# A Review on Machine Learning for Asset Management

**Pedro M. Mirete-Ferrer** 1,2,*,† 🆔, **Alberto Garcia-Garcia** 3,† 🆔, **Juan Samuel Baixauli-Soler** 4 🆔 and **Maria A. Prats** 4 🆔

1 Escuela Internacional de Doctorado Universidad de Murcia, Interuniversity Doctorate in Economics (DEcIDE), 30100 Murcia, Spain

2 Faculty of Economics and Business Administration (ICADE), Universidad Pontificia Comillas, 28015 Madrid, Spain

3 Departamento de Tecnología Informática y Computación, Universidad de Alicante, 03690 Alicante, Spain; agarcia@dtic.ua.es

4 Facultad de Economía y Empresa, Universidad de Murcia, 30100 Murcia, Spain; samuel@um.es (J.S.B.-S.); mprats@um.es (M.A.P.)

* Correspondence: pmmirete@comillas.edu or pedromanuel.miretef@um.es

† These authors contributed equally to this work.

**Abstract:** This paper provides a review on machine learning methods applied to the asset management discipline. Firstly, we describe the theoretical background of both machine learning and finance that will be needed to understand the reviewed methods. Next, the main datasets and sources of data are exposed to help researchers decide which are the best ones to suit their targets. After that, the existing methods are reviewed, highlighting their contribution and significance in the analyzed financial disciplines. Furthermore, we also describe the most common performance criteria that are applied to compare such methods quantitatively. Finally, we carry out a critical analysis to discuss the current state-of-the-art and lay down a set of future research directions.

## 1. Introduction

During the last 70 years, financial economists over the world have dedicated great efforts to model and forecast stock returns, trying to understand the patterns behind them. This has been the greatest challenge of the research focused on Asset Management over the last few decades. As Cochrane (2011) exposed, "In the beginning, there was chaos; practitioners thought one only needed to be clever to earn high returns. Then came the CAPM. Every clever strategy to deliver high average returns ended up delivering high market betas as well. Then anomalies erupted, and there was chaos again. The "value effect" was the most prominent anomaly". In the last 10 years, the huge amount of anomalies found has driven the academic professionals to call the phenomenon as a "factor zoo".

In contrast to this "anomaly-challenging" branch of literature, a growing amount of work indicates remarkable investment performance based on signals generated by various Machine Learning (ML) methods. With the recent advancement in financial technology (Fintech), there is an increasing trend of employing ML technologies to find new signals on price movements and build investment systems that can beat human fund managers from the perspective of practical investment management. ML routines in academic research have been implicitly motivated by the American Finance Association (AFA) presidential address of Cochrane (2011). This author suggested that, in the presence of a large number of noisy and highly correlated return predictors, there is a need for other methods beyond cross-sectional regressions and portfolio sorts. Indeed, ML uses "regularization" approaches to choose models, moderate over-fitting biases, find complicated patterns and hidden linkages, and handle high-dimensional predictor sets and more flexible functional forms, as exposed by Gu et al. (2020).

To illustrate the gap that still exists between academic finance and the financial industry, with regards to the interest of the academic finance community for ML techniques,

Huck (2019) carried out a very interesting experiment using the academic papers database: until 2017, the search for "machine/deep learning" via the EBSCO database, produces no reference at all in "The Journal of Finance", the leading academic finance journal, and only one reference in "The Journal of Financial Economics". If we expand the analysis to some other academic financial or economic journals, using an additional financial term as "stock returns", there can be found 21 references, mainly in "Quantitative Finance". Other 64 references can be placed in non-financial journals. In the last four years, although the interest of academic finance for ML techniques has apparently grown due to the movement initiated by econometricians and statiscians, as mentioned by Heaton et al. (2016), the result of this experiment is not so different, giving the result of 21 references within financial or economic journals, and a hundred additional references outside them.

As Heaton et al. (2017) points out, the main difference between the use of ML tools in Finance and other areas of Science is that, in Finance, "the emphasis is not in replicating tasks that humans already do well. Unlike recognizing an image or responding appropriately to verbal requests, humans have no innate ability to, for example, select a stock that is likely to perform well in some future period". For this reason, the usefulness of ML tools for financial purposes should be searched somewhere else. Specifically, they are extremely powerful in selection problems since, at their basis, they are the best and most rapid way to compute any function mapping data, and that is what returns, prices, economic data, accounting data, and so forth.

Our work aims to find a balance between finance, statistics, and computational disciplines by improving the analysis of recent literature from the standpoint of financial economics using ML approaches, allowing academics and practitioners to locate their areas of interest without gaps. We will formulate the asset management problem and break it down into disciplines, providing the reader with enough background knowledge to either get familiar with the area or acquire a standard terminology. The range of fields of our paper will span all the topics related to asset/portfolio management, from asset pricing and factor investing, more linked to economic and financial variables, to price forecasting and algorithmic trading, more concentrated on price and volume data. Finally, we will discuss the reviewed methods and datasets and provide useful insight in shape of future research directions and open challenges in the field.

The paper is organized as follows. In Section 2, we provide an overview of the most recent works of literature review in the field of ML applications to Finance. Furthermore, we justify the usefulness of our approach versus the rest of review papers. In Section 3, we expose the methodology and terminology to be used in the rest of the paper, as well as a review of the basic theoretical background about Asset Management and ML techniques. Section 4 gives a description of the different datasets used in the research. Section 5 provides a detailed description of the current state of the art in the application of ML to Asset Management, using a double outlook to classify the recent literature: the financial field of application and the ML empirical approach used. A discussion upon the evidence presented in previous sections is presented in Section 6. Lastly, we conclude in Section 7 with the final remarks.

## 2. Related Works

In this section, we will review the recent literature that, in form of survey papers, analyse the recent contributions regarding ML applications to Finance in general (General Surveys) and, later, to specific areas of Finance regarding Asset Management: price forecasting, asset pricing and portfolio management.

Finance is an extremely diverse field in Economics, that includes so diverse disciplines as Asset and Portfolio Management, Risk Assessment, Fraud Detection or Financial Regulation, among many others. The use of machine learning techniques in all these fields, in the recent years, has been increasingly relevant.

In the sample selected in this work, spanning the last five years, we have been able to find fifteen papers, as it is shown in Table 1. From general to specific, the compilation

papers adopting the scope of general applications of machine learning to finance, in its widest extension, is relatively popular, with four papers in the last five years. Among the specific review papers, which are focused on some of the application areas of Finance, the discipline related to price forecasting and time-series prediction counts for more than a half of the total, with eight papers. We can also find one general review about asset pricing and value investing. And finally, the asset and portfolio management discipline and its applications has a main role in two papers of revision, one of them exclusively focused on online portfolio selection. In most cases (12 out of 15), as it can be seen in the second column, the publishing journals were computer science outlets.

**Table 1.** List of review articles in the last five years, both general and specific financial areas. General Survey (GS), Price Forecasting (PF), Asset Pricing (AP), Portfolio Management (PM), Stock Market (SM), Exchange Rates forecasting (FX), Interest Rate forecasting (IR), Criptocurrencies (CC), Commodity Prices (CP), Derivatives (DV), Real Estate (RE). Applied Soft Computing Journal (1), Expert Systems with Applications (2), Artificial Intelligence Review (3), Frontiers of Business Research in China (4), International Journal of Electricity and Computer Engineering (5), International Journal on Emerging Technologies (6), Financial Markets and Portfolio Management (7), ACM Computing Surveys (8), ICEFR 2019 (9). (∗) for financial journals or conferences.

| Author | Journal/Venue | Period | Citations | References | Area | Market |
|---|---|---|---|---|---|---|
| Ozbayoglu et al. (2020) | 1 | 2014–2018 | 44 | 196 | GS | SM, CC, DV |
| Huang et al. (2020) | 4 * | 2014–2018 | 13 | 51 | GS | SM, FX, CP |
| Tkáč and Verner (2016) | 1 | 1994–2015 | 184 | 425 | GS | SM, FX, IR, DV |
| Cavalcante et al. (2016) | 2 | 2009–2015 | 282 | 144 | GS | SM, FX, DV |
| Sezer et al. (2020) | 1 | 2005–2019 | 162 | 216 | PF | SM, FX, CP |
| Henrique et al. (2019) | 2 | 1991–2017 | 101 | 98 | PF | SM |
| Bustos and P.-Quimbaya (2020) | 2 | 2014–2018 | 24 | 87 | PF | SM |
| Kamley et al. (2016) | 5 | 2000–2015 | 13 | 68 | PF | SM |
| Jiang (2021) | 2 | 2017–2019 | 39 | 234 | PF | SM |
| Nti et al. (2019) | 3 | 2017–2019 | 36 | 207 | PF | SM |
| Xing et al. (2018) | 3 | 1998–2016 | 157 | 153 | PF | SM |
| Durairaj and Mohan (2019) | 6 | 1999–2019 | 15 | 46 | PF | SM, FX, IR, DV |
| Weigand (2019) | 7 * | 1994–2018 | 8 | 49 | AP | SM, IR, DV, RE |
| Emerson et al. (2019) | 9 * | 2015–2018 | 1 | 81 | PM | SM |
| Li and Hoi (2014) | 8 | 1991–2013 | 137 | 246 | PM | SM |

## 2.1. General Surveys

A general survey about the recent applications of Deep Learning in Finance was conducted in Ozbayoglu et al. (2020). The authors chose to separate the most significant section of computational intelligence for finance research, which is focused on financial time-series forecasting, from this analysis. Their findings show that financial text mining, algorithmic trading, risk assessment, sentiment analysis, portfolio management and fraud detection are among the most-studied areas of finance research, and that Recurrent Neural Network (RNN)-based models (in particular, Long-short Term Memory (LSTM)), Convolutional Neural Networks (CNNs) and Deep Multi-layer Perceptron (DMLP) have been used extensively in implementations.

We can find that another recent paper, Huang et al. (2020), also focused on Deep Learning techniques with a general perspective of Finance and Banking. They defined seven domains of financial applications, which are: exchange rate prediction, stock market prediction, stock trading, banking default risk and credit, portfolio management, macroeconomic prediction, and oil price prediction. In their classification there is a predominance of price forecasting fields, and there exists a mixture of financial and banking disciplines.

The last two references that can be considered as review papers from a general financial perspective that were written in 2016. The paper by Tkáč and Verner (2016) is generally from the point of view of the fields of application, but specific in the methods used by the papers reviewed, that must have been focused on neural networks. More than 400 articles are classified according to the year of publication, application area, type of neural network, learning algorithm, benchmark method, citations, and journals. Lastly,

in Cavalcante et al. (2016), the financial review of the literature was organized according to their primary goal or main contribution to computational intelligence applied to financial markets literature. Following this criteria, papers are classified in application fields such as preprocessing mechanisms, forecasting models, and text mining and forecasting mechanisms.

*2.2. Price Forecasting*

As we previously mentioned, price forecasting has been the most prolific area of application of machine learning algorithms in the finance industry over the last few years. We have been able to find and classify six review papers fully devoted to price forecasting techniques, most of them focused on the stock market. As Ozbayoglu et al. (2020) describes, the most substantial portion of computational intelligence for finance research is dedicated to financial time-series forecasting, and this statement can be fully transferred to the field of review papers.

From the same authors, there exists a review paper fully focused on price forecasting Sezer et al. (2020). Their motivation was to provide a comprehensive literature review on Deep Learning (DL) studies for financial time-series forecasting implementations. They not only categorized the studies according to their intended forecasting implementation areas, such as index, forex, and commodity forecasting, but also grouped them based on their DL model choices, such as CNNs, Deep Belief Networks (DBNs), and LSTM. They found that RNN-based models (in particular, LSTM) are the most commonly used models, whilst Natural Language Processing (NLP), semantics and text mining-based hybrid models ensembled with time-series data might be more common in the near future.

The survey paper by Henrique et al. (2019) proposes the use of bibliographic survey techniques to classify around 60 papers about this issue. The authors reach two relevant conclusions: (1) the prevalence of Support Vector Machines (SVMs) and Neural Networks (NNs) as the most commonly used techniques, and (2) the relevance that the use of data from developing markets might have in the future. A more updated, but similar approach, comes from Bustos and P.-Quimbaya (2020), which aims to perform an updated systematic review of the forecasting techniques used in the stock market. The paper includes a long list of review papers in Finance, where the relevance of text mining techniques in this financial field emerges clearly. It addresses the classification of papers according to two criteria: the type of inputs (market information, technical indicator, economic indicators) and the technique or method used in the research.

Additionally involved in the stock price forecasting domain, we can find Kamley et al. (2016), which provides an overview of the machine learning techniques that have been used to forecast equity performance. The most distinct contribution of this paper has to do with the procedure to know how the prediction algorithms can be used to identify the most important variables in an equity market data set. The methods reviewed from the literature are Decision Trees (DTs), NNs, SVM, Genetic Algorithms (GAs), and Bayesian Networks (BNs). More recent are the contributions from Jiang (2021), specifically in the subfield of deep learning techniques, not only by categorizing the different data sources, various neural network structures, and commonly used evaluation metrics, but also their implementation and reproducibility. Finally, we can find Nti et al. (2019), where authors focused on fundamental and technical analyses, finding that SVM and Artificial Neural Network (ANN) are the most popular techniques of ML for stock market prediction.

The last two surveys about price forecasting are clearly more specific. In Xing et al. (2018), we can find a study about NLP, a pragmatic research viewpoint of computer linguistics that has grown in capability as a result of data availability and different methodologies developed over the last decade, and the subfield which aims to predict financial markets, the Natural Language-based Financial Forecasting (NLFF). This work includes recent attention in markets to sentiment analysis or event extraction. In Durairaj and Mohan (2019), a review of deep learning hybrids for financial time-series prediction can be found. According to the authors' definition, a hybrid forecasting model combines two or more

stand-alone forecasting models into a combined model in the hope to improve prediction accuracy and overcome the deficiencies of stand-alone models.

### 2.3. Value/Factor Investing

Accurate pricing or valuation of an asset can be considered as a fundamental area of research in Finance. As we will describe in Section 5, the asset pricing models have become one of the more prolific areas of study in Finance within the last 50 years. Although it is not a simple task to draw a division line between this area of study and others, mainly asset and portfolio management, we thought that it would be convenient to establish a separate field for all those papers and techniques which are aimed at finding the best representation models and forecasts for asset prices and returns (asset pricing models), but not involved in the process of combining those return forecasts with risk models in order to create an optimum portfolio, given an investor's constraints (portfolio management).

