating results. That mimics practical implementations because out-of-sample data are not
 985 available on the moment of the forecast. Comparison results presented by [Kim \(2003\)](#), for example, may be
 986 confronted simply by observing that the author selects the best SVM parameters based on accuracy results
 987 from the test set data, incurring in another example of data snooping bias.

988 Another research decision is whether the machine learning model is periodically updated or it remains
 989 unchanged through the testing phase. It is not possible to observe a best practice in this regard, for both
 990 approaches are valid in the reviewed literature of Tab. 17. Keeping the models unchanged once they are
 991 optimized has the advantage of low computational costs. Once the models are trained, they are promptly
 992 used as each testing sample becomes available. [Krauss et al. \(2017\)](#), [Pan et al. \(2017\)](#), [Yan et al. \(2017\)](#), [Pei et al.](#)
 993 ([2017](#)) and [Bezerra and Albuquerque \(2017\)](#) are examples of papers which consider fixed, unchanged models
 994 once they are trained. On the other hand, the machine learning models may be periodically updated once
 995 new data become available, a procedure called sliding window ([Gerlein et al., 2016](#), p. 198). Computational
 996 costs are higher because the models have to be retrained every time the sliding window moves. However,
 997 they are constantly adapting to new market conditions. Examples of papers which follow this last approach
 998 are [Chiang et al. \(2016\)](#), [Gorenc Novak and Velušček \(2016\)](#), [Tsai and Hsiao \(2010\)](#), [Li and Kuo \(2008\)](#) and
 999 [Thawornwong and Enke \(2004\)](#).

1000 Finally, after comparing predictions with the real, test samples, the selected performance measures are
 1001 calculated. Some authors further apply statistical tests to evaluate the significance of results. For instance,
 1002 [Zhong and Enke \(2017\)](#) and [Kara et al. \(2011\)](#) apply t-tests while [Kim and Han \(2000\)](#) and [Yu et al. \(2009\)](#)
 1003 apply McNemar's tests. Relying on statistical tests well established in the scientific literature improves
 1004 robustness of the research, given the authors correctly interpret the results. Although many works still

¹⁰⁰⁵ report results without them, statistical tests should be incorporated as best practices in the field of machine
¹⁰⁰⁶ learning financial time series prediction for significance and robustness.

¹⁰⁰⁷ 7. Conclusions about the review

¹⁰⁰⁸ This article provided previously available quantitative and objective methods for selecting the relevant
¹⁰⁰⁹ literature on a particular theme of scientific research. The literature available regarding any topic can be
¹⁰¹⁰ broad, and a complete coverage of all the published documents can be challenging or even impossible. A
¹⁰¹¹ systematic selection of the literature most relevant to a study is thus necessary, taking into account not only
¹⁰¹² the history of an area but its state of the art. Thus, this review described bibliographic survey techniques and
¹⁰¹³ used them in the systematic review of financial market predictions that use machine learning techniques.
¹⁰¹⁴ This review led to the summarization of the best procedures established by the scientific literature on
¹⁰¹⁵ the field in order to achieve robust results when researching financial markets prediction using machine
¹⁰¹⁶ learning.

¹⁰¹⁷ As this literature review addressed machine learning techniques, Sec. 2 briefly commented on popular
¹⁰¹⁸ models; for example, ANN, SVM, and RF. In relation to the broader survey of the literature, Sec. 3 described
¹⁰¹⁹ how to proceed to a search of databases using keywords and filters by subject. It is emphasized that
¹⁰²⁰ the quality of this initial search for articles determined the final quality of the results obtained by the
¹⁰²¹ bibliographic survey. Sec. 4, in turn, presented the results of the database of articles surveyed, validating
¹⁰²² it objectively with the use of Lotka's distribution, in addition to analysing the results regarding the most
¹⁰²³ productive authors and countries, the most-cited periodicals, and the potential targets for submitting new
¹⁰²⁴ works for publication.

¹⁰²⁵ Moving on to the actual review of the most important articles about financial market prediction using
¹⁰²⁶ machine learning, this work commented on the following: the most-cited articles, those with the greatest
¹⁰²⁷ bibliometric coupling and the highest co-citation frequencies, the most recently published articles, and those
¹⁰²⁸ that are part of the main path of the knowledge flow of the literature studied. It is emphasized that these
¹⁰²⁹ were objective and clear methods of surveys, independent of the experience of the researcher, serving not
¹⁰³⁰ only for initial studies in research but also as validation of knowledge for experienced specialists.

¹⁰³¹ Finally, this work proposed a classification of the 57 articles reviewed, based on the markets addressed,
¹⁰³² the type of index predicted, the variables used as inputs for the models, and the type of prediction sought.
¹⁰³³ Additionally, the prediction methods used and the main performance measures used by each article were
¹⁰³⁴ summarized. There is extensive use of data from the North American market, in addition to the application
¹⁰³⁵ of neural and SVM networks. Similarly, most of the predictions relate to market indices. Among the possible
¹⁰³⁶ conclusions about the classification proposed here, it is to be expected that new proposed models will be
¹⁰³⁷ compared to the benchmarks of neural and SVM networks, as is the use of data from the North American

¹⁰³⁸ market. The use of new models in financial market prediction continues providing research opportunities,
¹⁰³⁹ as does the exploration of the behaviour of predictions in developing markets, such as those of the BRICS.

¹⁰⁴⁰ **References**

- ¹⁰⁴¹ Abu-Mostafa, Y. S., Atiya, A. F., 1996. [Introduction to financial forecasting](#). Applied Intelligence 6 (3), 205–213.
- ¹⁰⁴² Adya, M., Collopy, F., 1998. [How efective are neural networks at forecasting and prediction? A review and evaluation](#). Journal of Forecasting 17 (1), 481–495.
- ¹⁰⁴⁴ Al Nasseri, A., Tucker, A., de Cesare, S., 2015. [Quantifying StockTwits semantic terms' trading behavior in financial markets: An effective application of decision tree algorithms](#). Expert Systems with Applications 42 (23), 9192–9210.
- ¹⁰⁴⁶ Ang, K. K., Quek, C., 2006. [Stock trading using RSPOP: A novel rough set-based neuro-fuzzy approach](#). IEEE Transactions on Neural Networks 17 (5), 1301–1315.
- ¹⁰⁴⁸ Armano, G., Marchesi, M., Murru, A., 2005. [A hybrid genetic-neural architecture for stock indexes forecasting](#). Information Sciences 170 (1), 3–33.
- ¹⁰⁵⁰ Atsalakis, G. S., Valavanis, K. P., 2009. [Surveying stock market forecasting techniques—Part II: Soft computing methods](#). Expert Systems with Applications 36 (3), 5932–5941.
- ¹⁰⁵² Ballings, M., den Poel, D. V., Hespeels, N., Gryp, R., 2015. [Evaluating multiple classifiers for stock price direction prediction](#). Expert Systems with Applications 42 (20), 7046–7056.
- ¹⁰⁵⁴ Barak, S., Arjmand, A., Ortobelli, S., 2017. [Fusion of multiple diverse predictors in stock market](#). Information Fusion 36 (1), 90–102.
- ¹⁰⁵⁵ Batagelj, V., 2003. [Efficient algorithms for citation network analysis](#). arXiv preprint cs/0309023.
- ¹⁰⁵⁶ Bezerra, P. C. S., Albuquerque, P. H. M., 2017. [Volatility forecasting via SVR-GARCH with mixture of Gaussian kernels](#). Computational Management Science 14 (2), 179–196.
- ¹⁰⁵⁸ Bollerslev, T., 1986. [Generalized autoregressive conditional heteroskedasticity](#). Journal of Econometrics 31 (3), 307–327.
- ¹⁰⁵⁹ Box, G. E., Jenkins, G. M., Reinsel, G. C., Ljung, G. M., 2015. [Time series analysis: forecasting and control](#), 3rd Edition. Vol. 1. John Wiley & Sons, Hoboken, New Jersey.
- ¹⁰⁶¹ Breiman, L., 2001. [Random Forests](#). Machine Learning 45 (1), 5–32.
- ¹⁰⁶² Campbell, J. Y., 1987. [Stock returns and the term structure](#). Journal of Financial Economics 18 (2), 373–399.
- ¹⁰⁶³ Cao, L., 2003. [Support vector machines experts for time series forecasting](#). Neurocomputing 51 (1), 321–339.
- ¹⁰⁶⁴ Cao, Q., Leggio, K. B., Schniederjans, M. J., 2005. [A comparison between Fama and French's model and artificial neural networks in predicting the Chinese stock market](#). Computers & Operations Research 32 (10), 2499–2512.
- ¹⁰⁶⁶ Cavalcante, R. C., Brasileiro, R. C., Souza, V. L., Nobrega, J. P., Oliveira, A. L., 2016. [Computational intelligence and financial markets: A survey and future directions](#). Expert Systems with Applications 55 (1), 194–211.
- ¹⁰⁶⁸ Chang, P.-C., Fan, C.-Y., 2008. [A hybrid system integrating a wavelet and TSK fuzzy rules for stock price forecasting](#). IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews) 38 (6), 802–815.
- ¹⁰⁷⁰ Chang, P.-C., Liu, C.-H., Lin, J.-L., Fan, C.-Y., Ng, C. S., 2009. [A neural network with a case based dynamic window for stock trading prediction](#). Expert Systems with Applications 36 (3, Part 2), 6889–6898.
- ¹⁰⁷² Chen, A.-S., Leung, M. T., Daouk, H., 2003. [Application of neural networks to an emerging financial market: Forecasting and trading the Taiwan Stock Index](#). Computers & Operations Research 30 (6), 901–923.
- ¹⁰⁷⁴ Chen, H., Xiao, K., Sun, J., Wu, S., 2017. [A Double-Layer Neural Network Framework for High-Frequency Forecasting](#). ACM Transactions on Management Information Systems (TMIS) 7 (4), 11:2–11:17.
- ¹⁰⁷⁶ Chen, T.-l., Chen, F.-y., 2016. [An intelligent pattern recognition model for supporting investment decisions in stock market](#). Information Sciences 346 (1), 261–274.