According to our classification criteria, only one paper can be considered as a literature review in the field of asset pricing/value investing. It is the work by Weigand (2019), where the author suggests that ML can aid future study, but that academics should be skeptical of these approaches because ML has flaws and is still relatively new to asset pricing. This literature review summarizes the research in fields like equities, global equity indices, derivatives, real estate appraisals, and bonds, as well as mortgage delinquencies. The authors conclude that these ML approaches show potential in overcoming deficiencies of traditional methods in empirical asset pricing, characterized by the general problems of prediction, variable selection, as well as the definition of a suitable functional form.

### 2.4. Portfolio Management

Finally, some of the survey papers about ML applied to Finance in the last few years are devoted to portfolio optimization. Despite the relevance that this discipline has reached in recent years in the banking industry, we have been able to find, in the last five years, only one survey, and we had to expand our focus and go back seven years to find the second. In Emerson et al. (2019), the authors evaluated current and potential applications of ML to the investment process. In particular, they included the development of ML applications for return forecasting, portfolio construction, and risk modelling, which are framed under the general category of Quantitative Finance and, more specifically, the investment process. They identified 67 papers from different perspectives, focused on those three core areas regarding the investment process. The analysis and discussion of the recent literature about them was then made individually for each of the previously identified papers.

The second reference about portfolio management, previously mentioned, was written seven years ago. In Li and Hoi (2014), the authors provided a comprehensive survey of online portfolio algorithms selection, the asset management field that sequentially selects a portfolio over a set of assets in order to achieve certain targets. The authors concentrate on a survey of multiperiod/sequential portfolio selection work due to the sequential (online) character of this assignment, in which the aim is to maximize the expected log return across a sequence of trading periods, and the portfolio is rebalanced to a set allocation at the conclusion of each trading session.

### 2.5. Our Proposal in Context

After analyzing the different survey papers written in the last five years, we have reached the following conclusions:

- Finance is an extremely diverse field in Economics, that includes such diverse disciplines as Asset/Portfolio Management, Risk Assessment, Fraud Detection or Financial Regulation. The use of ML techniques in all these fields, in the recent years, has been increasingly relevant. Most of the literature reviews have tried to span every single discipline within this extensive and heterogeneous group, so that the analysis has been revealed, in some occasions, as relatively superficial. In some other cases, the classification has been quite arguable because of the dominance of price forecasting fields,

and the mixture of financial and banking disciplines. Maybe due to this extensive strategy adopted by their peers, some authors have adopted the reverse approach, focusing the revision on a single ML methodology, achieving an arguable result extensive in financial disciplines but excessively focused on a single ML technique.

- Within the financial disciplines related to asset management, price forecasting has been the favourite field of review papers in the last five years. Favored by this wide coverage, researchers and practitioners that are interested in this topic can decide which path they should take according to currently existing literature.
- In terms of popularity, asset pricing is the antithesis of price forecasting. Only one review paper Weigand (2019) is concentrated on this financial field, which can be considered as the conceptual basis of the rest of the disciplines.
- Similarly, portfolio management reviews are clearly underrated. Just one paper Emerson et al. (2019) in the last five years can be considered as a general survey about this topic.

According to previous analysis points, we think our paper bridges a gap in the area of review papers and makes the following useful contributions:

- It covers a very understated field in the surveys' literature, despite the growing relevance that this financial discipline has reached in recent years in the banking industry. The range of fields of our paper spans all the topics related to asset/portfolio management, from asset pricing and factor investing, more linked to economic and financial variables, to price forecasting and algorithmic trading, more concentrated on price and volume data.
- Our paper tries to find an equilibrium point between finance, statistics and computational fields, enhancing the analysis of the recent literature from the perspective of financial economics, as well as ML methods, so that researchers and practitioners can find their areas of interest without gaps. We dedicate a section to analyzing the quality and homogeneity of datasets, a somewhat neglected aspect in finance, unlike other scientific fields. Similarly, we make an analysis of the different approaches adopted in the recent literature regarding the performance criteria.
- Beyond the gathering, processing, and classification of papers and information, our work gives responses to several research questions, such as areas of interest to the financial and ML community, degree of maturity of the existing research in each of the application areas, areas with more promising potential from academic and industrial perspectives, and suggestions about future directions for ML research in asset/portfolio management.

## 3. Theoretical Background

This section includes all the methodology and terminology to be used in the rest of the paper. ML applications in Finance is a field that can be included in the discipline of Quantitative Finance. On the one hand, it includes a review of the basic financial theoretical background, from the definitions of the different financial disciplines to the empirical asset pricing models, passing by the well-known theoretical asset pricing models as Capital Asset Pricing Model (CAPM), APT, and so forth. On the other hand, this section introduces the main theoretical principles about ML, mainly those connected with asset and portfolio management in a wide sense. Finally, it includes a review and description of the main indicators of performance used in the literature of ML applications to Asset Management, very related to reproducibility issues.

### 3.1. Financial Background

According to Snow (2020), the Asset Management discipline can be divided "into four streams: portfolio construction, risk management, capital management, infrastructure and deployment, and sales and marketing". We will focus our work on the first stream, portfolio construction, which is intimately tied with the investment process. Portfolio

construction is divided, in turn, into four areas: price forecasting, event prediction, value investing, and weight optimization.

The first three areas can be included in the field of trading strategies, but they differ in the type of data used and the outcome they are trying to predict. Price strategies include technical analysis, and macro global and statistical arbitrage (also called pairs trading), and the price plays a starring and lonely role. Right in the opposite corner, we can place value investing strategies, where the relationship between value and price is what generates the investment opportunities. Examples of these types of strategies are risk parity, factor investing, and fundamental investing. Finally, event strategies can be considered as part of a relatively new stream: an event-driven strategy refers to an investment strategy in which an investor attempts to profit from a stock mispricing that may occur during or after a corporate event. The theoretical basis of all of them, one way or another, can be found in the Asset Pricing theory, although in the case of price strategies, the relevance of the concept of value decays notably.

The fourth area, weight optimization, can be considered independently from the other three areas. It comprises the use of mathematical or statistical techniques to solve optimization and simulation problems in finance, like optimal execution, position sizing, and portfolio optimization.

Due to the extreme complexity and variety of disciplines included in the term of Asset Management, we have decided to define our own structure of disciplines to face the gathering, processing, classification and analysis of recent literature. We will divide the asset management disciplines into three main streams: value investing and price forecasting, as trading strategies topics, and portfolio management as an independent area which involves optimization. Value investing will comprise all the disciplines which use asset pricing models to select the most valuable assets to invest in. The most relevant and recent example of this type of discipline is factor investing. Price forecasting, on the other hand, will comprise all the financial areas focused on the best prediction of asset prices. We will open a special category in Section 5 for algorithmic trading, the most relevant area of price strategies in recent literature. Lastly, following the proposal from Snow (2020), weight optimization will be defined as an independent discipline that, for our purposes, will be categorized as portfolio management.

In short, these will be the three financial disciplines that will be used in the rest of the article:

1. Value/factor investing. Investment strategies which use asset pricing models to select the most valuable assets to invest in.
2. Price Forecasting. Investment strategies focused on the best prediction of asset prices. Algorithmic trading can be considered as a special case of price strategy.
3. Portfolio Management. Mathematical and statistical techniques which solve optimization and simulation problems in investment management.

### 3.1.1. Value/Factor Investing

Over the past three decades, hundreds of financial research articles have been dedicated to the study of asset pricing models with a dual purpose: on the one hand, to analyze the behavior of asset prices and, on the other hand, to try to find variables that contain information about them. The advances in financial research about asset pricing and quantitative methods about factor investing have been crucial for the exceptional development of the Asset Management industry in general and the enhancement of the investment processes, specifically.

Throughout this section, we will analyze both theoretical and empirical models. This body of knowledge is usually applied to cross-sectional data, that is, data from different financial assets at the same period of time. When this occurs, we will use the sub-index $i$. When we apply the models to data which vary over time, that is, time-series data, we will use the sub-index $t$. The coexistence of both dimensions, cross-section and time-series, is traditionally solved in Econometrics through panel data methodology but, specifically

in factor investing, this double dimension is faced through the methodology of Fama and Macbeth (1973).

SDF as General Source of Asset Pricing Models

Regardless of the type of asset pricing models we use, all of them can be deduced from the general Stochastic Discount Factor (SDF) Model, also known as the Euler Equation. The SDF, also denominated as the pricing kernel, allows one to relate the current price of an asset to its future payoffs. The roots of this type of representation are based on the Arrow–Debreu model of general equilibrium, and its application to option pricing (Cox and Ross (1976) and Ross (1978)), along with the Asset Pricing Theory (APT) of Ross (1976).

Using the same notations as Cochrane (2000), $x_{t+1}$ will represent an asset pay-off at date $t + 1$, $p_t$ the asset price at date $t$, and $m_{t+1}$ the SDF at date $t + 1$. The fundamental equality states that:

$$P_t = E_t(m_{t+1}X_{t+1}) \tag{1}$$

Equation (1) indicates that the asset price is equal to the conditional expected value of the future payoff multiplied by the SDF, which is a random variable whose realizations are always greater than zero. The expected price can be understood, hence, as a weighted average of future payoffs in different states of the economy, where each state has a probability of occurrence defined by its SDF. Consequently, the SDF is simply a discount factor that transforms expected payoffs tomorrow into value today, but in a world of uncertainty; if there is no uncertainty or if investors are risk-neutral, the SDF is merely a constant that transforms expected payoffs tomorrow into value today. [1]

The asset pricing model which allows to make a direct translation of the SDF architecture into a model with observable variables is a consumption-based Model where the SDF is derived from the optimality conditions for a single agent. The very well-known Consumption-based Capital Asset Pricing Model (CCAPM) from Breeden (1979) (see Section 3.1.1) can be regarded as an equilibrium version with multiple investors of this seminal model. [2]

The major contributions of the stochastic discount factor approach are its simplicity and universality. Instead of using three apparently different theories for different types of assets, for example, bonds, stocks, and derivatives, we can consider them as just special cases of the same theory.

Theoretical Factor Models

In Cochrane (2000), the author demonstrated that formulating a factor model in terms of the SDF is equal to a "beta representation" of expected returns, which is more common across different factor model formulations. Therefore, SDF model can be translated into a factor model, that is, a linear function of some risk factors, which discounts uncertain payoffs differently across different states of the world. Factor models correct the limitations of the unobservable SDF and the restrictions of a consumption-based approach, bearing in mind that consumption is a macroeconomic aggregate that is provided with lags and is continuously revised.

A factor model suggests the existence of a scalar $a$ and a K × 1 vector $b$ as loadings for a K × 1 vector $f$ of factor values, such that the quantity described by:

$$m = a + b_i^T f_i \qquad i = 1, 2, ..., K \tag{2}$$

is an SDF. A factor $f_k$ (the $k$-th factor) that has a non-zero loading $b_k$ is known as a "pricing factor".

Let us assume that $f$ is a K × 1 random vector. Thus, a scalar $a$ and a K × 1 vector $b$ exist such that Equation (2) prices all assets if, and only if, a scalar $\kappa$ and a K × 1 vector $\Lambda$ exist such that for each asset $i$, the expected return $R_i$ is:

$$E[R_i] = \kappa + \Lambda_i^T \beta_i \qquad i = 1, 2, ..., K \tag{3}$$

where the K × 1 vector $\beta_i$ is the vector of multivariate regression coefficients of $R_i$ on $f$ with a constant.

As Cochrane (2000) explains in his book, "$b_j$ (vector $b$) coefficients asks whether factor $j$ helps to price assets given the other factors. $b_j$ gives the multiple regression coefficient of $m$ on $f_j$ given the other factors. On the contrary, $\lambda_j$ (components of vector $\Lambda$) asks whether factor $j$ is priced, or whether its factor-mimicking portfolio carries a positive risk premium". $\lambda_j$ gives the single regression coefficient of $m$ on $f_j$. Therefore, when factors are correlated between them, one should test $b_j = 0$ to see whether to include a factor $j$ given the other factors rather than test $\lambda_j = 0$. Multiple authors consider that, when we are able to find a factor whose risk premium is non-zero ($\lambda_j \neq 0$), we can call it a "priced factor" or characteristic and, when we are able to find a factor that has a non-zero coefficient $b_j$, we can call it a "pricing factor".

When confronted with a factor model, a fundamental question arises: which are the genuine pricing factors? We will go over the theoretical factor models that can be used to solve this question.

- Static CAPM. The CAPM of Sharpe (1964) and Lintner (1965) is an equilibrium model in which the excess return on the market portfolio is the only pricing component. As a result, the model predicts that every asset's projected excess return is proportionate to its market beta.

$$E[R_{i,t}] = R_t^d + \beta_i^m (E[R_t^m] - R_t^d), \qquad (4)$$

$R_t^d$ being the risk-free rate, $R_t^m$ the return on the market portfolio, and $\beta_{im}$ the beta of the asset $i$ with respect to the market.

- Intertemporal CAPM. By allowing for various time horizons and preferences among investors, Intertemporal Capital Asset Pricing Model (ICAPM) of Merton (1973) relaxes some assumptions of the static CAPM. Asset risk premia are linear functions of the market beta and other betas in terms of factors. As a result, the market factor is really not the exclusive determinant of pricing any longer.

$$\mu_{i,t} - R_t^d = \beta_{i,t}^m (\mu_t^m - R_t^d) + \sum_{k=1}^{K-1} \beta_{i,t}^k (\mu_{X,t}^k - R_t^d) \qquad (5)$$

where $\mu_t^m$ is the expected return on the market portfolio, and $\mu_{X,t}^k$ is the expected return on the portfolio that maximizes the squared correlation with the $k$-th state variable.

- Consumption-Based CAPM. The ICAPM in its seminal form calls for determining the variables that influence the opportunity set's evolution. In Breeden (1979), the author introduced a CCAPM that substitutes the multiple betas in the decomposition of expected returns by a single beta, which reflects changes in aggregate consumption. The assumptions are the same as in Merton's ICAPM.

  Breeden (1979) shows that expected excess returns are given by:

$$\mu_{i,t} - R_t^d = \frac{\mu_t^C - R_t^d}{\beta_{C,t}^C} \beta_{i,t}^C \qquad i = 1, 2, ..., N \qquad (6)$$

where $\beta_{i,t}^C$ is the beta of asset $i$ with respect to aggregate consumption $C$, $\beta_{C,t}^C$ is the beta of the portfolio that maximizes the squared correlation with changes in aggregate consumption, and $\mu_t^C$ is the expected return on this portfolio.

- Arbitrage Pricing Theory. APT from Ross (1976), along with CAPM, is one of the most influential theories on asset pricing. The APT varies from the CAPM in that its assumptions are less restrictive. The APT concentrates on return factor decomposition: statistical description of asset returns as linear combinations of $K$ common factors and a random disturbance serves as its foundation.

  If $X_i$ indicates the pay-off of an asset, then we have:

$$X_i = E[X_i] + \beta_i^T f + \epsilon_i, \qquad i = 1, 2, ..., N \tag{7}$$

where the idiosyncratic return is uncorrelated from the factors. Such a decomposition of factors is always satisfied, since it is always possible to make a regression for a payoff on a given set of factors. Statistically, it is necessary to assume that there is no autocorrelation in the residuals $\epsilon_i$ across assets.

Empirical Factor Models

Early work on the Sharpe–Lintner CAPM tended to be broadly supportive. In the 1970s, the classic studies of Black et al. (1972) and Fama and Macbeth (1973), found that, empirically, the high-beta stocks tended to have higher average returns than low-beta stocks and that the relation was roughly linear. During the 1980s and 1990s, researchers began to look at other characteristics of stocks besides their betas. All those deviations from the original CAPM were called "anomalies", due to the fact that there did not exist any kind of theoretical model able to explain the existence of those kinds of factors. Empirically, all these anomalies may be explained more efficiently by utilizing multifactor models in which the factors are chosen based on empirical evidence, rather than theoretical support.

- Size and value factors. The empirical evidence that small-cap equities perform better than large-cap equities is known as the size effect. Van Dijk (2011) made an extensive survey of 30 years of research in equity returns. As the author recognizes, this additional factor in CAPM was primarily introduced by Fama and French (1992) with their three-factor model, and since then, "there has been a vigorous debate on whether the size premium is a compensation for systematic risk".
  The size factor is represented as the excess return of small caps over large caps. Fama and French (1992) introduced a "Small–Minus–Big" (SMB) portfolio, which is a zero-investment portfolio built as the difference between the average return on three small-cap portfolios and that on three large-cap portfolios, which has been ordered, previously, according to the book-to-market ratio (Value, Neutral and Growth). This previous filter is defining the third factor, "High-Minus-Low" (HML) of their model, which is another zero-investment portfolio built as the difference between the average return on two value stock portfolios and that on two growth stock portfolios, according to the size quantiles (Big and Small).
  The three-factor model by Fama and French (1992) was expressed as follows:

$$E[R_{i,t}] = R_t^d + \beta_i^m(E[R_t^m] - R_t^d) + \beta_i^{SMB}SMB_t + \beta_i^{HML}HML_t + \epsilon_{i,t} \tag{8}$$

  $R_t^d$ being the risk-free rate, $R_t^m$ the return on the market portfolio, $\beta_i^m$ the beta of the asset $i$ with respect to the market factor, $\beta_i^{SMB}$ the beta of the asset $i$ with respect to the size factor, and $\beta_i^{HML}$ the beta of the asset $i$ with respect to the value factor.
- Momentum factor. Beyond the size and value factors, the momentum factor is the most prevalent factor in the literature. Momentum can be defined as the rate of acceleration of a security's price, and simply, it refers to the inertia of a price trend to continue either rising or falling for a particular length of time. The trading strategies related to this effect, as we will see in the next Section 3.1.2, also called "trend following", seek to capitalize on momentum to enter a trend as it is picking up steam. In statistical terms, the momentum effect characterizes by the existence of serial autocorrelation. The paper by Carhart (1997) can be considered as a study about this matter which, in the last few years, has received more acknowledgement from the investment industry. The author provides evidence in this research which focuses on the mutual fund business, that strong previous performance does not necessarily imply future returns, but that the contrary might be true (if the performance is based on loading up on spe-

cific risk factors). This paper presents a four-factor model with the three factors from Fama and French (1992), plus a new factor which represents the momentum effect:

$$E[R_{i,t}] = R_t^d + \beta_i^m (E[R_t^m] - R_t^d) + \beta_i^{SMB} SMB_t + \beta_i^{HML} HML_t + \beta_i^{WML} WML_t + \epsilon_{i,t} \quad (9)$$

$\beta_i^{WML}$ being the beta of the asset $i$ with respect to the momentum factor. The WML factor is defined as the excess return of an equally-weighted portfolio for 30% of past winners over an identical portfolio of the 30% past losers ("Winners-Minus-Losers").

- Profitability and Investment factors. Following the release of the five-factor model from Fama and French (2015), these two components have lately gained a lot of traction in stock investment techniques. This model was expressed as follows:

$$E[R_{i,t}] = R_t^d + \beta_i^m (E[R_t^m] - R_t^d) + \beta_i^{SMB} SMB_t +$$
$$+ \beta_i^{HML} HML_t + \beta_i^{RMW} RMW_t + \beta_i^{CMA} CMA_t + \epsilon_{i,t} \quad (10)$$

$\beta_i^{RMW}$ being the beta of asset $i$ with respect to the profitability factor, and $\beta_i^{CMA}$ the beta of asset $i$ with respect to the investment factor.

The procedure to estimate these two new factors is similar to previous factors in Fama and French (1992). Stocks are first sorted according to a measure of profitability or investment. The profitability factor is the excess return of robust profitability stocks over weak profitability ones ("Robust-Minus-Weak" or RMW factor), while in the case of the investment factor, it is defined as the excess return of high-investment stocks over low-investment ones ("Conservative-Minus-Aggressive" or CMA factor). The authors choose as measures the operating profit after interest expenses and the growth of total assets, respectively.

### 3.1.2. Price Forecasting

The price forecasting discipline can be defined as the financial area which aims to predict the future price of an asset using market data and its transformations. Some examples of this kind of strategy are trend trading strategies, where one takes a position in the asset only after predicting a change in trend, and statistical arbitrage, which seeks mispricing by detecting asset relationships and/or potential anomalies, believing the anomaly will return to normal. Algorithmic Trading, a financial discipline on the rise in recent years, can be included within trend-following methods, so that, from a theoretical perspective, they can be connected to price forecasting methods.

In this discipline, asset prices play a central and lonely role, while asset pricing models and economic/financial fundamentals stay in a very secondary term. Sometimes, the academic approach to price forecasting has not taken into account risk-related issues or, directly, has not been interested in finding explanatory variables of risk, but in the best prediction of prices. Shortly, the studies have been more interested in the predicting power of modeling than in the explanatory power and economic meaning of the variables included in the model. In the case of technical and trend strategies, the main theoretical topic has to do with the accomplishment of Efficient Market Hypothesis (EMH) that we will expose in the next few lines. Although statistical arbitrage strategies are also very related to EMH and autocorrelation issues, they are connected, to some extent, to theoretical issues previously explained in the Value Investing section, related to asset pricing and anomaly-seeking.

Regardless of the type of approach, time-series data analysis is prevalent in this financial discipline, unlike the previous one, value investing, in which cross-sectional data were clearly predominant.

The EMH states that asset prices always reflect all available information. A direct implication is that it is impossible to "beat the market" consistently on a risk-adjusted basis, since market prices should only react to new information, which can be considered as a crucial implication for the asset management industry, mainly in the case of price forecasting. The EMH is linked to the random walk theory, which defends that the best prediction for tomorrow's price is the current price.

The validity of the random walk hypothesis and, as a result, the unpredictability of asset returns were early concerns for jobs that faced the existence of the momentum effect in stock prices (see previous section). An example of this kind of approach, which uses time-series data, is Lo and MacKinlay (1988). From the perspective of cross-sectional data, we can find Jegadeesh (1990), where authors examine the performance of stock selection techniques based on prior monthly returns. Despite the fact that this anomaly was discovered after the size and value effects, already exposed in previous section, institutional investors have a long and reliable history with the momentum approach Asness et al. (2014).

The fundamental model of asset pricing, the SDF model, may explain how efficient markets are (or are not) connected to random walk theory. This model, as mentioned in the preceding section, makes mathematical predictions about a stock's price assuming that there is no arbitrage, that is, assuming that there is no risk-free way to successfully trade. Formally, if arbitrage is impossible, the model predicts that a stock's price will be the discounted value of its future price. We may rewrite the Equation (1) if we imagine we live in a world without dividends or, to be more limiting, if we assume we are functioning in the short term and no dividend is paid:

$$P_t = E_t(m_{t+1} P_{t+1}) \tag{11}$$

Note that this equation does not generally imply a random walk. However, if we assume the SDF is constant [3], we have:

$$P_t = m E_t(P_{t+1}) \tag{12}$$

Taking logs and assuming that Jensen's inequality term is negligible, we have:

$$\log P_t = \log m + E_t(\log P_{t+1}) \tag{13}$$

which implies that the log of stock prices follows a random walk (with a drift).

Regardless the type of asset pricing model we use, the EMH will be satisfied only under some restrictive assumptions. As Cochrane (2000) points out, "'If investors are risk-neutral, returns are unpredictable, and prices follow martingales (random walk process). In general, prices scaled by marginal utility are martingales, and returns can be predictable if investors are risk averse and if the conditional second moments of returns and discount factors vary over time. This is more plausible at long horizons".

As Timmermann and Granger (2004) points out, "The EMH is a backbreaker for forecasters", because in its crudest form it effectively says that returns from speculative assets are unforecastable. That may appear to be the conclusion of the narrative from an intellectual standpoint. Despite the strength of the argument, it does not appear to be fully compelling to many forecasters. The reason is that, despite its simplicity, the EMH is surprisingly difficult to test and considerable care has to be exercised in empirical tests.

### 3.1.3. Portfolio Management

Portfolio management is the financial discipline of which its basic goal is to establish the optimal weight of each asset by simultaneously maximizing expected return and minimizing risk. Better portfolio optimization models have a more efficient frontier, which can help investors get a greater expected return while taking the same risk. As a result, developing a more efficient portfolio optimization model has become a hot issue in the field of investment management. The main framework which can be considered as a foundation of portfolio management is the seminal work from Markowitz (1952), also known as the Modern Portfolio Theory.

The modern portfolio theory (MPT), often known as mean-variance optimization (Mean-variance optimization (MVO)), is a mathematical framework for constructing an asset portfolio that maximizes expected return for a given degree of risk. It is a formalization and extension of the concept of diversification in investment, which holds that having a

variety of financial assets is less hazardous than owning only one. Its basic concept is that an asset's risk and return should not be assessed by itself, but rather on how it adds to the overall risk and return of a portfolio. The volatility of asset prices, measured in terms of standard deviation of returns, is used as a risk proxy.

For a portfolio consisting of $m$ assets with expected returns $\mu_i$, let $w_i$ be the weight of the portfolio's value invested in asset $i$ such that $\sum_{i=1}^{m} w_i = 1$, and let $w = (w_1, ..., w_m)^T$, $\mu = (\mu_1, ..., \mu_m)$, $\mathbf{1} = (1, ..., 1)^T$. The portfolio return has mean $w^T \mu$ and variance $w^T \sum w$, where $\sum$ is the covariance matrix of the asset returns; see Lai and Xing (2008). Given a target value $\mu_*$ for the mean return of a portfolio, Markowitz characterizes an efficient portfolio by its weight vector $w_{eff}$ that solves the optimization problem:

$$w_{eff} = \arg \min_w w^T \sum w$$
$$\text{subject to} : w^T \mu = \mu_*, \ w^T \mathbf{1} = 1, \ w \geqslant 0 \tag{14}$$

If portfolio return variance, rather than standard deviation, were plotted horizontally, the inverse of the slope of the frontier would be $q$ at the point on the frontier where the inverse of the slope of the frontier would be $q$. On $q$, the entire frontier is parametric. Markowitz (1952) developed a specific procedure for solving the above problem, called the critical line algorithm, that can handle additional linear constraints, upper and lower bounds on assets, and which has been proven to work with a semi-positive definite covariance matrix.

### 3.2. Machine Learning Background

As we have previously described, since the financial crisis of 2008, many quantitative factor models have failed and many conventional factors have become unprofitable Feng et al. (2019). As a result, some market players are seeking for alternatives to standard stock-picking methods.

Many practitioners started building hand-crafted models that can dynamically learn from past data, as popular quantitative elements have become less credible. For many years, investors have relied on econometric approaches, but only a handful have had success using dynamic models based only on these techniques. This might be due to various factors: (1) inherent noise in financial data, (2) the fact that factors can be multicolinear, and (3) that connections between variables and returns can be changeable, non-linear, and contextual. These properties make estimating any dynamic connections between possible predictors and expected returns problematic for traditional models.

Given those limitations of hand-crafted methods, the use of ML techniques that automatically learn the best features from data has become widespread to avoid them. Gu et al. (2021) made the effort to concentrate all the elements involved in the ML techniques to give a context-specific definition of ML: "A diverse collection of high-dimensional models for statistical prediction, combined with so-called regularization methods for model selection and the mitigation of overfitting, and efficient algorithms for searching among a vast number of potential model specifications".

These kind of methods have proven to be extremely successful in other disciplines (like image processing or NLP). Part of that success is due to their ability to generalize to unseen data, learn from noisy distributions, and automatically learn features for complex data. In other words, they excel in those areas where classical domain-specific approaches crafted by experts have found limitations. This meteoric rising of ML applications is mainly due to the fact that, over the last decade, a series of enhancements have enabled recent advances in practical ML and unlocked its utility: increased computing power, increased data availability, and novel optimization techniques and architectural breakthroughs.

The most promising ML applications in finance are on the buy-side, focusing on finding predictive signal among the noise and capturing alphas [4]. Some example applications are: time-series forecasting, market segmentation, regime-switching detection and, of course, asset management, which is the main issue of this work.

Nevertheless, as Arnott et al. (2019) points out, it is important to bear in mind the very special characteristic of financial markets, they reflect the actions of people, which may be influenced by others' actions and by the findings of past research. Unlike other scientific disciplines, research can influence future actions of economic agents. In many respects, the problems that ML faces are just a continuation of the long-standing concerns that quantitative finance experts have always confronted. While investors must exercise caution, perhaps even more caution than in previous implementations of quantitative methods, these new tools have a wide range of financial applications.

ML has been successfully applied to virtually any existing scenario where data are available and useful information can be learned from it. However, different techniques must be applied depending on the problem at hand. There are three main paradigms which are characterized by the nature of the problem:

### 3.2.1. Supervised Learning

The dataset $D$ contains samples $x_i$ together with their expected outputs $y_i$, such as a dog image together with its label "dog", so $D = \{(x_1, y_1), ..., (x_n, y_n)\}$. Therefore, a function $f$ maps each sample to its output $f(x_i) = y_i$. The goal of supervised learning is to find the best function approximation $g \approx f$ that meets two important criteria: (1) minimizes the difference between known samples $g(x_i) \sim f(x_i) = y_i$, and (2) generalizes properly to samples $x_o$ which are not part of $D$.