- 1078 Chen, Y.-S., Cheng, C.-H., Tsai, W.-L., 2014. **Modeling fitting-function-based fuzzy time series patterns for evolving stock index**
 1079 **forecasting**. Applied Intelligence 41 (2), 327–347.
- 1080 Chiang, W.-C., Enke, D., Wu, T., Wang, R., 2016. **An adaptive stock index trading decision support system**. Expert Systems with
 1081 Applications 59 (1), 195–207.
- 1082 Chiu, S. L., 1994. **Fuzzy model identification based on cluster estimation**. Journal of Intelligent & Fuzzy Systems 2 (3), 267–278.
- 1083 Dash, R., 2017. Performance analysis of an evolutionary recurrent Legendre Polynomial Neural Network in application to FOREX
 1084 prediction. Journal of King Saud University-Computer and Information SciencesIn Press.
- 1085 Dietterich, T. G., 1998. **Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms**. Neural Computation
 1086 10 (7), 1895–1923.
- 1087 Donaldson, R. G., Kamstra, M., 1999. **Neural network forecast combining with interaction effects**. Journal of the Franklin Institute
 1088 336 (2), 227–236.
- 1089 Egghe, L., 2006. **Theory and practise of the g-index**. Scientometrics 69 (1), 131–152.
- 1090 Elman, J. L., 1990. **Finding structure in time**. Cognitive Science 14 (2), 179–211.
- 1091 Engle, R. F., 1982. **Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation**. Economet-
 1092 rica: Journal of the Econometric Society 50 (4), 987–1007.
- 1093 Enke, D., Thawornwong, S., 2005. **The use of data mining and neural networks for forecasting stock market returns**. Expert Systems
 1094 with Applications 29 (4), 927–940.
- 1095 Fama, E. F., 1991. **Efficient capital markets: II**. The Journal of Finance 46 (5), 1575–1617.
- 1096 Fernandez-Rodriguez, F., Gonzalez-Martel, C., Sosvilla-Rivero, S., 2000. **On the profitability of technical trading rules based on artificial**
 1097 **neural networks: Evidence from the Madrid stock market**. Economics Letters 69 (1), 89–94.
- 1098 Gerlein, E. A., McGinnity, M., Belatreche, A., Coleman, S., 2016. **Evaluating machine learning classification for financial trading: An**
 1099 **empirical approach**. Expert Systems with Applications 54 (1), 193–207.
- 1100 Göçken, M., Özçalıcı, M., Boru, A., Dosdoğru, A. T., 2016. **Integrating metaheuristics and Artificial Neural Networks for improved**
 1101 **stock price prediction**. Expert Systems with Applications 44 (1), 320–331.
- 1102 Gorenc Novak, M., Velušček, D., 2016. **Prediction of stock price movement based on daily high prices**. Quantitative Finance 16 (5),
 1103 793–826.
- 1104 Hájek, P., Olej, V., Myskova, R., 2013. **Forecasting stock prices using sentiment information in annual reports—a neural network and**
 1105 **support vector regression approach**. WSEAS Transactions on Business and Economics 10 (4), 293–305.
- 1106 Hassan, M. R., Nath, B., Kirley, M., 2007. **A fusion model of HMM, ANN and GA for stock market forecasting**. Expert Systems with
 1107 Applications 33 (1), 171–180.
- 1108 Henrique, B. M., Sobreiro, V. A., Kimura, H., 2018. **Building direct citation networks**. Scientometrics 115 (2), 817–832.
- 1109 Hornik, K., 1991. **Approximation capabilities of multilayer feedforward networks**. Neural Networks 4 (2), 251–257.
- 1110 Hornik, K., Stinchcombe, M., White, H., 1989. **Multilayer feedforward networks are universal approximators**. Neural Networks 2 (5),
 1111 359–366.
- 1112 Hsu, M.-W., Lessmann, S., Sung, M.-C., Ma, T., Johnson, J. E., 2016. **Bridging the divide in financial market forecasting: Machine learners**
 1113 **vs. financial economists**. Expert Systems with Applications 61 (1), 215–234.
- 1114 Huang, C.-L., Tsai, C.-Y., 2009. **A hybrid SOFM-SVR with a filter-based feature selection for stock market forecasting**. Expert Systems
 1115 with Applications 36 (2), 1529–1539.
- 1116 Huang, W., Nakamori, Y., Wang, S.-Y., 2005. **Forecasting stock market movement direction with support vector machine**. Computers &
 1117 Operations Research 32 (10), 2513–2522.
- 1118 Hummon, N. P., Doreian, P., 1989. **Connectivity in a citation network: The development of DNA theory**. Social Networks 11 (1), 39–63.
- 1119 Kamstra, M., Donaldson, G., 1996. **Forecasting combined with neural networks**. Journal of Forecast 15 (1), 49–61.
- 1120 Kara, Y., Boyacioglu, M. A., Baykan, Ö. K., 2011. **Predicting direction of stock price index movement using artificial neural networks**

- 1121 and support vector machines: The sample of the Istanbul Stock Exchange. Expert Systems with Applications 38 (5), 5311–5319.
- 1122 Kessler, M. M., 1963. Bibliographic coupling between scientific papers. Journal of the Association for Information Science and Technology
1123 14 (1), 10–25.
- 1124 Kim, K., 2003. Financial time series forecasting using support vector machines. Neurocomputing 55 (1-2), 307–319.
- 1125 Kim, K.-j., Han, I., 2000. Genetic algorithms approach to feature discretization in artificial neural networks for the prediction of stock
1126 price index. Expert Systems with Applications 19 (2), 125–132.
- 1127 Kimoto, T., Asakawa, K., Yoda, M., Takeoka, M., 1990. Stock market prediction system with modular neural networks. In: Neural
1128 Networks, 1990., 1990 IJCNN International Joint Conference on. IEEE, pp. 1–6.
- 1129 Krauss, C., Do, X. A., Huck, N., 2017. Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P
1130 500. European Journal of Operational Research 259 (2), 689–702.
- 1131 Kumar, D., Meghwani, S. S., Thakur, M., 2016. Proximal support vector machine based hybrid prediction models for trend forecasting
1132 in financial markets. Journal of Computational Science 17 (1), 1–13.
- 1133 Kumar, M., Thenmozhi, M., 2014. Forecasting stock index returns using ARIMA-SVM, ARIMA-ANN, and ARIMA-random forest
1134 hybrid models. International Journal of Banking, Accounting and Finance 5 (3), 284–308.
- 1135 Laboissiere, L. A., Fernandes, R. A., Lage, G. G., 2015. Maximum and minimum stock price forecasting of Brazilian power distribution
1136 companies based on artificial neural networks. Applied Soft Computing 35 (1), 66–74.
- 1137 Lage Junior, M., Godinho Filho, M., 2010. Variations of the kanban system: Literature review and classification. International Journal of
1138 Production Economics 125 (1), 13–21.
- 1139 Lahmiri, S., 2014a. Improving forecasting accuracy of the S&P500 intra-day price direction using both wavelet low and high frequency
1140 coefficients. Fluctuation and Noise Letters 13 (01), 1450008.
- 1141 Lahmiri, S., 2014b. Entropy-based technical analysis indicators selection for international stock markets fluctuations prediction using
1142 support vector machines. Fluctuation and Noise Letters 13 (02), 1450013.
- 1143 Lahmiri, S., Boukadoum, M., 2015. An Ensemble System Based on Hybrid EGARCH-ANN with Different Distributional Assumptions
1144 to Predict S&P 500 Intraday Volatility. Fluctuation and Noise Letters 14 (01), 1550001.
- 1145 Leigh, W., Purvis, R., Ragusa, J. M., 2002. Forecasting the NYSE composite index with technical analysis, pattern recognizer, neural
1146 network, and genetic algorithm: a case study in romantic decision support. Decision Support Systems 32 (4), 361–377.
- 1147 Leung, M. T., Daouk, H., Chen, A.-S., 2000. Forecasting stock indices: a comparison of classification and level estimation models.
1148 International Journal of Forecasting 16 (2), 173–190.
- 1149 Li, S.-T., Kuo, S.-C., 2008. Knowledge discovery in financial investment for forecasting and trading strategy through wavelet-based
1150 SOM networks. Expert Systems with Applications 34 (2), 935–951.
- 1151 Liu, J. S., Lu, L. Y., 2012. An integrated approach for main path analysis: Development of the Hirsch index as an example. Journal of the
1152 American Society for Information Science and Technology 63 (3), 528–542.
- 1153 Liu, J. S., Lu, L. Y., Lu, W.-M., Lin, B. J., 2013. Data envelopment analysis 1978–2010: A citation-based literature survey. Omega 41 (1),
1154 3–15.
- 1155 Malkiel, B. G., 2003. The Efficient Market Hypothesis and Its Critics. Journal of Economic Perspectives 17 (1), 59–82.
- 1156 Malkiel, B. G., Fama, E. F., 1970. Efficient Capital Markets: A Review of Theory and Empirical Work. The Journal of Finance 25 (2),
1157 383–417.
- 1158 Mariano, E. B., Sobreiro, V. A., Rebelatto, D. A. N., 2015. Human development and data envelopment analysis: A structured literature
1159 review. Omega 54 (1), 33–49.
- 1160 Mo, H., Wang, J., 2017. Return scaling cross-correlation forecasting by stochastic time strength neural network in financial market
1161 dynamics. Soft Computing 1 (1), 1–13.
- 1162 Nayak, R. K., Mishra, D., Rath, A. K., 2015. A Naïve SVM-KNN based stock market trend reversal analysis for Indian benchmark
1163 indices. Applied Soft Computing 35 (1), 670–680.