The following are the the most prominent classical supervised learning methods employed in finance:

- Least Squares Levenberg (1944): A method typically employed to find a linear regression by finding the best fit in the least squares framework, that is, minimizing the sum of squared residuals.
- LASSO Tibshirani (1996): A form of linear regression that is characterized for using shrinkage. This means that LASSO performs L1 regularization to penalize the absolute value of the magnitude of the coefficients. As a result, typically a sparse set of coefficients is produced by helping reduce overfitting and model complexity. Ridge regression works in a similar fashion but by enforcing L2 penalties, which does not produce sparse models.
- Regression Trees Elith et al. (2008): A decision tree is an ML architecture that uses a flowchart-like structure to arrive to infer a result by taking tests over input variables. Each node of the tree is a test on an input variable and depending on the outcome, the flow continues in one branch of the tree or another until the flow reaches the leaves where final outputs are given. Regression trees are just an extension of decision trees where the target value to predict takes the form of a continuous value.
- Random Forest Breiman (2001): A classification or regression method that works by constructing multiple decision trees at training times. That multitude of trees constitutes an ensemble to produce a final prediction (e.g., in the case of classification by voting and in the case of regression by averaging the outputs of all trees).
- Support Vector Machines Cortes and Vapnik (1995): Binary classifiers that map the training samples to points in another space to maximize the gap between the two categories. They can also perform non-linear classification using specific kernels which map those inputs to high-dimensional feature spaces where non-linear decision boundaries can be tackled. Intuitively, an SVM finds a hyperplane that optimally separates the decision boundary by maximizing the distance between one class and another. They can also be used for regression with the appropriate modifications (namely, Support Vector Regressions (SVRs)).

Apart from the classical algorithms, we also briefly explain the modern architectures that are more popular mainly due to the advent of deep learning:

- Neural Networks (NNs) are the most basic architecture, usually composed of individual perceptrons which are arranged into multiple layers—usually with non-linear

activation functions interleaved—of varying width. A perceptron Rosenblatt (1958) is a function $f$ that maps an input $x$ to generate an output $z$ in the following way:

$$z = f(x) = \sigma(wx + b), \tag{15}$$

where $w$ is a vector of weights, $b$ is a bias, and $\sigma$ is an activation function.

In its most simple form, the activation function is just a threshold and the perceptron is just a binary classifier. Note that the bias simply shifts the decision boundary away from the origin. Single-layer perceptrons can be combined together to form a Multi-layer Perceptron (MLP). This architecture is usually composed of three layers: the input layer as before, a hidden layer, and an output one. The input layer remains as before, but the hidden and output ones can be composed by an arbitrary number of nodes (also named neurons). Each of those nodes is a single-layer perceptron that uses a non-linear activation function. Deep Neural Network (NN) (also called fully connected networks) often refer to MLPs with more hidden layers $l$. In this general case, the output of a certain neuron $i$ of a layer $l$ can be defined as follows:

$$z_i^l = \sigma^l(w_i^l * z^{l-1} + b_i^l) \tag{16}$$

Vanilla NNs are capable of learning any non-linear function (they are universal approximators) given enough network complexity, but they face a number of challenges: (1) due to their fully connected nature, they require a huge number of parameters, (2) are usually harder to train, (3) they lose spatial information of the input, and (4) there is no built-in mechanism for capturing sequential data.

- Convolutional Neural Networks (CNNs) LeCun et al. (1998) uses learnable kernel filters to extract relevant features from the inputs by applying the convolution operation with them. They are especially useful with structured data and in those cases where spatial information is important. Typically, they are currently applied to process 2D images (although a convolution can be applied to any dimensionality). For instance, for the 1D case, we can formulate the output of a single neuron in a CNN as follows:

$$z_i^l = \sigma \sum_k w_k^l x_{i-k}^{l-1}, \tag{17}$$

where $w_k^l$ is a vector of weights, also named the kernel, with $k$ elements. This kernel is convoluted over the adequate portion of the input $x_{i-k}$ and passed through a non-linear activation function to compute the output of that neuron. As we can observe, a CNN shares this kernel across the whole layer and, by doing so, it does not fully connect each neuron from the previous layer to the next ones. Furthermore, each layer can have more than one kernel; each one is convoluted individually to produce a separate output. These outputs are often referred to as feature maps and are stacked to form a multi-channeled output.

  With regards to fully connected NNs, they sport some advantages: (1) as we mentioned, by convolving the input with filters of predefined size instead of being fully connected, they capture spatial features, and (2) by not being fully connected, but instead sharing kernel weights across the whole input, they require way less parameters and thus are easier to train and less prone to overfitting.

- Recurrent Neural Networks (RNNs) are specifically designed to deal with sequence data and learn from temporal information. Although internally they can be shaped either as traditional NNs or CNNs, they usually add recurrent connections in their layers, which helps take into account the state from previous sequence elements or temporal instants. They therefore can capture sequential information and share parameters across different timesteps (in a similar fashion as CNNs do spatially). Typically, the most general topology is a fully recurrent RNN where the outputs of all neurons are connected to the inputs of all for them. Each one multiplies the current inputs and previous outputs through an activation function. Other relevant topologies

are Gated Recurrent Unit Network (GRU) and the widely spread Long-short Term Memory (LSTM).

GRUs Cho et al. (2014) features two gating mechanisms: update and reset. The update gate is responsible for determining the amount of previous information that will flow to the next step. The reset gate decides which information from the past timestep to neglect for the current state.

LSTMs Hochreiter and Schmidhuber (1997) features three gating mechanisms: input, output, and forget. This triple gate system allows the architecture to model long- and short-term dependencies properly.

As a matter of fact, all vanilla RNNs, GRUs, or LSTMs are able to model arbitrary time dependencies. The problem is, however, computational and numerical: due to the nature of the training process, the required gradients to learn can easily explode (turn to infinity) or vanish (go to zero) preventing any learning. GRUs are a step forward in comparison with vanilla RNNs and the additional gates from LSTMs help even more to control the information flow to avoid those problems.

### 3.2.2. Unsupervised Learning

The dataset $D$ contains samples $x_i$ without their expected outputs, such as, a set of images of dogs, cats, and mouses but without labels whatsoever, so $D = \{x_0, ..., x_n\}$. In this case, we do not know how $f$ behaves, and therefore we cannot learn the output mapping $f(x_i) = y_i$. Given this setting, the goal of unsupervised learning is to learn a function $g$ that finds patterns or trends in the dataset. For instance, $g$ could be a function that clusters samples in $D$ based on their similarity according to certain features.

Here, we describe the most common unsupervised learning methods, which are often employed as pre-processing steps:

- k-Means Clustering: Usually employed as a pre-processing technique to reduce the number of data points by summarizing them according to their mean expectations. In other words, it takes a number of samples ($n$) and aims to partition them in some sets ($k$, where $k < n$) so that the variance within each cluster is minimized. The most common algorithm is the iterative or naïve k-means Lloyd (1982).
- Principal Components Analysis (PCA) Pearson (1901): Another common pre-processing technique to reduce the number of features while preserving their variance. It does so by computing the principal components of the input data and then using them to perform a change of basis.

It is important to remark Generative Adversarial Networks (GANs) Creswell et al. (2018): An architecture in which two networks (generator and discriminator) compete (adversarial); the goal of the generator is to produce samples able to fool the discriminator whilst the discriminator's role is to detect false examples from the generator. In other words, given a training set, this architecture is able to generate new data that statistically resembles the originally provided. Although initially proposed as a form of generative model for Unsupervised Learning, it has now impacted all paradigms.

### 3.2.3. Reinforcement Learning

The dataset $D$ does not even exist beforehand, but new samples $x_i$ arrive or are generated on the fly, such as, the current state of a chess board in a match after a movement has been performed. A function $g$ which produces an output $z_i$ given a sample and modifies the current state to produce another sample $x_{i+1}$. At any state or at certain times, we can measure how good this $g$ is behaving according to a predefined criteria. The goal is to learn a $g$ that will maximize (or minimize) such criteria.

Reinforcement Learning (RL) can be typically subdivided into two main categories Sutton and Barto (2018): model-based and model-free. The first builds an internal model of the possible states, transitions, and outcomes in the environment; the latter does not use any model but rather learns actions/transitions directly from experience at the expense of statistical efficiency.

As we will review later, the most common reinforcement learning algorithm for finance is Q-learning Watkins and Dayan (1992) and its deep counterpart Deep Q-learning Hester et al. (2018). Q-learning learns a so-called quality or action-value function, which describes how good is it to take a particular action in a determined state. To do so, a table of state-action pairs is kept; this table assigns a scalar reward that defines the quality of the action at a given state. During training, actions are performed either randomly or by looking at the best one in the table. By analyizing the reward after each action, the state-action table can be updated based on the old reward values and the new ones. Deep Q-learning keeps the same procedure, but makes use of deep neural networks to represent the state-action table; it is typically applied in problems in which the option space is so big that defining a state-action table would be too complex and computationally expensive.

Another successful trend of RL is Recurrent Reinforcement Learning (RRL) Li et al. (2015). RRL combines Supervised Learning with Reinforcement Learning typically by employing a RNN to learn the representation of hidden states for the RL algorithm, which is normally a deep Q-learning network to obtain the policy that maximizes the reward.

### 3.3. Performance Criteria

As we will see in Section 5, about Methods, one of the most interesting areas of analysis has to do with the degree of heterogeneity of the performance criteria measures used by researchers in the last five years. The utilization of different measures of return has coexisted with diverse measures of dispersion or volatility and, in some other cases, with ratios of joint measure of risk adjusted return. This is one of the reasons why the reproducibility of this type of research might be questioned. In the next lines we will make a review and description of the main indicators of performance used in the literature of ML applications to Asset Management: returns, risk/returns ratios, goodness of fit/prediction, risk of loss measures, statistical significance, and accuracy of predictions.

#### 3.3.1. Returns

Average returns appear as the first type of performance measure. Depending on the data frequency, we can find daily, monthly and annual returns as measures of profitability of the different investment strategies the authors propose using ML applications.

- As summary measure of return we can define the annualized rate of return as the Compounded Annual Growth Rate (CAGR) of the portfolio value between two periods separated $n$ years:

$$\text{CAGR} = \left[ \frac{P_{t+n}}{P_t} \right]^{(1/n)} - 1 \qquad (18)$$

$P_t$ being the investment value in period $t$ and $n$ the number of years between the two periods we want to compare.

- Sometimes, the returns are measured in terms of Excess Returns. That means that the portfolio return is measured in terms of comparison with the risk-free asset or, in general, a benchmark asset that is used as reference. The arithmetic excess return can be expressed as follows:

$$R_A^E = \overline{R}_p - \overline{R}_b \qquad (19)$$

where $\overline{R}_p$ is the portfolio return and $\overline{R}_b$ the benchmark return. We can also define the excess return as a geometric measure:

$$R_G^E = \frac{\overline{R}_p + 1}{\overline{R}_b + 1} - 1 \qquad (20)$$

#### 3.3.2. Risk/Return Ratios

In many cases, the authors who propose new ML techniques in order to find a better fitting or prediction capacity in financial strategies have opted by considering the risk involved in achieving certain returns. Since the 1960s, investors and researchers have

known how to quantify and measure risk with the variability of returns, basically using ratios, but no single measure actually looked at both risk and return together.

- The most popular ratio to measure portfolio performance is the Sharpe Ratio (SR). Conceived by Bill Sharpe, this measure closely follows his work on the CAPM and, by extension, uses total risk to compare portfolios to the Capital Market Line (CML). It compares the portfolio return with the risk involved in achieving this return, in form of total risk, measured through the return standard deviation, as follows:

$$\text{SR} = \frac{\overline{R}_p - \overline{R}_b}{\sigma_p} \tag{21}$$

  $\sigma_p$ being the standard deviation of portfolio returns.

- Differential Return (DR), by contrast, results in an excess return for the portfolio manager that considers risk in the form of standard deviation (the variability of past returns). It is a sort of a modified Sharpe ratio. Here is the formula:

$$\text{DR} = \overline{R}_p - \left[ \frac{\overline{R}_b - \overline{\text{RFR}}}{\sigma_{R_b}} * \sigma_p \right] - \overline{\text{RFR}} \tag{22}$$

  where RFR is the risk-free rate of return.

- When we used the systematic risk measured by the CAPM, instead of total risk, we are referring to Treynor Ratio (TR):

$$\text{TR} = \frac{\overline{R}_p - \overline{R}_b}{\beta_p} \tag{23}$$

  where $\beta_p$ is the covariance between portfolio returns and market returns according to CAPM.

- We can find another very popular ratio, the Calmar Ratio (CR), which can be defined as the ratio between the CAGR and its Maximum Drawdown (MDD) which, at the same time, measures the maximum observed loss from a peak to a trough of a portfolio, before a new peak is attained, and can be considered as an indicator of a downside risk over a specified time period.

$$\text{CR} = \frac{\text{CAGR}}{\text{MDD}} \tag{24}$$

  A very similar approach is achieved by the Sterling Ratio (STR).

- Finally, the Certainty Equivalent Return (CEQ) considers the risk-free return for an investor with quadratic utility and risk aversion parameter $\lambda$ compared to the risky portfolio and is given by the following equation:

$$\text{CEQ} = (\mu - \text{RFR}) - \frac{\lambda}{2}\sigma^2 \tag{25}$$

### 3.3.3. Goodness of Fit/Prediction

A statistical model's goodness of fit defines how well it fits a collection of data. The disparity between actual values and predicted values under the model in issue is often summarized by goodness of fit measures. Very similarly, goodness of prediction refers to discrepancy between observed values and the values predicted by the model. Every goodness of fit statistic can be defined/used for prediction purposes, and vice versa.

- For linear regression models, R-squared is a goodness-of-fit metric. This statistic shows the percentage of variance in the dependent variable that the independent factors account for when taken jointly. The strength of the link between your model

and the dependent variable is measured by R-squared, which is defined between 0 and 1. It can be calculated as follows:

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \tag{26}$$

where $y_i$ is the actual i-observation of dependent variable $y$, $\hat{y}_i$ its estimated value, and $\bar{y}$ its mean value. When the estimated values are substituted by the predicted ones, we are talking about Out-of-the-Sample R Squared (OOS R2), a measure of goodness of prediction.

- Mean Absolute Percentage Error (MAPE) is one of the most commonly used performance indicators to measure forecast accuracy. It can be defined as the sum of the individual absolute errors divided by the observed value (each period separately). It is the average of the percentage errors.

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^{n} \left| \frac{y_i - \hat{y}_i}{y_i} \right| \tag{27}$$

where $n$ is the number of forecast periods. It is a quite well-known indicator among researchers, despite being a poor accuracy indicator. As it can be seen in the formula, MAPE divides each error individually by the observed value, so it is skewed.