- ¹¹⁶⁴ Oliveira, N., Cortez, P., Areal, N., 2017. **The impact of microblogging data for stock market prediction: Using Twitter to predict returns, volatility, trading volume and survey sentiment indices.** Expert Systems with Applications 73 (1), 125–144.
- ¹¹⁶⁵
- ¹¹⁶⁶ Ortega, L., Khashanah, K., 2014. **A Neuro-wavelet Model for the Short-Term Forecasting of High-Frequency Time Series of Stock Returns.** Journal of Forecasting 33 (2), 134–146.
- ¹¹⁶⁷
- ¹¹⁶⁸ Pai, P.-F., Lin, C.-S., 2005. **A hybrid ARIMA and support vector machines model in stock price forecasting.** Omega 33 (6), 497–505.
- ¹¹⁶⁹ Pan, Y., Xiao, Z., Wang, X., Yang, D., 2017. **A multiple support vector machine approach to stock index forecasting with mixed frequency sampling.** Knowledge-Based Systems 122 (1), 90–102.
- ¹¹⁷⁰
- ¹¹⁷¹ Patel, J., Shah, S., Thakkar, P., Kotecha, K., 2015. **Predicting stock and stock price index movement using Trend Deterministic Data Preparation and machine learning techniques.** Expert Systems with Applications 42 (1), 259–268.
- ¹¹⁷²
- ¹¹⁷³ Pei, A., Wang, J., Fang, W., 2017. **Predicting agent-based financial time series model on lattice fractal with random Legendre neural network.** Soft Computing 21 (7), 1693–1708.
- ¹¹⁷⁴
- ¹¹⁷⁵ Rodríguez-González, A., García-Crespo, Á., Colomo-Palacios, R., Iglesias, F. G., Gómez-Berbís, J. M., 2011. **CAST: Using neural networks to improve trading systems based on technical analysis by means of the RSI financial indicator.** Expert systems with Applications 38 (9), 11489–11500.
- ¹¹⁷⁶
- ¹¹⁷⁷
- ¹¹⁷⁸ Saam, N., Reiter, L., 1999. **Lotka's law reconsidered: The evolution of publication and citation distributions in scientific fields.** Scientometrics 44 (2), 135–155.
- ¹¹⁷⁹
- ¹¹⁸⁰ Schumaker, R. P., Chen, H., 2009. **Textual analysis of stock market prediction using breaking financial news: The AZFin text system.** ACM Transactions on Information Systems (TOIS) 27 (2), 12.
- ¹¹⁸¹
- ¹¹⁸² Seuring, S., 2013. **A review of modeling approaches for sustainable supply chain management.** Decision Support Systems 54 (4), 1513 – 1520, rapid Modeling for Sustainability.
- ¹¹⁸³
- ¹¹⁸⁴ Small, H., 1973. **Co-citation in the scientific literature: A new measure of the relationship between two documents.** Journal of the Association for Information Science and Technology 24 (4), 265–269.
- ¹¹⁸⁵
- ¹¹⁸⁶ Son, Y., Noh, D.-j., Lee, J., 2012. **Forecasting trends of high-frequency KOSPI200 index data using learning classifiers.** Expert Systems with Applications 39 (14), 11607–11615.
- ¹¹⁸⁷
- ¹¹⁸⁸ Tay, F. E., Cao, L., 2001. **Application of support vector machines in financial time series forecasting.** Omega 29 (4), 309–317.
- ¹¹⁸⁹
- ¹¹⁹⁰ Thawornwong, S., Enke, D., 2004. **The adaptive selection of financial and economic variables for use with artificial neural networks.** Neurocomputing 56 (1), 205–232.
- ¹¹⁹¹
- ¹¹⁹² Thawornwong, S., Enke, D., Dagli, C., 2003. **Neural networks as a decision maker for stock trading: A technical analysis approach.** International Journal of Smart Engineering System Design 5 (4), 313–325.
- ¹¹⁹³
- ¹¹⁹⁴ Timmermann, A., Granger, C. W., 2004. **Efficient market hypothesis and forecasting.** International Journal of Forecasting 20 (1), 15–27.
- ¹¹⁹⁵
- ¹¹⁹⁶ Tsai, C.-F., Hsiao, Y.-C., 2010. **Combining multiple feature selection methods for stock prediction: Union, intersection, and multi-intersection approaches.** Decision Support Systems 50 (1), 258–269.
- ¹¹⁹⁷
- ¹¹⁹⁸ Tsaih, R., Hsu, Y., Lai, C. C., 1998. **Forecasting S&P 500 stock index futures with a hybrid AI system.** Decision Support Systems 23 (2), 161–174.
- ¹¹⁹⁹
- ¹²⁰⁰ Vapnik, V., 1995. **The nature of statistical learning theory.** Springer Heidelberg, New York, New York.
- ¹²⁰¹
- ¹²⁰² Wang, J.-J., Wang, J.-Z., Zhang, Z.-G., Guo, S.-P., 2012. **Stock index forecasting based on a hybrid model.** Omega 40 (6), 758–766.
- ¹²⁰³
- ¹²⁰⁴ Wang, Y.-F., 2002. **Predicting stock price using fuzzy grey prediction system.** Expert Systems with Applications 22 (1), 33–38.
- ¹²⁰⁵
- ¹²⁰⁶ Wang, Y.-F., 2003. **Mining stock price using fuzzy rough set system.** Expert Systems with Applications 24 (1), 13–23.
- ¹²⁰⁷
- ¹²⁰⁸ Weng, B., Ahmed, M. A., Megahed, F. M., 2017. **Stock market one-day ahead movement prediction using disparate data sources.** Expert Systems with Applications 79 (1), 153–163.
- ¹²⁰⁹
- ¹²¹⁰ Xiao, Y., Xiao, J., Lu, F., Wang, S., 2013. **Ensemble ANNs-PSO-GA Approach for Day-ahead Stock E-exchange Prices Forecasting.** International Journal of Computational Intelligence Systems 6 (1), 96–114.
- ¹²¹¹
- ¹²¹² Yan, D., Zhou, Q., Wang, J., Zhang, N., 2017. **Bayesian regularisation neural network based on artificial intelligence optimisation.**

- 1207 International Journal of Production Research 55 (8), 2266–2287.
- 1208 Yoon, Y., Swales Jr, G., Margavio, T. M., 1993. **A comparison of discriminant analysis versus artificial neural networks.** Journal of the
1209 Operational Research Society 44 (1), 51–60.
- 1210 Yu, L., Chen, H., Wang, S., Lai, K. K., 2009. **Evolving least squares support vector machines for stock market trend mining.** IEEE
1211 Transactions on Evolutionary Computation 13 (1), 87–102.
- 1212 Zhang, G., Patuwo, B. E., Hu, M. Y., 1998. **Forecasting with artificial neural networks: The state of the art.** International Journal of
1213 Forecasting 14 (1), 35–62.
- 1214 Zhang, N., Lin, A., Shang, P., 2017. **Multidimensional k-nearest neighbor model based on EEMD for financial time series forecasting.**
1215 Physica A: Statistical Mechanics and its Applications 477 (1), 161–173.
- 1216 Zhong, X., Enke, D., 2017. **Forecasting daily stock market return using dimensionality reduction.** Expert Systems with Applications
1217 67 (1), 126–139.
- 1218 Zweig, M. H., Campbell, G., 1993. **Receiver-operating characteristic (ROC) plots: A fundamental evaluation tool in clinical medicine.**
1219 Clinical Chemistry 39 (4), 561–577.

Algorithm 1: Algorithm for the construction of networks of direct citations.

```

1 Initialize list_articles;
2 for  $i \leftarrow 1$  to Total (list_articles) do
3   article  $\leftarrow$  list_articles[ $i$ ];
4   list_references  $\leftarrow$  References (article);
5   for  $j \leftarrow 1$  to Total (list_references) do
6     for  $k \leftarrow 1$  to Total (list_articles) do
7       if list_references[ $j$ ] = Title (list_articles[ $k$ ]) then
8         | Trace an edge starting at list_articles[ $k$ ] and finishing at list_articles[ $i$ ];
9       end
10      end
11    end
12 end

```

Algorithm 2: Algorithm for finding the main path of the literature via SPC.

```

1 Initialize list_sources;
2 Initialize list_sinks;
3 for  $i \leftarrow 1$  to Total (list_sources) do
4   source  $\leftarrow$  list_sources[ $i$ ];
5   for  $j \leftarrow 1$  to Total (list_sinks) do
6     sink  $\leftarrow$  list_sinks[ $j$ ];
7     paths  $\leftarrow$  Paths (source, sink);
8     foreach (path in paths) do
9       | Add 1 to the weight of each edge that is part of path;
10      end
11    end
12  end
13 Return the path between the source and the sink with the largest sum of weights.

```

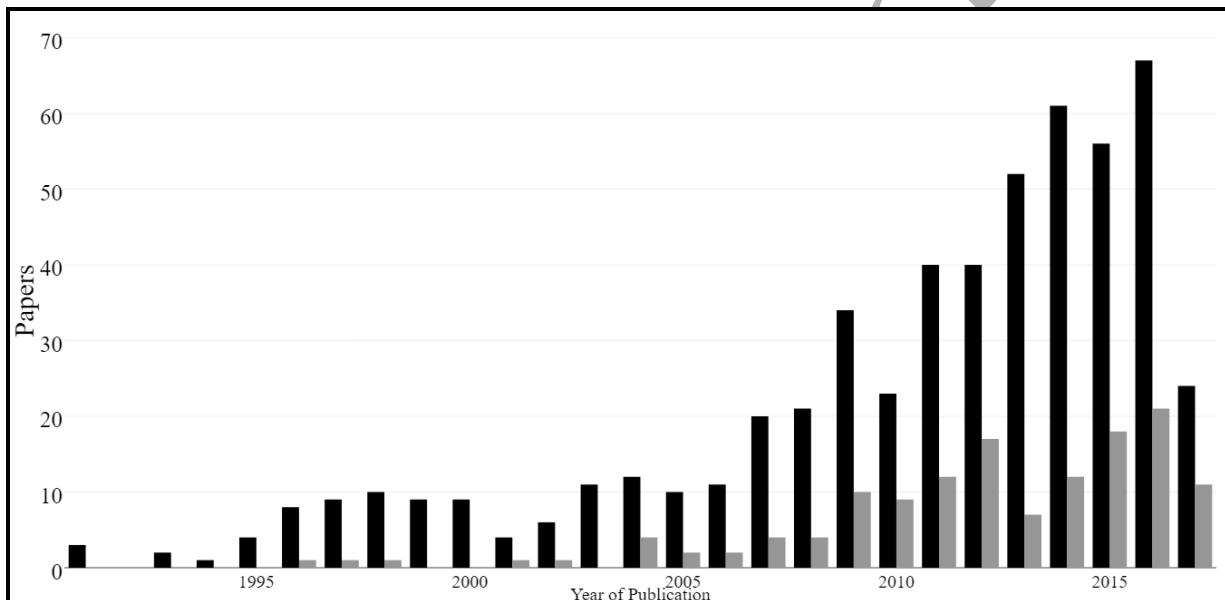


Figure 1: Frequency of publication of the articles considered in the bibliometrics, between 1991 and 2017, in dark bars. The lighter bars, for comparison, illustrate the publications about credit risk that apply machine learning.

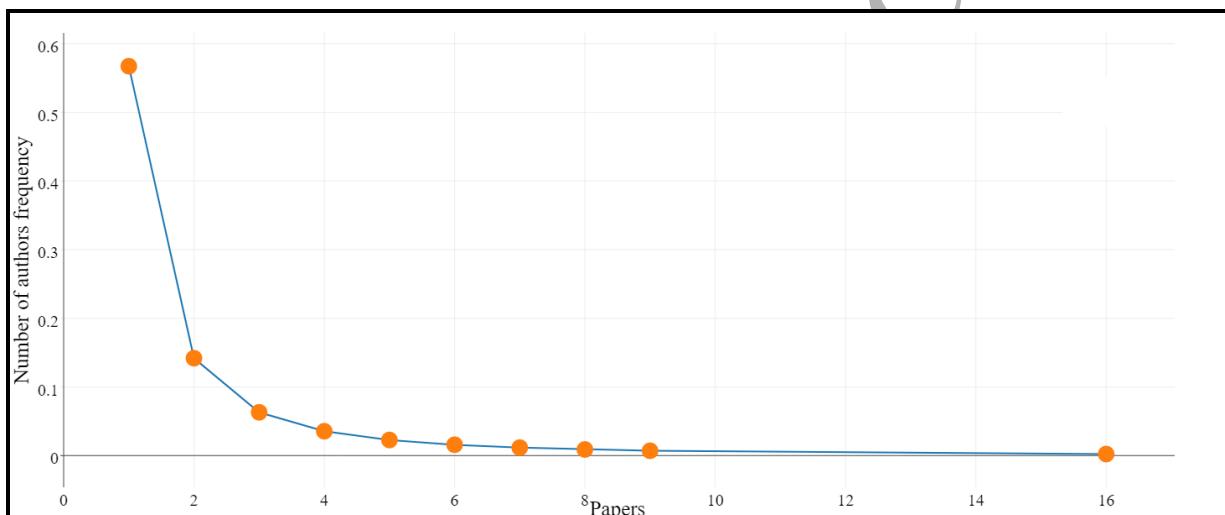


Figure 2: Frequency of publication, by author, in the database of articles considered in the bibliometrics. The continuous line represents Lotka's theoretical distribution, whereas the circles mark the distribution observed in the bibliographic survey of this present work.

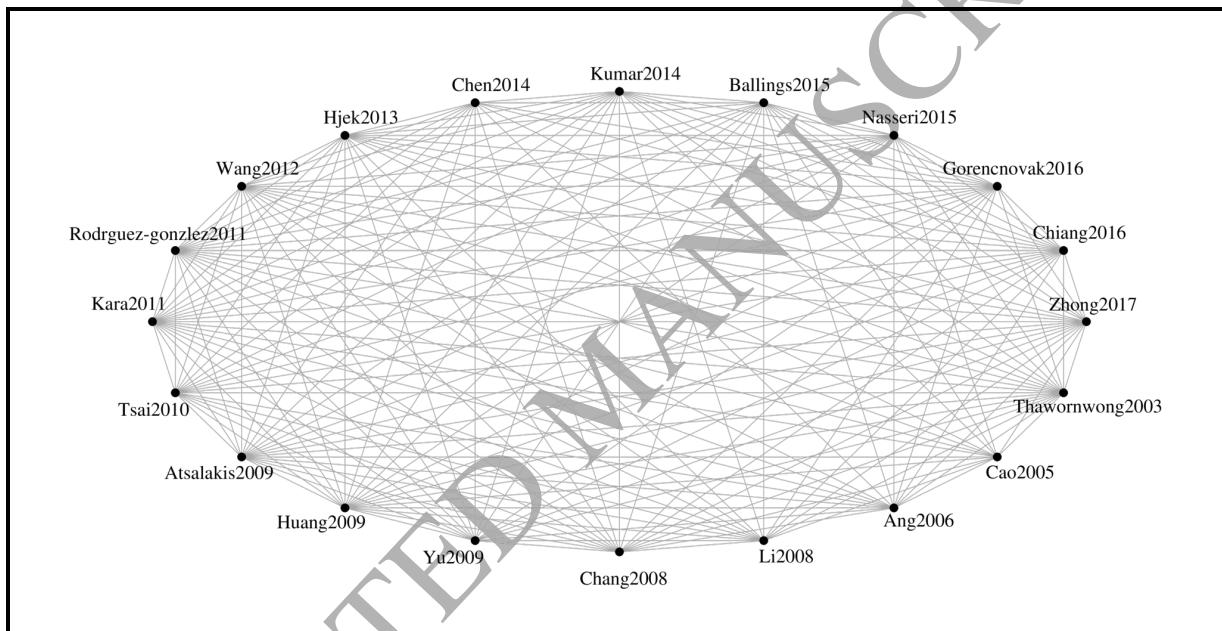


Figure 3: Bibliographic coupling of the 20 articles with the highest degrees of relationship. Each line represents a coupling relationship between the articles. The interconnections make a dense, almost fully-mashed network of papers with similar references.

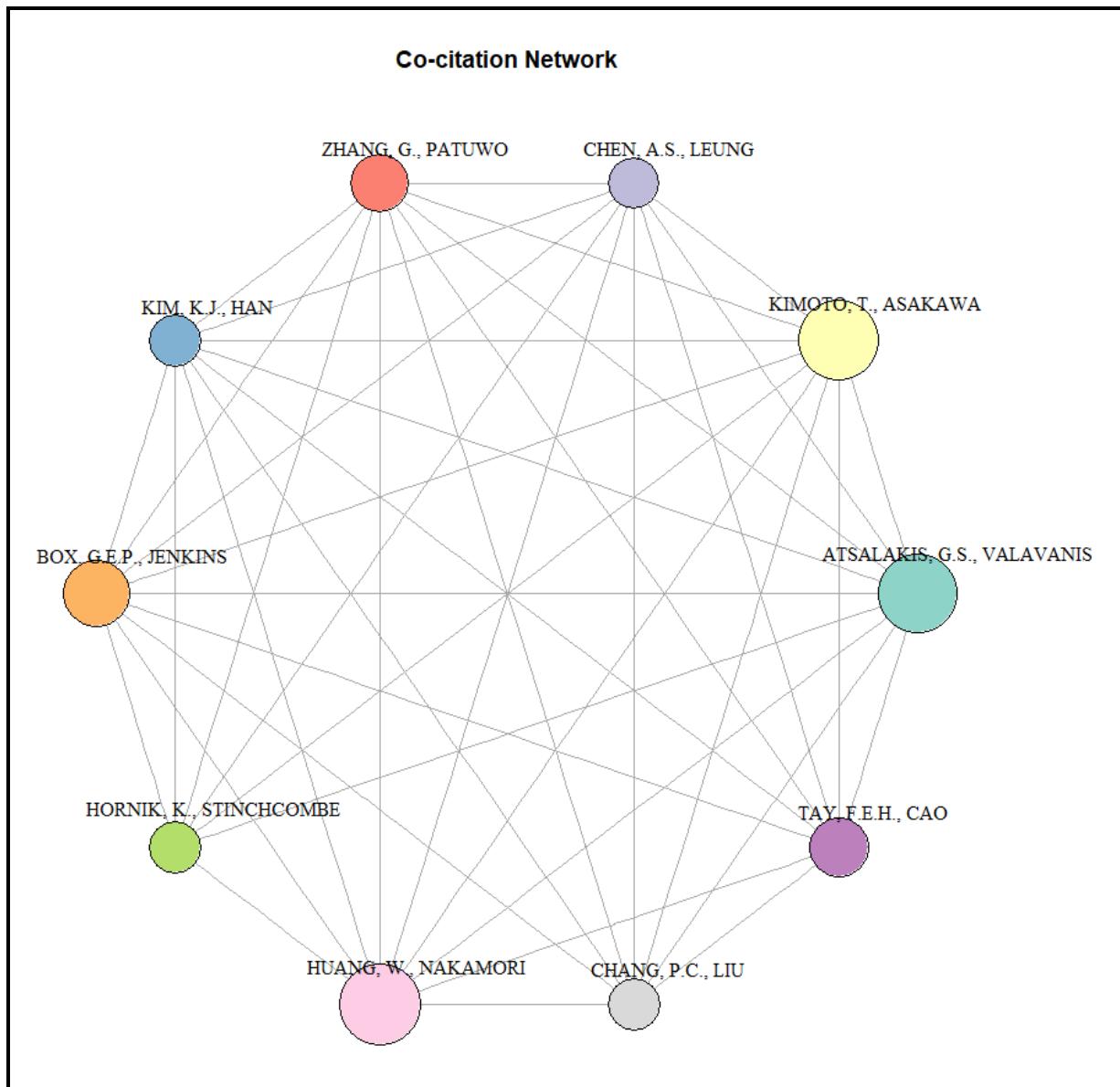


Figure 4: Network of co-citations for the 10 most-related authors and co-authors.