- Mean Absolute Error (MAE) is a very useful performance indicator to measure forecast accuracy. As the name implies, it is the mean of the absolute error.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i| \tag{28}$$

It solves the problem of skewness of the previous indicator but, in return, it is not scaled, so it depends on the magnitude of the dependent variable.

- Root Mean Squared Error (RMSE) is a frequently used measure of the differences between values (sample or population values) predicted by a model or an estimator and the values observed.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2} \tag{29}$$

The RMSE serves to aggregate the magnitudes of the errors in predictions for various data points into a single measure of predictive power. RMSE is a measure of accuracy, to compare forecasting errors of different models for a particular dataset and not between datasets, as it is scale-dependent. Actually, many algorithms (especially for ML) are based on the Mean Squared Error (MSE), which is directly related to RMSE.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \tag{30}$$

### 3.3.4. Risk of Loss Measures

Risk of loss measures are typically used by firms and regulators in the financial industry to gauge the amount of assets needed to cover possible losses. But they are also very common in financial research to compare the risk associated with market and credit positions, mainly in the Risk Management discipline.

- Value at Risk (VaR) is a metric for calculating investment risk. It calculates how much a set of assets would lose (with a specified probability) in a specific time period, such as a day, under typical market conditions. According to Abad et al. (2014), the VaR is thus a conditional quantile of the asset return loss distribution. Let $r_1$, $r_2$, $r_3$, ..., $r_n$ be identically distributed independent random variables representing the financial

returns. Use $F(r)$ to denote the cumulative distribution function, $F(r) = Pr(r_t < r \mid \Omega_{t-1})$ conditionally on the information set $\Omega_{t-1}$ that is available at time $t - 1$. Assume that $r_t$ follows the stochastic process:

$$
\begin{aligned}
r_t &= \mu + \epsilon_t \\
\epsilon_t &= z_t \sigma_t \quad z_t \sim iid(0,1)
\end{aligned}
\tag{31}
$$

where $\sigma_t^2 = E(z_t^2 \mid \Omega_{t-1})$ and $z_t$ has the conditional distribution function $G(z)$, $G(z) = Pr(z_t < z \mid \Omega_{t-1})$. The VaR with a given probability $\alpha \in (0,1)$, denoted by VaR($\alpha$), is defined as the $\alpha$ quantile of the probability distribution of financial returns:

$$
\begin{aligned}
F(\text{VaR}(\alpha)) &= Pr(r_t < \text{VaR}(\alpha)) = \alpha \quad \text{or} \\
\text{VaR}(\alpha) &= inf(v \mid Pr(r_t \le v) = \alpha)
\end{aligned}
\tag{32}
$$

This quantile can be estimated in two different ways: (1) inverting the distribution function of financial returns, $F(r)$ and (2) inverting the distribution function of innovations $G(z)$, in which case is also necessary to estimate $\sigma_t^2$.

$$
\text{VaR}(\alpha) = F^{-1}(\alpha) = \mu + \sigma_t G^{-1}(\alpha)
\tag{33}
$$

- Conditional Value at Risk (CVaR), also known as Expected Shortfall (ES), is a risk measure derived from the previous one. The ES at the $\alpha$% level is the expected return on the portfolio in the worst $\alpha$% of cases. ES is an alternative to VaR that is more sensitive to the shape of the tail of the loss distribution.

The estimation of risk measures has recently gained a lot of attention, partly because of the backtesting issues of VaR and CVaR related to elicitability. As Pitera and Schmidt (2018) mention, "once the parameters of a model need to be estimated, one has to take additional care when estimating risks". The typical estimations approaches, very often, introduce a bias which leads to a systematic underestimation of risk.

3.3.5. Statistical Significance

Many times, researchers are interested in checking whether the variables, factors or characteristics included in the models they propose are statistically significant. In order to achieve this result, they make the usual statistical hypothesis tests based in the t-student distribution. According to this, t-student statistic and $p$-value are the two most common measures to check the statistical significance.

A result has statistical significance in hypothesis testing when it is extremely improbable to have occurred given the null hypothesis. The significance level of the study rejecting the null hypothesis, represented by $\alpha$, is the probability of the study rejecting the null hypothesis if this is true.

- Most times, hypothesis testing of statistical significance can be run using a t-student distribution. The t-statistic can be expressed this way:

$$
t = \frac{\overline{X} - \mu}{\hat{\sigma}/\sqrt{n}}
\tag{34}
$$

where $\overline{X}$ is the sample mean from a sample $x_1, x_2, \ldots, x_n$, of size $n$, $\hat{\sigma}$ is the estimate of the standard deviation of the population, and $\mu$ is the population mean. When the t-statistic value is higher, in absolute value, than the critical value of the t-student distribution given a significance level $\alpha$, the null hypothesis can be rejected, and it can be affirmed than the coefficient or loading is statistically significant, and the variable, factor or characteristic can be considered as relevant.

- The *p*-value in hypothesis significance testing is the probability of getting test findings that are at least as extreme as the actual results, assuming that the null hypothesis is valid. A tiny *p*-value indicates that under the null hypothesis, such an extreme observed result would be very implausible. *p*-values of statistical tests are commonly reported in academic articles in a variety of quantitative domains.

$$p\text{-value} = Pr(T \geq t | H_0) \tag{35}$$

for a one-sided left-tail test, being $H_0$ the null hypothesis. In a formal significance test, the null hypothesis $H_0$ is rejected if the *p*-value is less than a predefined threshold value $\alpha$, which is referred to as the significance level. The meaning is equivalent to that in which the t-student statistic is higher than the critical value at a given significance level.

### 3.3.6. Accuracy of Predictions

In the case of price forecasting and algorithmic trading techniques, performance of the selected classifiers is evaluated using different evaluation metrics. Since the problem is a multi-class classification problem and the distribution of classes is not uniform, therefore it is very common to use accuracy primary classification metrics, namely, precision, recall, and the F-measure.

- Accuracy is a classification metric for evaluating classifiers and can be expressed as:

$$\text{Accuracy} = \frac{\text{\# correct predictions}}{\text{Total \# predictions}} \tag{36}$$

- Precision is the skill of the model to classify samples accurately and can be calculated as follows:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \tag{37}$$

where TP is the true-positive rate and FP is the false-positive rate of the algorithm.
- Recall shows the skill of the model to classify the maximum possible samples, and is represented by the following equation:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \tag{38}$$

where FN is the false-negative rate of the algorithm.
- F-measure describes both precision and recall and can be represented as follows:

$$\text{F-measure} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \tag{39}$$

## 4. Datasets

According to Arnott et al. (2019), one crucial limitation of ML applications in Finance involves data availability. On the one hand, it has been hard to find standardized data sources for finance. On the other hand, deep ML methods usually require large datasets where complex patterns can be extracted while avoiding overfitting LeCun et al. (2015); those kind of datasets are hard to generate in finance. Therefore, data plays a key role in any ML application to asset management.

The work on asset management papers related to ML techniques may be categorized based on the kind of inputs used. A significant number of the articles examined employ structured type inputs, for which processing techniques already exist and the relevance of which has been thoroughly researched. The more current ones permit the usage of unstructured data, which is more difficult to analyse and extract valuable data from. In this section, we review both sources of data.

## 4.1. Structured Data

When we talk about structured data we are referring to data which is organized and fits tidily into spreadsheets and relational databases. This kind of data is quantitative and is often displayed as numbers, dates, values, and strings. It is usually stored in rows and columns.

Most of the articles revised in this work use this type of structure in their information, which is usually open and prepared for API programming interfaces. The most common is the time-series of historical prices, with different frequencies. The preferred periodicity is monthly, but in other cases we can find yearly or daily data. In the case of algorithmic trading, some High Frequency Trading (HFT) models use intraday data.

### 4.1.1. Stock Values

Generally, this information is public and free and can be downloaded from the pages of the stock markets. Besides, some companies like Bloomberg and Reuters provide paid services with additional information related to stock prices. In some articles, daily stock information is used, which consists of opening price, closing price, maximum and minimum, as well as the volume of transactions or negotiation. In the case of intraday information, it is very usual to find data related to bid-ask spread.

### 4.1.2. Macroeconomic Indicators and Financial Information

Taking into account that a large list of papers are focused on asset pricing and factor investing, it is very usual to use corporate information from financial reports to estimate the fair value of the stock. These financial reports are the balance sheet, the income statement and the cash flow statement. Related to them, we can find calculated indicators called financial ratios, which summarize the company's financial and economic situation numerically. Some of these ratios are Debt to Equity, Price to Earnings or Enterprise Value to Operating Earnings (EBITDA). This information is usually provided by private companies like Bloomberg or Thompson Reuters.

Many times, fundamental analysis uses macroeconomic indicators to understand how the stock prices are correlated to changes in variables outside the company. Firstly, because depending on the health of a country's economy, one can estimate the earnings growth of a company, and secondly, because the consumption is on the basis of the majority of asset pricing models. These economic indicators are usually free and published by governments and public institutions. However, they have an important pitfall: their low frequency and the relevant lag with regard to the current prices.

### 4.1.3. Technical Indicators

This category of data is very focused on algorithmic trading, since they are useful for predicting the stock market direction. Technical indicators are not based on economic or financial models. Traders that utilize technical analysis use heuristic or pattern-based signals generated by the price, volume, and/or open interest of a security or contract. Technical analysts utilize indicators to forecast future price changes by evaluating previous data.

There are two types of technical indicators: trend indicators and oscillators. The former ones are focused on identifying movement directions, and the latter ones on finding the turning points in the time-series.

All of then can be calculated using the information of prices obtained in the first section, or directly downloaded from the same information sources. Bloomberg is the main source for this type of data.

## 4.2. Unstructured Data

Unstructured data, unlike structured ones, have no predefined construction or systematization. It comes in text form, audio, images, videos and can be challenging to analyze.

The utilization of unstructured data increases the complexity of stock market forecasting, but at the same time opens a whole new world of possibilities. In order to be

utilized as an input for the model, this data must be preprocessed and transformed to categorical or numerical data. Text mining algorithms, which extract news segments or views from social networks and may create numerical representations, are required for textual unstructured inputs.

The news analysis is usually taken from media sources, but sometimes from the same company. In the case of social networks, we are talking about a very new, challenging and complex world. In this case, the main problem is the enormous volume of information as well as the computational challenges.

News feeds, social media, earnings call transcripts, multiple CRM platforms, email, call notes, and other unstructured data sources are common in Finance. The attractiveness and added value come from a substantially better information base, which includes unstructured data for decision-making that is both relevant and timely.

*4.3. Analysis*

In order to get a more complete perspective about the type of datasets which are used in the research about ML applied to Finance, we will select a small sample of articles within the three areas of research we have used in previous sections. To obtain this sample, we will apply a double filter: firstly, to be published in a Q1–Q2 journal, and, secondly, to be classified in the first quartile in terms of number of citations.

4.3.1. Value/Factor Investing

In this particular field of financial research, it is very common to use public and shared datasets, many times constructed, updated and fed by the own authors. The most usual data frequency is monthly.

- Kozak et al. (2018). The methodology used is a very good example of how this type of datasets can be a good way to generate results that can be globally interpreted. In this work, the authors use, firstly, the 5 × 5 size and book-to-market (B/M) sorted portfolios of Fama and French (1993). Secondly, they use 15 anomaly long-short strategies defined as in Novy-Marx and Velikov (2016) and the underlying 30 portfolios from the long and short sides of these strategies. These two datasets are available in two websites fed by the own authors, with the aim to contribute to future and reproducible research. In the first case, datasets are provided by the Kenneth French's website, which provides downloadable and updated files of Fama/French factors from 1926. In the second case, the author shares the datasets used in his 2016 article. The French's webpage can be considered as one of the best examples of publicly available databases that has become as meeting point of asset pricing researchers.
- Feng et al. (2020). In this paper we can find another excellent example of how using publicly available datasets from other authors. In their article, the authors firstly download all workhorse factors in the U.S. equity market from Ken French's data library. Then they add several published factors directly from the authors' websites, including liquidity from Pastor and Stambaugh (2003) (Stambaugh's website), the q-factors from Hou and Zhang (2015), and the intermediary asset pricing factors from He et al. (2016). In addition to these 15 publicly available factors, they follow Fama and French (1993) to construct 135 long-short value-weighted portfolios as factor proxies, using firm characteristics surveyed in Hou et al. (2017) and Green et al. (2016).
- Gu et al. (2021). Another example of academic database, but in this case accessible by subscription, is Center for Research in Security Prices (CRSP) US Stock Databases, an affiliate of University of Chicago, which contain daily and monthly market and corporate action data for over 32,000 active and inactive securities with primary listings on the NYSE, NYSE American and NASDAQ. The research-quality data created by this transformational project spawned a vast amount of scholarly research from several generations of academics. Today, nearly 500 leading academic institutions in 35 countries rely on CRSP data for academic research and to support classroom

instructions. The CRSP value-weighted index is one of the most usual equity market benchmarks used in financial research.

4.3.2. Portfolio Management

After analyzing the most cited papers within this discipline, we can conclude that, in most of cases, the datasets are composed exclusively by price data. Just in some cases financial data or news can be found in the articles researched. Sometimes, the closing price is joined by other indicators of prices, as maximum, minimum or opening price. Most times, the assets utilized are stock indices constituents and, eventually, Exchange Traded Fund (ETF) prices. The data sources are Datastream, Bloomberg and, in some specific cases, Ken French's website.

- Heaton et al. (2017). Weekly returns data for the components of the biotechnology IBB index in the period 2012–2016. They train the learner without knowledge of the actual component weights. Their goal is to find a selection of investments that outperforms the official index.
- Krauss et al. (2017). Monthly and daily returns data for the components of S&P 500 in the period 1989–2015. The data source is the Thomson Reuters Datastream. The goal was to build portfolios following a statistical arbitrage strategy with better performance than the benchmark index.
- Ban et al. (2018). Ken French's website mentioned in the section of asset pricing. They collect monthly excess returns for three different data sets, composed of 5, 10 and 49 industry portfolios, in the period 1994–2013.
- Almahdi and Yang (2017). A five-asset portfolio using five of the most commonly traded ETFs from different asset categories. They extract the weekly closing prices for each of the five assets from Yahoo Finance website, for the period 2011–2015.
- Lee et al. (2019). Data over a period of 22 years from 1995 to 2016, sourced from Thomson Reuters Datastream database. The dataset is composed by weekly closing prices of 10 global equity indices.
- Paiva et al. (2019). Opening, closing, maximum and minimum daily prices of the components of the Brazilian index Ibovespa, from 2001 to 2016. They were sourced from the Bloomberg terminal.

4.3.3. Price Forecasting

In this third discipline, the datasets used in the most cited articles are quite diverse, depending on the methodology implemented.