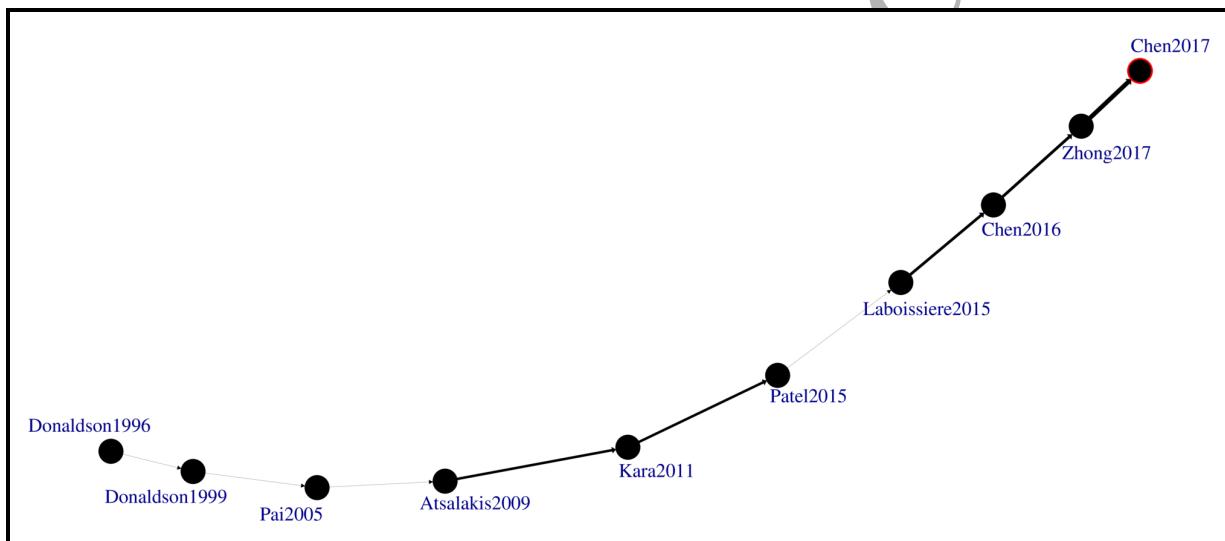


Figure 5: Main path followed by the literature. The edges have a thickness proportional to the weight assigned by Alg. 2.
Source: Henrique et al. (2018).

Characteristic	Value
Number of articles.	547
Periodicals.	243
Number of keywords.	1336
Period of the publications.	1991 – 2017
Average number of citations per article.	17.63
Authors.	1.151
Authors with single-author articles.	43
Articles per author.	0.475
Authors per article.	2.1

Table 1: Description of the database of articles used in the bibliometrics.

Author	Articles in the Database
Wang, J.	16
Cheng, C-H.	9
Wei, L-Y.	9
Enke, D.	8
Dash, PK.	7
Chang, P-C.	6
Chen, H.	6
Lahmiri, S.	6
Sun, J.	6
Dhar, J.	5
Hu, Y.	5
Kamstra, M.	5
Li, H.	5
Zhang, Y.	5
Zhang, Z.	5
Bekiros, SD.	4
Bisoi, R.	4
Chen, T-L.	4
Chen, Y.	4
Donaldson, RG.	4

Table 2: The 20 authors with the highest number of published articles in the database of articles searched.

Number of Articles	Authors	Distribution of Frequency
1	964	0.837533
2	126	0.109470
3	35	0.030408
4	11	0.009557
5	6	0.005213
6	4	0.003475
7	1	0.000869
8	1	0.000869
9	2	0.001738
16	1	0.000869

Table 3: Authors added to the most productive according to the h index.

Country	Number of Articles	Frequency
China.	88	0.16635
Taiwan.	72	0.13611
India.	60	0.11342
USA.	55	0.10397
Korea.	24	0.04537
Iran.	18	0.03403
United Kingdom.	16	0.03025
Spain.	15	0.02836
Greece.	14	0.02647
Singapore.	14	0.02647
Italy.	12	0.02268
Australia.	11	0.02079
Brazil.	11	0.02079
Hong Kong.	10	0.01890
Canada.	9	0.01701
Germany.	9	0.01701
Japan.	8	0.01512
Turkey.	8	0.01512
Lithuania.	6	0.01134
Malaysia.	5	0.00945

Table 4: The 20 countries with the highest number of articles according to the database of articles searched.

Author	<i>h</i> index	<i>g</i> index	Citations	Articles
Wang, J.	9	16	296	16
O, J.	6	11	122	13
Cheng, C-H.	6	9	243	9
Wei, L-Y.	7	9	231	9
Chen, Y.	3	9	86	9
Enke, D.	5	8	295	8
Chen, H.	5	8	438	8
Zhang, Z.	3	7	121	7
Dash, PK.	3	6	41	7
Chang, P-C.	6	6	214	6
Sun, J.	5	6	184	6
Wang, Y.	3	6	292	6
Kim, S.	5	6	89	6
Hu, Y.	3	5	30	5
Kamstra, M.	5	5	311	5
Li, H.	5	5	184	5
Zhang, Y.	2	5	60	5
Liu, M.	3	5	35	5
Bekiros, S.	3	5	27	5
Chen, T.	5	5	187	5

Table 5: Authors added to the most productive according to the *g* index.

Author	<i>h</i> index	<i>g</i> index	Citation	Articles
Chen, T-L.	4	4	182	4
Donaldson, RG.	4	4	270	4
Fan, C-Y.	4	4	177	4
Lu, C-J.	4	4	104	4
Quek, C.	4	4	101	4
Lin, C.	4	4	362	4
Lahmiri, S.	3	4	20	6

Table 6: Authors added to the most productive according to the *h* index.

Country	Number of Articles	Frequency
China.	88	0.16635
Taiwan.	72	0.13611
India.	60	0.11342
USA.	55	0.10397
Korea.	24	0.04537
Iran.	18	0.03403
United Kingdom.	16	0.03025
Spain.	15	0.02836
Greece.	14	0.02647
Singapore.	14	0.02647
Italy.	12	0.02268
Australia.	11	0.02079
Brazil.	11	0.02079
Hong Kong.	10	0.01890
Canada.	9	0.01701
Germany.	9	0.01701
Japan.	8	0.01512
Turkey.	8	0.01512
Lithuania.	6	0.01134
Malaysia.	5	0.00945

Table 7: The 20 countries with the highest number of articles according to the database of articles searched.

Country	Number of Citations	Mean Number of Citations per Article
Taiwan.	2429	33.736
USA.	1913	34.782
Korea.	1211	50.458
China.	900	10.227
Greece.	340	24.286
Singapore.	336	24.000
Canada.	282	31.333
India.	239	3.983
Spain.	234	15.600
Italy.	226	18.833
Australia.	190	17.273
United Kingdom.	173	10.812
Turkey.	152	19.000
Iran.	142	7.889
Brazil.	141	12.818
Germany.	133	14.778
Hong Kong.	106	10.600
Thailand.	49	9.800
Norway.	39	39.000
Slovenia.	35	17.500

Table 8: The 20 countries with the highest number of citations according to the database searched.

Periodical	Number of Articles
<i>Expert Systems with Applications.</i>	76
<i>Neurocomputing.</i>	21
<i>Applied Soft Computing Journal.</i>	20
<i>Decision Support Systems.</i>	14
<i>Neural Computing and Applications.</i>	11
<i>Neural Network World.</i>	11
<i>Journal of Forecasting.</i>	8
<i>Studies in Computational Intelligence.</i>	8
<i>International Journal of Applied Engineering Research.</i>	7
<i>Journal of Theoretical and Applied Information Technology.</i>	7
<i>Mathematical Problems in Engineering.</i>	7
<i>Soft Computing.</i>	7
<i>Information Sciences.</i>	6
<i>Journal of Information and Computational Science.</i>	6
<i>Knowledge-Based Systems.</i>	6
<i>Lecture Notes in Computer Science.</i>	6
<i>Physica A: Statistical Mechanics and its Applications.</i>	6
<i>Applied Intelligence.</i>	5
<i>Fluctuation and Noise Letters.</i>	5
<i>Computational Economics.</i>	4

Table 9: The 20 journals with the highest number of articles in the database searched.

Keyword	Articles that Use the Keyword
<i>Neural networks.</i>	59
<i>Forecasting.</i>	45
<i>Data mining.</i>	36
<i>Stock market.</i>	34
<i>Artificial neural networks.</i>	30
<i>Neural network.</i>	29
<i>Artificial neural network.</i>	28
<i>Prediction.</i>	25
<i>Genetic algorithm.</i>	22
<i>Machine learning.</i>	21
<i>Stock price forecasting.</i>	19
<i>Time series.</i>	19
<i>Feature selection.</i>	17
<i>Support vector machine.</i>	17
<i>Technical analysis.</i>	17
<i>Support vector machines.</i>	16
<i>Stock market prediction.</i>	15
<i>Support vector regression.</i>	15
<i>Stock prediction.</i>	14
<i>Genetic algorithms.</i>	13

Table 10: The 20 keywords most used in the database of articles searched.

References	Title	Journal	Citations	Citations per Year
Kim (2003).	<i>Financial Time Series Forecasting Using Support Vector Machines.</i>	<i>Neurocomputing.</i>	546	39.00
Kim and Han (2000).	<i>Genetic Algorithms Approach to Feature Discretization in Artificial Neural Networks for the Prediction of Stock Price Index.</i>	<i>Expert Systems with Applications.</i>	279	16.41
Pai and Lin (2005).	<i>A Hybrid ARIMA And Support Vector Machines Model in Stock Price Forecasting.</i>	<i>Omega.</i>	278	23.17
Atsalakis and Valavanis (2009).	<i>Surveying Stock Market Forecasting Techniques - Part II: Soft Computing Methods.</i>	<i>Expert Systems with Applications.</i>	224	28.00
Schumaker and Chen (2009).	<i>Textual Analysis of Stock Market Prediction Using Breaking Financial News: The AZFin Text System.</i>	<i>ACM Transactions on Information Systems.</i>	186	23.25
Chen et al. (2003).	<i>Application of Neural Networks to an Emerging Financial Market: Forecasting and Trading the Taiwan Stock Index.</i>	<i>Computers and Operations Research.</i>	184	13.14
Wang (2002).	<i>Predicting Stock Price Using Fuzzy Grey Prediction System.</i>	<i>Expert Systems with Applications.</i>	163	10.87
Enke and Thawornwong (2005).	<i>The Use of Data Mining and Neural Networks for Forecasting Stock Market Returns.</i>	<i>Expert Systems with Applications.</i>	152	12.67
Armano et al. (2005).	<i>A Hybrid Genetic-Neural Architecture for Stock Indexes Forecasting.</i>	<i>Information Sciences.</i>	130	10.83
Leigh et al. (2002).	<i>Forecasting the NYSE Composite Index with Technical Analysis, Pattern Recognizer, Neural Network, and Genetic Algorithm: A Case Study in Romantic Decision Support.</i>	<i>Decision Support Systems.</i>	130	8.67
Hassan et al. (2007).	<i>A Fusion Model Of HMM, ANN and GA for Stock Market Forecasting.</i>	<i>Expert Systems with Applications.</i>	124	12.40
Tsaih et al. (1998).	<i>Forecasting S&P 500 Stock Index Futures with a Hybrid AI System.</i>	<i>Decision Support Systems.</i>	121	6.37
Leung et al. (2000).	<i>Forecasting Stock Indices: A Comparison of Classification and Level Estimation Models.</i>	<i>International Journal of Forecasting.</i>	118	6.94
Wang (2003).	<i>Mining Stock Price Using Fuzzy Rough Set System.</i>	<i>Expert Systems with Applications</i>	117	8.36
Kamstra and Donaldson (1996).	<i>Forecast Combining with Neural Networks.</i>	<i>Journal of Forecasting</i>	115	5.48
Kara et al. (2011).	<i>Predicting Direction of Stock Price Index Movement Using Artificial Neural Networks and Support Vector Machines: The Sample of the Istanbul Stock Exchange.</i>	<i>Expert Systems with Applications.</i>	104	17.33
Yu et al. (2009).	<i>Evolving Least Squares Support Vector Machines for Stock Market Trend Mining.</i>	<i>IEEE Transactions on Evolutionary Computation.</i>	104	13.00
Fernandez-Rodriguez et al. (2000).	<i>On the Profitability of Technical Trading Rules Based on Artificial Neural Networks: Evidence from the Madrid Stock Market.</i>	<i>Economics Letters.</i>	100	5.88
Huang and Tsai (2009).	<i>A Hybrid SOFM-SVR with a Filter-Based Feature Selection for Stock Market Forecasting.</i>	<i>Expert Systems with Applications.</i>	99	12.38
Yoon et al. (1993).	<i>A Comparison of Discriminant Analysis Versus Artificial Neural Networks.</i>	<i>Journal of the Operational Research Society.</i>	99	4.12

Table 11: The 20 articles most cited in the compiled database. The number of citations refers to the citations in the entire Scopus database.

References	Title	Journal	Citations
Bollerslev (1986).	<i>Generalized Autoregressive Conditional Heteroscedasticity.</i>	<i>Journal of Econometrics.</i>	23
Kim (2003).	<i>Financial Time Series Forecasting Using Support Vector Machines.</i>	<i>Neurocomputing.</i>	18
Enke and Thawornwong (2005).	<i>The Use of Data Mining and Neural Networks for Forecasting Stock Market Returns.</i>	<i>Expert Systems with Applications.</i>	12
Elman (1990).	<i>Finding Structure in Time.</i>	<i>Cognitive Science.</i>	10
Engle (1982).	<i>Autoregressive Conditional Heteroscedasticity with Estimator of the Variance of United Kingdom Inflation.</i>	<i>Econometrica.</i>	9
Huang et al. (2005).	<i>Forecasting Stock Market Movement Direction with Support Vector Machine.</i>	<i>Computers and Operations Research.</i>	9
Thawornwong and Enke (2004).	<i>The Adaptive Selection of Financial and Economic Variables for Use with Artificial Neural Networks.</i>	<i>Neurocomputing.</i>	9
Campbell (1987).	<i>Stock Returns and the Term Structure.</i>	<i>Journal of Financial Economics.</i>	8
Chen et al. (2003).	<i>Application of Neural Networks to an Emerging Financial Market: Forecasting and Trading the Taiwan Stock Index.</i>	<i>Computers and Operations Research.</i>	8
Pai and Lin (2005).	<i>A Hybrid ARIMA and Support Vector Machines Model in Stock Price Forecasting.</i>	<i>Omega.</i>	8
Malkiel and Fama (1970).	<i>Efficient Capital Markets: A Review of Theory and Empirical Work.</i>	<i>Journal of Finance.</i>	7
Hornik et al. (1989).	<i>Multilayer Feedforward Networks are Universal Approximators.</i>	<i>Neural Networks</i>	7
Leung et al. (2000).	<i>Forecasting Stock Indices: A Comparison of Classification and Level Estimation Models.</i>	<i>International Journal of Forecasting.</i>	7
Tsaih et al. (1998).	<i>Forecasting S&P 500 Stock Index Futures with a Hybrid AI System.</i>	<i>Decision Support Systems.</i>	7
Zhang et al. (1998).	<i>Forecasting with Artificial Neural Networks: The State of the Art</i>	<i>International Journal of Forecasting.</i>	7
Abu-Mostafa and Atiya (1996).	<i>Introduction to Financial Forecasting.</i>	<i>Applied Intelligence.</i>	6
Adya and Collopy (1998).	<i>How Effective are Neural Networks at Forecasting and Prediction? A Review and Evaluation.</i>	<i>Journal of Forecasting.</i>	6
Atsalakis and Valavanis (2009).	<i>Surveying Stock Market Forecasting Techniques-Part II: Soft Computing Methods.</i>	<i>Expert Systems with Applications.</i>	6
Chiu (1994).	<i>Fuzzy Model Identification Based on Cluster Estimation.</i>	<i>Journal of Intelligent and Fuzzy Systems.</i>	6
Hornik (1991).	<i>Approximation Capabilities of Multilayer Feedforward Networks.</i>	<i>Neural Networks.</i>	6

Table 12: The 20 articles most cited by the articles in the compiled database.

Note: It should be noted that the articles in this table may not be part of the initially compiled database.

Reference	Title	Journal
Weng et al. (2017).	<i>Stock Market one-day ahead Movement Prediction using Disparate Data Sources.</i>	<i>Expert Systems with Applications.</i>
Zhang et al. (2017).	<i>Multidimensional k-Nearest Neighbor Model based on EEMD for Financial Time Series Forecasting.</i>	<i>Physica A: Statistical Mechanics and its Applications.</i>
Barak et al. (2017).	<i>Fusion of Multiple Diverse Predictors in Stock Market.</i>	<i>Information Fusion.</i>
Krauss et al. (2017).	<i>Deep Neural Networks, Gradient-boosted Trees, Random Forests: Statistical Arbitrage on the S&P 500.</i>	<i>European Journal of Operational Research.</i>
Oliveira et al. (2017).	<i>The Impact of Microblogging Data for Stock Market Prediction: Using Twitter to Predict Returns, Volatility, Trading Volume and Survey Sentiment Indices.</i>	<i>Expert Systems with Applications.</i>
Yan et al. (2017).	<i>Bayesian Regularisation Neural Network Based on Artificial Intelligence Optimisation.</i>	<i>International Journal of Production Research.</i>
Pan et al. (2017).	<i>A Multiple Support Vector Machine Approach to Stock Index Forecasting with Mixed Frequency Sampling.</i>	<i>Knowledge- Based Systems.</i>
Pei et al. (2017).	<i>Predicting Agent-based Financial Time Series Model on Lattice Fractal with Random Legendre Neural Network.</i>	<i>Soft Computing.</i>
Bezerra and Albuquerque (2017).	<i>Volatility Forecasting via SVR-GARCH With Mixture of Gaussian Kernels.</i>	<i>Computational Management Science.</i>
Mo and Wang (2017).	<i>Return Scaling Cross-Correlation Forecasting by Stochastic Time Strength Neural Network in Financial Market Dynamics.</i>	<i>Soft Computing.</i>

Table 13: The 10 most recent articles in the compiled database.