- Nikou et al. (2019). In some cases, historical closing prices are the only reference, where the data used include the daily closing price of iShares MSCI UK ETF, also collected from the Yahoo Finance site.
- Zhong and Enke (2019). In other cases we can find financial and economic factors -as in asset pricing models-. In this paper, the dataset includes the daily direction (up or down) of the closing price of the SPDR S&P 500 ETF as the output, along with 60 financial and economic factors as input features. The daily data is collected from 2518 trading days from June 2003. The data sources are public and free (e.g., finance.yahoo.com).
- Khan et al. (2020). Lastly, we can find some examples of financial news and social media data, perfect examples of unstructured data. The source of stock historical daily prices is the same, Yahoo Finance, but the downloaded data have seven features, from date to closing price, passing by traded volume. Given the methodology used in the article, financial news data are also needed, as well as social media data. In the first case, the authors have used Business Insider because it contains a collection of stock market related news from the most famous world news websites, such as Reuters, Financial Times, and so forth. In the second case, they have utilized Twitter API, implemented in Python, to download desired tweets.

## 5. Methods

In this section, we will make an extensive review of all the literature regarding the use of ML techniques in Asset Management. In this sense, we will try to answer the following research questions:

- What financial application areas, within the asset management discipline, are of interest to the financial and ML community?
- In each of these application areas, which ML models/methods are preferred (and more successful)?
- Which are the most used performance metrics by the researchers?

### 5.1. Methodology

The revision of literature has been made classifying the papers according to four financial areas. The three first areas are the financial disciplines we have used in previous sections: asset pricing/value investing, price forecasting and portfolio management. The fourth one, algorithmic trading, might be classified, from a conceptual point of view, as an intermediate discipline between portfolio management—since it can be considered as an special type of trading strategy and price forecasting—as this trading strategy has as priority goal to forecast the future direction of prices. Nevertheless, the growing relevance of this new discipline, not only at a practitioner level, but also between the researchers, has driven to consider it as an independent area. Given its special characteristics, as we will expose in Section 5.5, this application field can be considered as one of the most promising financial areas to be supported by ML applications.

To identify relevant journal articles dealing with ML applications in the four Asset Management disciplines mentioned above, we followed a search process in EBSCOhost, Google Scholar, Science Direct, SpringerLink and Wiley Online Library databases for the period of 2015–2021 using combinations of keywords "machine learning", "deep learning", "neural networks" and "asset management", "portfolio management", "asset pricing", "asset returns", "stocks", "finance", "price forecasting", and "algorithmic trading". After searching through the databases, we reached a list of around 130 identified papers.

After this first preselection, each paper was assessed on quality. This was achieved by using a variety of quality indicators as the citation count and the impact factor of the journal. The arXiv and SSRN databases were also searched to ensure that the most up-to-date research papers were included in the sample. After this second filter we reached a list of 91 identified journals (see Table 2). Finally, and after the assessment of each one of the articles, we focused our research and commentaries in a final number of 60 articles, also summarized in the different tables throughout the article.

**Table 2.** Recurring themes and reference count from the literature review.

| Themes | References |
|---|---|
| Value Investing | 18 |
| Portfolio Management | 31 |
| Price Forecasting | 25 |
| Algorithmic Trading | 17 |
| TOTAL | 91 |

### 5.2. Value/Factor Investing

This financial area is very wide and can be considered as the starting point for the rest of interest areas within the Asset Management discipline. We will consider in this section papers and works regarding the search and location of factors, characteristics, patterns in securities prices which allow the investor to understand the drivers of asset prices, the way the expected returns compensate the assumption of risks, and how to outperform the market. The selection has been summarized in Table 3. The logic behind value/factor investing,

in general, is that a firm's financial performance is influenced by fundamentals/factors, whether latent and unobservable or connected to fundamental characteristics.

Despite being a decades-old academic topic, value/factor investing has gained popularity in line with the emergence of equity traded funds (ETFs) as investment vectors, as we discussed in Section 3.1.1. In the decade of 2010, both gained traction. The mutually advantageous feedback loop between practical financial engineering and academic research has encouraged both sides, which is not surprising.

Nevertheless, in the realm of traditional quantitative techniques, researchers have developed in recent years more sophisticated approaches to organize the so-called "factor zoo" and, more crucially, to detect false anomalies and evaluate alternative asset pricing model specifications due to the ever-increasing number of factors and their relevance in asset management. Hundreds of possible candidates have emerged from the search for factors that explain the cross-section of expected stock returns, as noted by Cochrane (2011) and more recently by Harvey and Liu (2019), David McLean and Pontiff (2016), and Hou et al. (2017).

For instance, Harvey and Liu (2019) used bootstrap on orthogonalized factors in their regressions to solve the problem of correlations among predictors. Fama and French (2018) compared asset pricing models through squared maximum Sharpe ratios, and Giglio and Xiu (2019) estimated factor risk premia using a three-pass method based on principal component analysis. It is obvious that there does not exist an infallible method, but the majority of new contributions in the field are interested in the search for robustness.

In all the previous cases, the decomposition of returns has been made using linear factor models, of course because of its simple interpretation. Nevertheless, beyond the problem of robustness of the estimations, there has been an eternal debate about whether firms returns are explained by their exposure to macroeconomic factors or simply by their intrinsic characteristics. Characteristics, rather than factor loadings, explain a higher share of variation in predicted returns, according to Chordia et al. (2019). On the other hand, adopting a theoretical model in which certain agents' needs are sentiment-driven, Kozak et al. (2018) reconciles factor-based theories of risk premia.

In all this immense sea of different approaches to the factor investing discipline, and given the exponential increase in data availability, ML techniques make their appearance with the aim to help avoid the mentioned limitations of classical approaches. As we mentioned in Section 3, the work from Gu et al. (2020) provided a detailed description of ML tools for empirical asset pricing and give their justification for the growing role of ML in financial research. They perform a comparative analysis of ML methods for "the canonical problem of empirical asset pricing: measuring asset risk premiums". The methods they compare are, between others, generalized linear models, dimension reduction tools, Boosted Regression Trees (BRTs), and Random Forests (RFs). In comparison with standard forecasting approaches, they discover that ML tools increase the description of predicted returns. They also highlight that all ML techniques agree on a limited set of main predictive signals, which include variants on momentum, liquidity, and volatility, and that BRT and NN are the top performing techniques. According to the authors, these findings suggest that enhanced risk premium measurement by ML can simplify the examination of asset pricing economic mechanisms, and that ML is a viable technique for new financial technology.

**Table 3.** Selection of papers for value/factor investing. Mean Absolute Percentage Error (MAPE), Out-of-the-Sample (OOS), Mean Standard Error (MSE), and Maximum Drawdown (MDD).

| Author | Target Market | Method | Performance Criteria |
|---|---|---|---|
| Tobek and Hronec (2020) | NYSE, Amex and NASDAQ common stocks | WLS, PWLS, RF, GBRT, NN | Average return, Sharpe ratio, MDD |
| Giglio and Xiu (2019) | US stocks, T-bonds, C-Bonds and currencies | PCA | R Squared, *p*-value |
| Kelly et al. (2018) | World stocks | IPCA | R Squared, *p*-value |
| Moritz and Zimmermann (2016) | US stocks | DT | Excess returns, R Squared, MSE |
| Kozak et al. (2019) | US stocks | Bayesian and Lasso Regressions | OOS R2, Sharpe ratio |
| Messmer (2017) | US stocks | DFNN | Sharpe ratio |
| Feng et al. (2018a) | NYSE, Amex and NASDAQ common stocks | DFNN | Sharpe ratio |
| Chen et al. (2020) | US stocks | DFNN, LSTM, GAN | Sharpe ratio |
| Feng et al. (2018b) | NYSE, Amex and NASDAQ common stocks | TensorFlow, SGD, AD | MSE, R Squared |
| Simonian et al. (2019) | US stocks | RF, ARL | R Squared, Annual return, Sharpe ratio |
| Sun (2020) | NYSE common stocks | Ordered-Weighted LASSO | SR, Mean returns |
| Freyberger et al. (2020) | NYSE, Amex and NASDAQ common stocks | LASSO | Sharpe ratio |
| Lu et al. (2019) | Chinese stocks | NN, MLP | Average return, Sharpe ratio |
| Feng and He (2019) | US stocks | Bayesian Hierarchical | OOS R2 |
| Feng et al. (2020) | US stocks | DS LASSO | SR, Mean returns, t-stat |
| Sugitomo and Minami (2018) | TOPIX 500 stocks | SVM, GBRT and NN | Average return, Sharpe ratio, RMSE |
| Avramov et al. (2021) | US stocks | NN3, FFN, LSTM, GAN | Average return, Sharpe ratio |
| Aw et al. (2019) | US stocks | NNs | Average return, Sharpe ratio |
| Gogas et al. (2018) | NYSE, Amex and NASDAQ common stocks | SVR | R Squared, MAPE |

The review of specific papers about application of ML techniques to Asset Pricing could start with the contribution mentioned above of Giglio and Xiu (2019). The authors use PCA to solve the problem of bias in the estimation of linear asset pricing models when some priced factors are omitted. They show that in a linear factor model, the risk premium of a factor may be identified independently of the rotation of the other control factors as long as they all cover the actual factor space.

In a similar way, Kelly et al. (2018) presented the Instrumented PCA, a new cross-sectional modeling technique for equity returns that accounts for latent factors and time-varying loadings by incorporating observable features that instrument for the unobservable dynamic loadings. In this sense, if IPCA tool identifies the corresponding latent factors, that will mean that the relationship between characteristics and expected returns is due to risk compensation. The other way round, if no such factors exist, the characteristic effect will be compensation without risk or "anomaly".

The closest work in methodology and application of ML techniques to Gu et al. (2020) comes from Tobek and Hronec (2020). The authors look at out-of-sample returns on over a hundred equities anomalies that have been recorded in scholarly literature. They then demonstrate that ML approaches that combine all anomalies into a single mispricing signal are valuable all around the world and can persist in a liquid universe of stocks.

Among others, the techniques used are Weighted Least Squares (WLS), Penalized Weighted Least Squares (PWLS), RFs, Gradient Boosted Regression Trees (GBRTs) and NNs.

Moritz and Zimmermann (2016) also used an ML approach to look at the cross-section of stock returns. In the context of portfolio sorting, they utilize tree-based models to link information from previous returns to future returns. The authors demonstrate that the traditional linear Fama–MacBeth framework does not take use of all of the data's significant information, and that their ML approach is more robust.

There are two aforementioned contributions which deserve additional analysis. In Kozak et al. (2018), the authors contribute with their study to the everlasting fight between factors and characteristics, on the one hand, and risk and behavioural explanations to mispricing, on the other hand. They point out that traditional factor models' efforts to summarize the cross-section of stock returns using a sparse number of characteristic-based factors was futile. Moreover, there is just not enough redundancy across the large variety of potential predictors for such a basic model to price the cross-section appropriately. As a result, a SDF model requires a large number of characteristic-based factors to be loaded. ML techniques, and more specifically, the unsupervised statistical technique PCA, helps in this process, which might be useful in future study on the economic interpretation of the SDF. In Kozak et al. (2019), the authors' method achieves robust out-of-sample performance by imposing an economically motivated prior on SDF coefficients that shrinks contributions of low-variance principal components of the candidate characteristics-based factors. In other words, if the characteristic-based models doesn't work well with a very low number of factors, a SDF formed from a small number of principal components performs well.

Without a doubt, the most frequent technique of ML in the literature are NNs. A first example of this kind of approach can be found in Messmer (2017). Based on a very large set of firm characteristics, they use DL techniques to predict the US cross-section of stock returns. Specifically, they train a deep NN and, after applying a network optimization strategy, he finds that deep NN learned long-short portfolios can generate attractive risk-adjusted returns in comparison with a linear model. This result highlights the relevance of non-linearities in the relationship between firm characteristics and expected returns. In the same line of study using DL techniques, Feng et al. (2018b) designed a deep NN with the aim to minimize pricing errors. As inputs they use firm characteristics, they generate risk factors as intermediate features, and finally fit the cross-sectional returns as outputs. Another example of deep NN can be found in Chen et al. (2020), where the authors combine three different deep neural network structures in a novel way: a NN to capture non-linearities, a recurrent LSTM network to find a small set of economic state processes, and a GAN to identify the portfolio strategies with the most unexplained pricing information estimate. The primary contributions of this study include the use of the fundamental non-arbitrage condition as a criteria function, the use of an adversarial technique to design the most informative test assets, and the extraction of economic states from a large number of macroeconomic time-series.

The procedure of sorting securities, based on firm characteristics, very usual in factor investing literature, is the starting point of the work by Feng et al. (2018a), which uses multi-layer deep networks to augment traditional long-short factor models.

Another notable architectures are RFs, which are one of the most used classical techniques of ML in recent years. For instance, Simonian et al. (2019) showed how to use RFs to produce factor frameworks that improve upon more traditional models in terms of their ability to account for non-linearities and interaction effects among variables, as well as their higher explanatory power. In combination with Association Rule Learning (ARL), they are able to produce viable trading strategies.

Sun (2020) proposed a new ML method, the Ordered and Weighted LASSO (OWL), which circumvents complications from correlations between the different factors in the traditional approach. This method can identify and group correlated factors while shrinking off redundant ones. Using Monte Carlo simulations, he shows that OWL outperforms

Least Absolute Shrinkage and Selection Operator (LASSO), specially when factors are highly correlated.

Very similarly, Freyberger et al. (2020) suggested a non-parametric approach for determining which features give incremental information for the cross-section of expected returns. They select features and evaluate how they impact expected returns non-parametrically using the adaptive group LASSO. This technique can manage a high number of factors, has a flexible form, and is not affected by outliers.

Lu et al. (2019) tried to extract factors according to the definition from Barra team from MSCI company. They utilize Smart Beta Index technique to construct factor indexes to reflect performance and style on the market they analyze, and they bring NNs into the work of cross-section factor integration. Doing so, their experimental results show that the index that compiled based on factors integration by NNs, specifically with MLP, exhibits better profitability and stability.

In Feng and He (2019), we can find a Bayesian Hierarchical (BH) approach. This market-timing method uses heterogeneous time-varying coefficients driven by lagging fundamental factors to jointly estimate conditional expected returns and residual covariance matrix, allowing for estimation risk in portfolio analysis. The BH approach also allows to model different assets separately while sharing information across assets. According to the authors' conclusions, the BH approach outperforms alternative methods in terms of prediction for the US market. At the same time, they were able to identify the most important factors in the past decade: size, investment and short-term reversal.