Reference	Title	Journal
Kumar and Thenmozhi (2014).	<i>Forecasting Stock Index Returns Using ARIMA-SVM, ARIMA-ANN, and ARIMA-Random Forest Hybrid Models.</i>	<i>International Journal of Banking, Accounting and Finance.</i>
Ballings et al. (2015).	<i>Evaluating Multiple Classifiers for Stock Price Direction Prediction.</i>	<i>Expert Systems with Applications.</i>
Al Nasseri et al. (2015).	<i>Quantifying Stocktwits Semantic Terms' Trading Behavior in Financial Markets: An Effective Application of Decision Tree Algorithms.</i>	<i>Expert Systems with Applications.</i>
Gorenc Novak and Velušček (2016).	<i>Prediction of Stock Price Movement Based on Daily High Prices.</i>	<i>Quantitative Finance.</i>
Chiang et al. (2016).	<i>An Adaptive Stock Index Trading Decision Support System.</i>	<i>Expert Systems with Applications</i>
Zhong and Enke (2017).	<i>Forecasting Daily Stock Market Return Using Dimensionality Reduction.</i>	<i>Expert Systems with Applications.</i>
Thawornwong et al. (2003).	<i>Neural Networks as a Decision Maker for Stock Trading: A Technical Analysis Approach.</i>	<i>International Journal of Smart Engineering System Design.</i>
Cao et al. (2005).	<i>A Comparison Between Fama and French's Model and Artificial Neural Networks in Predicting the Chinese Stock Market.</i>	<i>Computers and Operations Research.</i>
Ang and Quek (2006).	<i>Stock Trading Using RSPOP: A Novel Rough Set-Based Neuro-Fuzzy Approach.</i>	<i>IEEE Transactions on Neural Networks.</i>
Li and Kuo (2008).	<i>Knowledge Discovery in Financial Investment for Forecasting and Trading Strategy Through Wavelet-Based SOM Networks.</i>	<i>Expert Systems with Applications.</i>
Chang and Fan (2008).	<i>A Hybrid System Integrating a Wavelet and TSK Fuzzy Rules for Stock Price Forecasting.</i>	<i>IEEE Transactions on Systems, Man and Cybernetics Part C: Applications and Reviews.</i>
Yu et al. (2009).	<i>Evolving Least Squares Support Vector Machines for Stock Market Trend Mining.</i>	<i>IEEE Transactions on Evolutionary Computation.</i>
Huang and Tsai (2009).	<i>A Hybrid SOFM-SVR with a Filter-Based Feature Selection for Stock Market Forecasting.</i>	<i>Expert Systems with Applications.</i>
Atsalakis and Valavanis (2009).	<i>Surveying Stock Market Forecasting Techniques - Part II: Soft Computing Methods.</i>	<i>Expert Systems with Applications.</i>
Tsai and Hsiao (2010).	<i>Combining Multiple Feature Selection Methods for Stock Prediction: Union, Intersection, and Multi-Intersection Approaches.</i>	<i>Decision Support Systems.</i>
Kara et al. (2011).	<i>Predicting Direction of Stock Price Index Movement Using Artificial Neural Networks and Support Vector Machines: the Sample of the Istanbul Stock Exchange.</i>	<i>Expert Systems with Applications.</i>
Rodríguez-González et al. (2011).	<i>CAST: Using Neural Networks to Improve Trading Systems Based on Technical Analysis by Means of the RSI Financial Indicator.</i>	<i>Expert Systems with Applications.</i>
Wang et al. (2012).	<i>Stock Index Forecasting Based on a Hybrid Model.</i>	<i>Omega.</i>
Hájek et al. (2013).	<i>Forecasting Stock Prices Using Sentiment Information in Annual Reports - A Neural Network and Support Vector Regression Approach.</i>	<i>WSEAS Transactions on Business and Economics.</i>
Chen et al. (2014).	<i>Modeling Fitting-Function-Based Fuzzy Time Series Patterns for Evolving Stock Index Forecasting.</i>	<i>Applied Intelligence.</i>

Table 14: The 20 articles with the greatest bibliographic coupling among all the articles searched.

Reference	Title	Journal
Atsalakis and Valavanis (2009).	<i>Surveying Stock Market Forecasting Techniques—Part II: Soft Computing Methods.</i>	<i>Expert Systems with Applications.</i>
Box et al. (2015).	<i>Time Series Analysis: Forecasting and Control (Book).</i>	<i>John Wiley & Sons (Publisher).</i>
Chang et al. (2009).	<i>A Neural Network with a Case Based Dynamic Window for Stock Trading Prediction.</i>	<i>Expert Systems with Applications.</i>
Chen et al. (2003).	<i>Application of Neural Networks to an Emerging Financial Market: Forecasting and Trading the Taiwan Stock Index.</i>	<i>Computers & Operations Research.</i>
Hornik et al. (1989).	<i>Multilayer Feedforward Networks are Universal Approximators.</i>	<i>Neural Networks.</i>
Huang et al. (2005).	<i>Forecasting Stock Market Movement Direction with Support Vector Machine.</i>	<i>Computers & Operations Research.</i>
Kim and Han (2000).	<i>Genetic Algorithms Approach to Feature Discretization in Artificial Neural Networks for the Prediction of Stock Price Index.</i>	<i>Expert Systems with Applications.</i>
Kimoto et al. (1990).	<i>Stock Market Prediction System with Modular Neural Networks.</i>	<i>International Joint Conference on Neural Networks.</i>
Tay and Cao (2001).	<i>Application of Support Vector Machines in Financial Time Series Forecasting.</i>	<i>Omega.</i>
Zhang et al. (1998).	<i>Forecasting with Artificial Neural Networks: The State of the Art.</i>	<i>International Journal of Forecasting.</i>

Table 15: The 10 articles with the greatest co-citation relationship among those searched.

Reference	Title	Journal
Kamstra and Donaldson (1996).	<i>Forecast Combining with Neural Networks.</i>	<i>Journal of Forecasting.</i>
Donaldson and Kamstra (1999).	<i>Neural Network Forecast Combining with Interaction Effects.</i>	<i>Journal of the Franklin Institute.</i>
Pai and Lin (2005).	<i>A Hybrid ARIMA and Support Vector Machines Model in Stock Price Forecasting.</i>	<i>Omega.</i>
Atsalakis and Valavanis (2009).	<i>Surveying Stock Market Forecasting Techniques - Part II: Soft Computing Methods.</i>	<i>Expert Systems with Applications.</i>
Kara et al. (2011).	<i>Predicting Direction of Stock Price Index Movement Using Artificial Neural Networks and Support Vector Machines: The Sample of the Istanbul Stock Exchange.</i>	<i>Expert Systems with Applications.</i>
Patel et al. (2015).	<i>Predicting Stock and Stock Price Index Movement Using Trend Deterministic Data Preparation and Machine Learning Techniques.</i>	<i>Expert Systems with Applications.</i>
Laboissiere et al. (2015).	<i>Maximum and Minimum Stock Price Forecasting of Brazilian Power Distribution Companies Based on Artificial Neural Networks.</i>	<i>Applied Soft Computing Journal.</i>
Chen and Chen (2016).	<i>An Intelligent Pattern Recognition Model for Supporting Investment Decisions in Stock Market.</i>	<i>Information Sciences.</i>
Zhong and Enke (2017).	<i>Forecasting Daily Stock Market Return Using Dimensionality Reduction.</i>	<i>Expert Systems with Applications.</i>
Chen et al. (2017).	<i>A Feature Weighted Support Vector Machine and K- Nearest Neighbor Algorithm for Stock Market Indices Prediction.</i>	<i>Expert Systems with Applications.</i>

Table 16: Articles that are part of the main path of the literature searched.

Reference	Market/s	Asset/s	Predictive Variable/s	Prediction/s	Main Method/s	Performance Measure/s
Al Nasseri et al. (2015).	USA.	Index.	Text.	Direction.	Analysis of sentiment.	Return.
Ang and Quek (2006).	Singapore.	Stocks.	TA.	Prices.	Neural networks.	Return.
Armano et al. (2005).	USA, Italy.	Indices.	TA.	Prices.	Neural networks, GA.	Sharpe rate.
Ballings et al. (2015).	Europe.	Stocks.	Fundamentalist.	Direction.	Neural networks, SVM, kNN, and RF.	AUC.
Barak et al. (2017).	Iran	Stocks.	Fundamentalist.	Return and risk.	Neural networks, SVM, decision trees.	Accuracy.
Bezerra and Albuquerque (2017).	Brazil, Japan.	Indices.	Prices.	Volatility.	SVR, GARCH.	MAE, RMSE.
Cao et al. (2005).	China.	Stocks.	Fundamentalist.	Return.	Neural networks, CAPM.	MAD, MAPE, MSE.
Chang and Fan (2008).	Taiwan.	Index.	TA.	Prices.	kNN, DWT, fuzzy logic.	MAPE.
Chang et al. (2009).	Taiwan.	Stocks.	TA.	Direction.	Neural networks, CBR.	Return.
Chen et al. (2003).	Taiwan.	Index.	Fundamentalist.	Return.	Neural networks, GMM.	Return.
Chen et al. (2014).	Taiwan, Hong Kong.	Indices.	TA.	Prices.	Fuzzy logic.	RMSE.
Chen and Chen (2016).	USA, Taiwan.	Indices.	TA.	Return.	Pattern recognition.	Return.
Chen et al. (2017).	China.	Indices.	TA.	Direction.	SVM, kNN.	MAPE, RMSE, AUC.
Chiang et al. (2016).	Multiple.	Indices.	TA.	Direction.	Neural networks.	Accuracy, return.
Enke and Thawornwong (2005).	USA.	Index.	Fundamentalist.	Direction.	Neural networks.	RMSE.
Fernandez-Rodriguez et al. (2000).	Spain.	Index.	Prices.	Direction.	Neural networks.	Accuracy, Sharpe rate.
Gorenc Novak and Velušček (2016).	USA.	Stocks.	TA.	Direction.	SVM.	Return, Sharpe rate.
Hájek et al. (2013).	USA.	Stocks.	Fundamentalist.	Return.	Neural networks, SVR, analysis of sentiment.	MSE.
Hassan et al. (2007).	USA.	Stocks.	Prices.	Prices.	Neural networks, GA.	MAPE.

to be continued.