The authors of Feng et al. (2020) offered a selection methodology to systematically assess each new factor's contribution to asset pricing beyond what a high-dimensional collection of current factors explains. To evaluate the contribution of a component to explaining asset prices in a high-dimensional context, they offer combined cross-sectional asset pricing regressions with the double-selection LASSO of Belloni et al. (2014). This model selection phase closely resembles the existing literature's strategy to dealing with the proliferation of asset price factors (e.g., Kozak et al. (2018)): to take a large set of factors, to apply some dimension-reduction method (LASSO, Elastic Net (EN), PCA, etc.), and to interpret the resulting low-dimensional model as the SDF.

Sugitomo and Minami (2018) used a multi-factor model of Fama–French type as a starting point to test if the ML techniques are able to enhance portfolio performance. Specifically, they used a typical method, consisting of SVM, GBRT and NN, and verified the effectiveness and applicability of nonlinear methods in practical operation by comparing it with conventional linear models.

Avramov et al. (2021) investigated if ML techniques can remove acceptable economic constraints in empirical finance, which is a largely uncharted field. They investigated whether signals generated by ML procedures can withstand economic constraints both in the cross-section and the time-series. For instance, in the cross-section, they remove microcaps and distressed forms, and in the time-series, they look at the sensitivity of investment payoffs to market conditions with less arbitrage opportunities. They concentrate on two DL approaches that perform well with financial data in order to do this job. They first implement NNs with three hidden layers, and then, they incorporate non-arbitrage conditions into multiple connected NNs, including vanilla NNs, LSTMs, and GANs.

A ML factor model using NNs is developed by Aw et al. (2019). This model delivers a superior in-sample performance, but a mediocre out-of-sample performance versus a conventional factor model. The reason they point out for this underperformance is that the market noise during the training period overwhelmed the non-linear association uncovered in the ML process. Nevertheless, they defend that the rationality behind investor behaviour explains the ultimate success of new ML techniques in asset management.

In Gogas et al. (2018), the methodological approach used is SVR, a direct extension of SVM and the objective, to evaluate the effectiveness of four of the most popular models in asset pricing theory, the CAPM, the APT and the three- and five-factor models from Fama and French. They observe large improvements in comparison to the traditional

linear regression in terms of the main measures of goodness of fit: R-squared-adjusted and MAPE.

*5.3. Portfolio Management*

Portfolio management is the practice of selecting assets for a portfolio during a specified length of time. Even though the basic purpose is the same, there are somewhat several variations of this problem, as seen in other financial applications. Portfolio Management is a broad term that encompasses the following closely related topics: portfolio optimization, selection, construction and allocation.

In portfolio management, the concept of diversification is extremely relevant, as we can find in  Markowitz (1952). Additionally, very related to diversification, we find the concept of assets correlation. Goetzmann and Kumar (2008) argues that while investors are aware of the benefits of diversity, they build portfolios without properly taking the correlations into account. This is the fundamental reason why, while recent and sophisticated portfolio optimization approaches perform well in-sample, they typically perform poorly out-of-sample. For instance, DeMiguel et al. (2009) proves that equal-weighted allocation, which assigns equal weight to each asset, outperforms the whole range of frequently used portfolio optimization strategies. In the end, every optimization model requires the inversion of a positive-definite covariance matrix, which results in errors of such size that the benefits of diversification are completely neutralized. At this point, ML has an important role to play with regard to the simplification of the problem. The selection of papers on ML for portfolio management has been summarized in Table 4.

Ban et al. (2018) adapted two ML methods with regularization for portfolio optimization. The objective of this technique, known with the acronym Performance-based Regularization (PBR), is to restrict the sample variances of the estimated portfolio risk and return, guiding the solution towards one with less estimation error in the performance. The results show how this technique outperforms all other benchmarks in a proportion of two out of three using Fama–French datasets.

Rasekhschaffe and Jones (2019) explains some of the fundamental ideas behind ML and offer a simple example of how investors may use ML approaches to estimate cross-section stock returns while avoiding overfitting. Moreover, in order to demonstrate the benefits of ML techniques to make accurate forecasts, they emphasize the importance of mixing forecasts from several algorithms and training periods for diversification. GBRT, SVM, Adaptive Boosting (AB), Deep Neural Network (DNN) are some examples of the algorithms used. They demonstrate that, with sensible feature engineering and forecast combinations, ML algorithms can produce results that dramatically exceed those derived from simple linear techniques, such as Ordinary Least Squares (OLS).

Very similarly,  Huck (2019) presents a summary of the main techniques that can be implemented to manage a long-short portfolio. He uses three different types of ML tools: DBNs, RF y EN regression, because they are all able to perform classification tasks as demanded by the trading system he designs. After developing several independent statistical arbitrage strategies based on these three ML methods, the article describes how adding predictors is not a guarantee to increase the performance of the portfolios. Among the tools considered, the RF seems to generate the best performance portfolios.

**Table 4.** Selection of papers for portfolio management. Value at Risk (VaR), Conditional Value at Risk (CVaR), Maximum Drawdown (MDD), Certainty Equivalent Return (CEQ), Mean Absolute Percentage Error (MAPE), Mean Absolute Error (MAE), Out-of-the-Sample (OOS), Mean Standard Error (MSE), Maximum Drawdown (MDD).

| Author | Target Market | Method | Performance Criteria |
|---|---|---|---|
| Ban et al. (2018) | NYSE, Amex and NASDAQ common stocks | PBR | Sharpe ratio, Turnover |
| Rasekhschaffe and Jones (2019) | Stocks from 22 countries | GBRT, SVM, AB, DNN | Excess return, Alpha |
| Huck (2019) | US Stocks | DFN, RF, EN | VaR, Sharpe ratio, Max. Drawdown |
| Huotari et al. (2020) | S&P 500 stocks | ANN, EIIE | Sharpe ratio, *p*-value |
| Krauss et al. (2017) | S&P 500 stocks | DNN, GBRT, RF | Return distribution, VaR, Calmar ratio |
| Park et al. (2020) | US and Korean ETFs | LSTM, DNN, Q-Learning | Cum. return, Sharpe ratio, Turnover |
| Heaton et al. (2017) | IBB Index | Autoencoders | Validation error |
| López de Prado (2016) | Monte Carlo simulations | HRP | OOS variance |
| Yun et al. (2020) | World ETFs | PCA, LSTM | IR, MDD, VaR, CVaR |
| Raffinot (2017) | S&P 500 stocks | HRP | IR, Sharpe ratio, MDD |
| Jain and Jain (2019) | NIFTY 50 index | HRP | CVaR, Sharpe Ratio |
| Tristan and Chin Sin (2021) | Singapore Index | AHC-DTW clustering | Cum. Return, Sharpe ratio |
| Konstantinov et al. (2020) | World Assets | NN, LASSO regressions | Sharpe ratio, MDD, CEQ |
| Xue et al. (2018) | Shanghai ETFs | FFN, IMK-ELN | MAP, MDD, Sharpe ratio |
| Wang et al. (2020) | UK Stock Exchange 100 Index | LSTM+MVO | MSE, RMSE, MAPE, MAE, R2 |
| Ta et al. (2020) | S&P 500 stocks | LSTM+MVO | Sharpe ratio |
| Lee et al. (2019) | World equity indices | SVM | Directional accuracy |
| Song et al. (2017) | Selected US Stocks | ListNet and RankNet (NN) | Sharpe ratio |
| Vo et al. (2019) | S&P 500 stocks | LSTM+MVO | MAE, RMSE |
| Ma et al. (2020) | China Securities 100 Index | DMLP, LSTM, CNN | MAE, MSE, MDD |
| Ma et al. (2021) | China Securities 100 Index | DMLP, LSTM, CNN, SVR, RF | MAE, MSE, MDD |
| Almahdi and Yang (2017) | US and World ETFs | RRL | Sharpe ratio, Calmar ratio, Sterling ratio |
| Aboussalah and Lee (2020) | Selected US Stocks | SDDRRL | Total return |
| Paiva et al. (2019) | Ibovespa stocks | SVM | Average return, st.deviation |

In Huotari et al. (2020), the main goal was to look at how modern ML analytics can help with portfolio management, specifically by using an ANN-based system to automatically detect market anomalies using technical analysis and exploiting them to maximize portfolio returns by realizing excess returns. They used the Ensemble of Identical Independent Evaluators (EIIE) architecture described by Jiang et al. (2017) on a sample of 415 stocks from the S&P 500 Index and incorporated selected performance indicators for stock performance in the analysis. They used reinforcement learning to create an ANN-based deep-learning

(multi-layer ANN) agent model for portfolio management (trading model) for this study. A reward function drove the agent model, and the objective was to maximize predicted rewards over time.

Krauss et al. (2017) implemented and analyzed the effectiveness of several ML methods in the context of statistical arbitrage. Specifically, they used DNNs, GBRTs, RFs and, finally, a combination of them all. Each model was trained on lagged returns of all stocks in the S&P500, after elimination of survivor bias. The database is comprised of daily data. The empirical findings show that a simple ensemble of the three techniques produces a significant excess of out-of-sample returns.

Park et al. (2020) proposed a novel long-only portfolio trading strategy in which an intelligent agent is trained to identify an optimal trading action by using Deep Q-learning, on of the most popular Deep Reinforcement Learning (DRL) methods. Compared with the stochastic programming-based models (Monte Carlo simulations) and heuristic methods (technical analysis), the authors defend that the proposed model, using daily data for two different portfolio cases which comprises ETFs from the US stock market, is a superior trading strategy relative to benchmark strategies.

Heaton et al. (2017) presented a four-step algorithm for model construction and validation with special emphasis on building deep portfolios. In particular, they introduced DL hierarchical decision models and provided a smart indexing example by auto-encoding the IBB biotechnology index.

Yun et al. (2020) proposed a two-stage DL framework for portfolio management, which uses LSTM for the prediction model, in addition to a cost function that addresses both absolute and relative return. The proposed methods are evaluated with an ETF dataset, and the empirical results show that the DL two-stage methods outperform ordinary DL models.

In López de Prado (2016), we can find a very complete definition and explanation of hierarchical methods, that address the main pitfalls of the Critical Line Algorithm (CLA), the classical quadratic optimization procedure specifically designed by Markowitz in 1954 for inequality-constrained portfolio optimization problems, which was, at the time, a brilliant solution to the generic-purpose quadratic programming models that did not guarantee a correct solution after a known number of iterations. In particular, the author presented the Hierarchical Risk Parity (HRP) method, which is based on graph theory and ML techniques, and used the information in the covariance matrix without requiring its inversion or positive definitiveness. The reason is that this new approach replaces the covariance structure with a tree structure. By using Monte Carlo simulations, the author demonstrates that, despite the CLA method delivering the minimum-variance portfolio, the HRP produces lower out-of-sample variance portfolios.

Raffinot (2017) proposes a hierarchical clustering-based asset allocation method, which uses graph theory and ML techniques, in a very similar way to López de Prado (2016). Complete Linkage (CL), Average Linkage (AL) and Directed Bubble Hierarchical Tree (DBHT) are among the hierarchical clustering approaches described and tested, using three empirical datasets from US Stock Market. AL and DBHT prove to be the best clustering methods, and the clustered portfolios to achieve statistically better risk-adjusted performance than commonly used portfolio optimization techniques.

In Jain and Jain (2019) we can find research which is also focused on the out-of-sample performance of the portfolios, because it aims to test if there are any covariance matrix forecasting techniques that outperform both traditional risk-based and ML-based portfolios (such as HRP introduced by López de Prado (2016)) empirically. According to their results, HRP is less sensitive to bad specifications of covariance than minimum variance or maximum diversification portfolios, while it is less robust than an inverse volatility weighted portfolio. The authors used daily prices from the individual stocks comprising the NIFTY 50 from the Indian stock market.

The approach is very similar in Tristan and Chin Sin (2021), because it aimed to use unsupervised time-series clustering-based ML techniques to diversify portfolios and over-

come the varied outcomes of industry diversification. Specifically, they used shape-based clustering approach for diversification, the Agglomerative Hierarchical Clustering algorithm (AHC-DTW), applied to the daily prices of the top 82 stocks listed in the Singapore equity market, and was demonstrated to clearly outperform industry diversification.

Another example of the use of hierarchical-based techniques is the work by Konstantinov et al. (2020), that might be placed in an intermediate point between the financial areas of factor investing and asset allocation. The aim of their work is to approach and compare factor and asset allocation portfolios using both traditional and alternative allocation techniques, considering centrality and hierarchical-based networks, specifically LASSO. The monthly data used comes from the US stock market.

Xue et al. (2018) developed Incremental Multiple Kernel Extreme Learning Machine (IMK-ELM), which aims to enhance the efficiency of previous algorithms to make classification tasks in robo-advisors services. Specifically, the novel algorithm is able to handle heterogeneous customer information sets. The empirical results, reached through simulation, show that IMK-ELM outperforms other generic classification methods.

Wang et al. (2020) suggested a hybrid approach consisting of LSTM networks and a MVO model for optimum portfolio construction in combination with asset preselection, in which long-term dependencies of financial time-series data may be represented. In this sense, the empirical results show how LSTM clearly outperforms other ML techniques, such as SVM, RF and common DNNs. In the second stage, after selecting asset with higher returns according to that ML technique, the MVO model is applied for portfolio optimization. The monthly data used comes from the UK Stock Exchange 100 Index.

Ta et al. (2020) implemented ML techniques at a double level. First, they use LSTM, an special type of RNN, to forecast stock direction based on historical data. In the second level, and in order to build an efficient portfolio, they make use of multiple optimization techniques, including Equal-weighted modeling (EQ), Monte Carlo simulation (MCS) and MVO. The results show that the LSTM prediction model works efficiently by obtaining high accuracy from stock prediction, and generating portfolios which outperform those obtained using alternative techniques as Logistic Regression (LR) or SVM. The data used in this work are the 10-year daily historical stock prices of 500 large-cap stocks listed on the America Stock Exchange S&P500.

The approach in Lee et al. (2019) is very similar to those aforementioned in the sense of using a double-scale framework. In this case, the first level of the trading strategy is based on the effect and usefulness of networks indicators. Using a Vector Autorregression model (VAR) model, they forecasted global and regional stock markets' directions. Once these trend predictions were defined, they were used as inputs for determining portfolio strategies via several ML techniques, such as LR, SVM, and RFs. The research data are daily stock index prices from 10 different countries over a period of 22 years. The empirical results show that the prediction accuracy and profit performances are enhanced with network indicators, and that the SVM approach displays the best performances.

Song et al. (2017) focused their attention on the area of investors' sentiment. In particular, they showed that learning-to-rank algorithms are effective in producing reliable rankings of the best and worst performing stocks and, according to them, they are able to implement outperforming portfolio strategies which produce risk-adjusted returns superior to the benchmark. The algorithms used with weekly prices and financial news from US Stock market are called RankNet and ListNet, which are supervised learning approaches that relies on NNs and Gradient Descent Optimization (GDO) techniques.