Reference	Market/s	Asset/s	Predictive Variable/s	Prediction/s	Main Method/s	Performance Measure/s
Huang and Tsai (2009).	Taiwan.	Index.	TA.	Prices.	SVR.	MSE, MAE, MAPE.
Kara et al. (2011).	Turkey.	Index.	TA.	Direction.	Neural networks, SVM.	Accuracy.
Kamstra and Donaldson (1996).	Multiple.	Indices.	Prices.	Volatility.	Neural Networks.	MSE, MAE.
Kim and Han (2000).	Korea.	Index.	TA.	Direction.	Neural networks, GA.	Accuracy.
Kimoto et al. (1990).	Japan.	Index.	Fundamentalist.	Direction.	Neural networks.	MAPE.
Krauss et al. (2017).	USA.	Stocks.	Prices.	Returns.	Neural networks, RF, decision trees.	Return, Sharpe rate.
Laboissiere et al. (2015).	Brazil.	Stocks.	Indices.	Maximums, minimums.	Neural networks.	MAE, MAPE, RMSE.
Leigh et al. (2002).	USA.	Index.	Prices, volume.	Prices.	Neural networks, GA.	Return.
Leung et al. (2000).	USA, United Kingdom, Japan.	Indices.	Fundamentalist.	Return.	Neural networks, LDA, regressions.	Return.
Li and Kuo (2008).	Taiwan.	Indices.	Prices.	Prices.	DWT, SOM.	MSE, MAE.
Mo and Wang (2017).	China, USA.	Indices.	Prices.	Correlation.	Neural networks.	MAE, RMSE, MAPE.
Oliveira et al. (2017).	Multiple.	Indices.	Text.	Return, volume, volatility.	Neural networks, SVM, RF.	MAE.
Pai and Lin (2005).	USA.	Stocks.	Prices.	Prices.	SVM.	MAE, MAPE, MSE, RMSE.
Pan et al. (2017).	USA.	Index.	Fundamentalist, prices.	Prices.	SVM.	RMSE, MAE.
Patel et al. (2015).	India.	Indices, stocks.	TA.	Direction.	Neural networks, SVM, RF, NB.	Accuracy.
Pei et al. (2017).	China.	Index.	Prices.	Mean of the prices.	Neural networks.	RMSE, MAE, MAPE.
Rodríguez-González et al. (2011).	Spain.	Index, stocks.	TA.	Direction.	Neural networks.	Accuracy.
Schumaker and Chen (2009).	USA.	Stocks.	News.	Prices.	Textual analysis, SVR.	MSE.
Thawornwong et al. (2003).	USA.	Stocks.	TA.	Direction.	Neural networks.	Return, Sharpe rate.
Tsai and Hsiao (2010).	Taiwan.	Stocks.	Fundamentalist.	Direction.	Neural networks.	Accuracy.
Tsaih et al. (1998).	USA.	Index.	TA.	Direction.	Neural networks.	Accuracy.

to be continued.

Reference	Market/s	Asset/s	Predictive Variable/s	Prediction/s	Main Method/s	Performance Measure/s
Wang (2002).	Taiwan.	Stocks.	Prices.	Prices.	Fuzzy logic.	Accuracy.
Wang (2003).	Taiwan.	Stocks.	Prices.	Prices.	Fuzzy logic.	Accuracy.
Wang et al. (2012).	China, USA.	Indices.	Prices.	Prices.	Neural networks, GA.	Accuracy, MAE, RMSE, MAPE.
Weng et al. (2017).	USA.	Stocks.	TA, text.	Direction.	Neural networks, SVM, decision trees.	Accuracy, AUC, F-measure.
Yan et al. (2017).	China.	Index.	Prices.	Prices.	Neural networks.	MAE, MAPE, MSE.
Yoon et al. (1993).	USA.	Stocks.	Fundamentalist.	Return.	Neural networks, LDA.	Accuracy.
Yu et al. (2009).	USA.	Indices.	Fundamentalist, TA.	Direction.	SVM, GA.	Accuracy.
Zhang et al. (2017).	USA.	Indices.	Prices.	Prices.	kNN, ARIMA.	MAPE, MSE.
Zhong and Enke (2017).	USA.	Index.	Fundamentalist.	Direction.	Neural networks.	MSE, Sharpe rate.
Abu-Mostafa and Atiya (1996).	FOREX.	Currency.	Fundamentalist, "hints", TA.	Direction.	Neural networks.	Return.
Donaldson and Kamstra (1999).	USA.	Index.	Prices.	Return.	Neural networks, GARCH.	RMSE, MAE.
Enke and Thawornwong (2005).	USA.	Index.	Fundamentalist.	Direction.	Neural networks.	RMSE.
Huang et al. (2005).	USA, Japan.	Indices, currency.	Indices, currency	Direction.	SVM, neural networks, LDA.	Accuracy.
Kim (2003).	Korea.	Index.	TA.	Direction.	SVM, neural networks.	Accuracy.
Kumar and Thenmozhi (2014).	India.	Index.	Returns.	Return.	Neural networks, SVM, RF, ARIMA.	Accuracy, MAE, RMSE.
Tay and Cao (2001).	USA.	Indices.	Prices, TA.	Prices.	Neural networks, SVM.	MAE, MSE.
Thawornwong and Enke (2004).	USA.	Index.	Fundamentalist.	Direction.	Neural networks.	RMSE.

Table 17: Classification of the reviewed articles about financial market prediction using machine learning techniques.

Note: AUC: Area Under the receiver operating characteristic Curve; GMM: Generalized Methods of Moments; MAD: Mean Absolute Deviation.

Main Method	Number of References	References
Neural Networks.	42	Ang and Quek (2006), Armano et al. (2005), Ballings et al. (2015), Barak et al. (2017), Cao et al. (2005), Chang et al. (2009), Chen et al. (2003), Chiang et al. (2016), Enke and Thawornwong (2005), Fernandez-Rodriguez et al. (2000), Hájek et al. (2013), Hassan et al. (2007), Kara et al. (2011), Kamstra and Donaldson (1996), Kim and Han (2000), Kimoto et al. (1990), Krauss et al. (2017), Laboissiere et al. (2015), Leigh et al. (2002), Leung et al. (2000), Mo and Wang (2017), Oliveira et al. (2017), Patel et al. (2015), Pei et al. (2017), Rodríguez-González et al. (2011), Thawornwong et al. (2003), Tsai and Hsiao (2010), Tsaih et al. (1998), Wang et al. (2012), Weng et al. (2017), Yan et al. (2017), Yoon et al. (1993), Zhong and Enke (2017), Abu-Mostafa and Atiya (1996), Donaldson and Kamstra (1999), Enke and Thawornwong (2005), Huang et al. (2005), Kim (2003), Kumar and Thenmozhi (2014), Tay and Cao (2001), Thawornwong and Enke (2004), Lahmiri (2014a), Lahmiri and Boukadoum (2015).
SVM/SVR.	20	Ballings et al. (2015), Barak et al. (2017), Bezerra and Albuquerque (2017), Chen et al. (2017), Gorenc Novak and Velišek (2016), Hájek et al. (2013), Huang and Tsai (2009), Kara et al. (2011), Oliveira et al. (2017), Pai and Lin (2005), Pan et al. (2017), Patel et al. (2015), Schumaker and Chen (2009), Weng et al. (2017), Yu et al. (2009), Huang et al. (2005), Kim (2003), Kumar and Thenmozhi (2014), Tay and Cao (2001), Lahmiri (2014b).
RF/Decision Trees.	7	Ballings et al. (2015), Barak et al. (2017), Krauss et al. (2017), Oliveira et al. (2017), Patel et al. (2015), Weng et al. (2017), Kumar and Thenmozhi (2014).
Sentiment/Text Analysis.	5	Al Nasseri et al. (2015), Hájek et al. (2013), Schumaker and Chen (2009), Weng et al. (2017), Oliveira et al. (2017).
kNN.	4	Ballings et al. (2015), Chang and Fan (2008), Chen et al. (2017), Zhang et al. (2017).
ARIMA/GARCH.	4	Bezerra and Albuquerque (2017), Donaldson and Kamstra (1999), Zhang et al. (2017), Kumar and Thenmozhi (2014).
Fuzzy Logic.	4	Chang and Fan (2008), Chen et al. (2014), Wang (2002), Wang (2003).
LDA.	3	Leung et al. (2000), Yoon et al. (1993), Huang et al. (2005).
NB.	1	Patel et al. (2015).

Table 18: Main forecasting techniques applied by each reviewed reference.

Main Method	Number of References	References
GA.	11	Armano et al. (2005), Bezerra and Albuquerque (2017), Hassan et al. (2007), Kim and Han (2000), Leigh et al. (2002), Wang et al. (2012), Yu et al. (2009), Donaldson and Kamstra (1999), Göçken et al. (2016), Chen and Chen (2016), Tsai and Hsiao (2010).
DWT.	6	Li and Kuo (2008), Chang and Fan (2008), Chiang et al. (2016), Ortega and Khashanah (2014), Xiao et al. (2013), Lahmiri (2014a).
PCA.	3	Zhong and Enke (2017), Son et al. (2012), Tsai and Hsiao (2010).
PSO.	3	Yan et al. (2017), Chiang et al. (2016), Xiao et al. (2013).
SOM.	2	Li and Kuo (2008), Cao (2003).

Table 19: Main optimization techniques applied by each reviewed reference.

Author Contributions Section

Bruno Miranda Henrique (BMH), Vinicius Amorim Sobreiro (VAS) and Herbert Kimura (HK) participated in the development of the research. The first author conducted the study and the results were discussed initially with VAS and HK. Following the three authors developed the initial version of the manuscript. Then, VAS revised and improvement in the paper. Finally, all authors read and approved the final manuscript.

ACCEPTED MANUSCRIPT