The work by Vo et al. (2019) enters into the emerging field of Socially Responsible Investments (SRIs). The authors defend that traditional optimization methods for portfolio management are inadequate for this kind of investments, so they propose a new model called Deep Responsible Investment Portfolio (DRIP) that contains a Multivariate Bidirectional LSTM neural network to predict stock returns for the construction of a SRI portfolio using the MVO model. For the empirical application, they used daily closing prices of all individual stocks contained in the S&P500 from the past 30 years. The portfolios obtained

using this method had a high degree of accuracy and achieved much higher Environmental, Social and Governance (ESG) ratings compared with standard MVO models.

Ma et al. (2020) used the most common DL techniques to build prediction-based portfolio optimization models. These models start by using DNNs to forecast each stock's future performance. The risk of each stock is then calculated using DNNs predictive errors. Following that, portfolio optimization models are constructed by combining predicted returns and semi-absolute deviation of prediction errors. These models are contrasted against three equal-weighted portfolios, with stocks picked by DMLP, LSTM, and CNN, respectively. Additionally, two SVR-based portfolio models are utilized as benchmarks. As experimental data, this article uses component stocks of the China Securities 100 index in the Chinese stock market. The prediction-based portfolio model based on DMLP performs the best among these models.

The approach is very similar in Ma et al. (2021), where the authors propose two ML models, specifically RF and SVR, and three DL models, in particular DMLP, LSTM, and CNN, for stock preselection before portfolio construction. Therefore, they incorporates those results in advancing MVO models. As benchmarks, they utilize portfolio models with Autorregressive Integrated Moving Average (ARIMA) predictions. Once evaluated the models with daily data from the Chinese Stock market index, the results show that portfolio models with RF predictions are the best among the set of models used.

Almahdi and Yang (2017) applied RRL techniques to build an optimal variable weight portfolio allocation. For this purpose, they propose a RRL with a coherent risk adjusted performance objective function, based on the expected maximum drawdown, the Calmar ratio. They use as dataset five asset portfolios built with five of the most commonly traded ETFs. The maximized function using this method yields superior return performance than other techniques proposed in the existing literature.

Similarly, using RRL techniques, Aboussalah and Lee (2020) aimed to enter into the field of continuous action and multi-dimensional state spaces, and, hence, they propose the called Stacked Deep Dynamic Recurrent Reinforcement Learning (SDDRRL) to build a real-time optimal portfolio. As performance metric, they use the Sharpe ratio. The model was trained and tested with daily data from the S&P500 index, and the results showed that their model outperforms three different benchmarks, including MVO model.

In Paiva et al. (2019), again, the two-pass methodology is applied to get a portfolio selection model. First, they apply a ML technique to make stock price predictions, the SVR method. After that, they use the traditional scope of MVO for portfolio construction. They compare the results of this method with benchmarks applied to the daily prices of the individual assets included in the Ibovespa index, and the results are favourable for the proposed technique.

### 5.4. Price Forecasting

Return and price forecasting play a main role in modern portfolio theory, in asset pricing models and, from a practical point of view, in the asset management industry. As it was pointed out by Gu et al. (2020), because the primary objective of asset pricing is to understand the behavior of risk premiums, return prediction is economically significant. In academic finance, the terms expected return and risk premium have been usually interchanged.

As Henrique et al. (2019) pointed out, the challenge of predicting asset prices and the search for models and profitable systems is still attractive for both the academia and the financial professionals despite the strong presence of the efficient market hypothesis (EMH) by Malkiel and Fama (1970), which defends that most of the financial asset prices follow, statistically, a random walk process, and therefore are almost unpredictable.

The recent literature facing the use of ML techniques has been summarized in Table 5 and can be divided, roughly, into two levels. In the first level, the new ML technology has been applied to enhance forecasts made using traditional inputs, such as fundamental accounting data, macroeconomic data, or technical indicators. In other words, improving

the results of econometric and time-series analysis using fundamental and technical approaches. In the second level, ML has been used to extract new inputs form alternative data, such as sentiments from news data.

**Table 5.** Selection of papers for price forecasting. Mean Absolute Error (MAE), Mean Standard Error (MSE), Maximum Drawdown (MDD).

| Author | Target Market | Method | Performance Criteria |
|---|---|---|---|
| Kumar et al. (2018) | Selected Indian stocks | SVM, RF, KNN, NB, Softmax | Accuracy, F-measure |
| Lee and Kang (2020) | S&P 500 stocks | MLP, CNN | Total return, MDD |
| Cervelló-Royo and Guijarro (2020) | Nasdaq 100 stocks | GBM, RF, CNN | Average accuracy ratio |
| Nabipour et al. (2020) | Selected Indian stocks | DT, RF, KNN, LR, ANN, RNN, LSTM | Accuracy, F-measure |
| Zhong and Enke (2019) | S&P 500 ETFs | DNN, FFNN | MSE |
| Shen and Shafiq (2020) | Selected Chinese stocks | CNN, LSTM | Overall accuracy |
| Nikou et al. (2019) | iShares MSCI UK ETFs | ANN, SVM, RF, RNN, LSTM | MAE, MSE, RMSE |
| Minh et al. (2018) | S&P 500 index | RNN, TGRU, LSTM | Overall accuracy |
| Ding et al. (2015) | S&P 500 index | WB+CNN | Total return |
| Khan et al. (2020) | Selected US stocks | GNB, SVM, LR, MLP, KNN, GBM, RF | Accuracy, precision, recall, F-measure |

Within the first category, we can cite Kumar et al. (2018), where the authors analysed various Supervised Learning (SL) techniques for stock market prediction. Specifically, they consider SVM, RF, K-Nearest Neighbor (KNN), Naive Bayesian Classifier (NV), and Softmax. Five models were developed and their performances compared in predicting stock market trends. According to their results, the RF algorithm performed the best for large datasets, while NV showed the best performance for small datasets. Moreover, they found that the reduction of technical indicators reduces the accuracy of each algorithm.

Lee and Kang (2020) proposed a novel method for training NNs to forecast the future prices of stock indexes. The main contribution of their work is to use only the data of individual companies -instead of index data- to obtain sufficient amount of data for training NN for the prediction of stock indexes. Their experiments, focused on S&P 500, show that NN trained this way outperform NN trained on stock index data. Specifically, they obtain a 5–16% annual return before transaction costs during the period 2006–18.

To evaluate the predictive capacity of certain popular technical indicators in the technological NASDAQ index, in Cervelló-Royo and Guijarro (2020) the authors compared the performance of four ML algorithms: RF, Deep Feedforward Neural Network (DFNN), GBRT and Generalized Linear Model (GLM). The results show that the RF algorithm beats the other ML algorithms, forecasting the market trend 10 days ahead with an average accuracy level of 80%.

Nabipour et al. (2020) seeked the reduction of risk in trend prediction using ML and DL techniques, and applying 11 ML logarithms to data from the Tehran stock exchange. They used DT, RF, AB, eXtreme Gradient Boosting (EGB), SVM, NV, KNN, LR and ANN as ML algorithms and RNN and LSTM as DL ones. The analysis findings show that RNN and LSTM beat other prediction models for continuous data by a significant margin.

Zhong and Enke (2019) focused their analysis on daily stock market returns, specifically in the SPDR S&P ETF prices. To anticipate the daily direction of future stock market index returns, DNNs and ANNs were applied to the full preprocessed but untransformed dataset, as well as two datasets transformed through principal component analysis (PCA). The simulation findings demonstrate that DNNs with two PCA-represented datasets,

as well as numerous other hybrid machine learning methods, have considerably greater classification accuracy than those with the full untransformed dataset.

Shen and Shafiq (2020) propose a solution for the Chinese stock market prices prediction which consists of a feature engineering along with a fine-tuned system based on a LSTM model. The feature engineering applied are the Feature Expansion (FE) approaches with Recursive Feature Elimination (RFE), followed by PCA. This proposed solution outperforms the ML and ML-based models in similar previous works.

In Nikou et al. (2019), the authors want to examine how well ML models can forecast the daily close price data of the iShares MSCI United Kingdom ETF. Four models of ML algorithms are used in the prediction process. The results indicate that the RNN and LSTM DL methods are better in prediction than the other ML methods, and the SVM method is in the next rank with respect to ANN and RF methods, according to error prediction.

Within the second category, we can cite Minh et al. (2018). This study is focused on the financial news as potential factor which causes fluctuations in stock prices. The main contribution of the paper is to propose a novel framework to forecast directions of stock prices by using both financial news and sentiment dictionary, specifically a novel two-stream GRU and Stock2Vec, a sentiment word embedding trained on financial news dataset.

Ding et al. (2015) suggested a DL method for event-driven stock market prediction. Events are first retrieved from news content and represented as dense vectors using an Neural Tensor Network (NTN). Second, a Deep Convolutional Neural Network (DCNN) is utilized to predict both short- and long-term effects of events on stock price fluctuations. When compared to state-of-the-art baseline approaches, experimental findings demonstrate that their model can enhance S&P 500 index prediction and individual stock prediction by nearly 6%.

Khan et al. (2020) utilized algorithms to examine the influence of social media and financial news data on stock market forecast accuracy over a ten-day period. To increase prediction performance and quality, feature selection and spam tweets reduction are carried out on the data sets. Finally, DL is implemented to achieve maximum prediction accuracy, and certain classifiers are ensembled. The highest forecast accuracies of 80.53% percent and 75.16%, respectively, were reached utilizing social media and financial news, according to their findings. RF classifier is found to be consistent and highest accuracy of 83.22% is achieved by its ensemble.

Similarly, Khan et al. (2020) sought to know whether public opinion and political environment in Pakistan on any given day may influence stock market patterns in individual firms or the whole market. Ten ML algorithms were applied to the final data sets to predict the stock market future trend. The experimental findings suggest that the sentiment feature increases machine learning algorithm prediction accuracy by 0–3%, whereas the political situation feature improves algorithm prediction accuracy by around 20%. The Sequential Minimal Optimization (SMO) algorithm was identified to have the best performance.

*5.5. Algorithmic Trading*

It is not easy to classify this discipline. Although it might be classified within one of the areas aforementioned, mainly price forecasting, we have thought that it has some characteristics that may justify to be defined into an independent category.

Algorithmic Trading can be defined as "buy-sell decisions made solely by algorithmic models", as cited in Ozbayoglu et al. (2020). These decisions can be based on some simple rules, mathematical models, optimized processes or, as in the case of ML and DL, highly complex function approximation techniques. Market making, inter-market spreading, arbitrage, and pure speculation such as trend following are examples of methods employed in algorithmic trading. Many fall into the category of high-frequency trading (HFT), which is characterized by high turnover and high order-to-trade ratios.

In the case of trend-following methods, many times algorithmic trading applications are connected to price forecasting methods. Consequently, most price forecasting models that generate buy-sell signals based on their forecast are likewise classified as Algo-trading

systems. Obviously, we will refrain from analyzing those papers focused on price fore-casting again. However, most of times Algo-trading is coupled with technical analysis, which means that, from a conceptual point of view, this discipline is poorly connected with Finance theory and Financial Economics. Nonetheless, and given the great symbiosis between algorithmic systems and ML techniques, we will make a very brief review of the most interesting papers within this emerging and very popular field among financial market practitioners. This review has been condensed in Table 6.

**Table 6.** Selection of papers for algorithmic trading.

| Author | Target Market | Method | Performance Criteria |
| --- | --- | --- | --- |
| Sezer et al. (2017) | Dow 30 stocks | MLP-ANN | Overall accuracy |
| Troiano et al. (2018) | Dow 30 stocks | LSTM | Overall accuracy |
| Sirignano and Cont (2019) | Selected US stocks | LSTM | Overall accuracy |
| Tsantekidis et al. (2017) | Selected Finnish stocks | CNNs | Recall, precision, F1 |
| Sezer and Ozbayoglu (2018) | World ETFss | CNNs | Annualized returns |
| Niño et al. (2018) | Selected US stocks | CNNs | Directional accuracy |
| Tsantekidis et al. (2017) | Selected Finnish stocks | LSTM | Recall, precision, F1 |

Sezer et al. (2017) proposed a stock trading system based on optimized technical analysis parameters for creating buy-sell points using genetic algorithms. The optimized parameters were then used to a DMLP for buy-sell-hold predictions. Daily prices of Dow 30 stocks were used. The results show that this method enhances the stock trading performance.

Troiano et al. (2018) employed a LSTM based on market indicators, in particular, the Moving Average Convergence and Divergence (MACD) signals, to forecast the trend of the Dow 30 stocks' daily prices. Using also LSTM, Sirignano and Cont (2019) proposed a novel method that used limit order book flow and history information for the determination of the stock movements. The same approach can be found in Tsantekidis et al. (2017).

Due to their effectiveness in image classification problems, several research papers have focused on using CNNs-based models. To do so, however, the financial input data have to be converted into images, which demands some creative preprocessing. It is the case of Sezer and Ozbayoglu (2018), who presented a new method for converting financial time-series data containing technical analysis indicator outputs to 2-dimensional images and classifying these images using CNNs to derive trading signals. Using the Limit Order Book Data and transaction data, Niño et al. (2018) encoded financial time-series into an image-like representation, and get a very good performance in terms of directional accuracy. Tsantekidis et al. (2017) proposed a novel method that uses the last 100 entries form the limit order book to create a 2-dimensional image for the stock price prediction.

## 6. Discussion

After reviewing the selected datasets, methods, and performance criteria, we will take a step back to analyze them at a higher level of abstraction. This discussion will comprise three aspects: firstly, we will provide an overview of the state of the art of ML for asset management; second, we will highlight the most successful data, methods, and criteria for each financial discipline we reviewed; at last, we will lay down the existing challenges which might motivate further research.

### 6.1. Overview

The traditional approach to solving asset management issues has been the focus of academics and practitioners alike throughout the last 50 years. However, it has also exhibited many drawbacks. Next, we will describe shortly the main pitfalls and challenges that this discipline is currently addressing:

- Researchers are sometimes compelled to present incomplete results that are often refuted by additional studies due to the publication bias towards successful results

(see Harvey 2017). As a result, replication is critical, and many academic findings have a very short expiration date, especially if transaction costs are taken into account Cakici and Zaremba (2021).

- One of the main pitfalls of the traditional econometric approach has to do with the p-hacking. As it was demonstrated by Chen (2019), p-hacking alone cannot account for all the anomalies documented in the literature. One way to reduce the risk of spurious detection is to increase the hurdles (often, the t-statistics) but the debate whose title might be "the factor zoo" is still ongoing Harvey and Liu (2019).

- Because of its easy understanding, the decomposition of returns into linear factor models is extremely useful. Nonetheless, there is an eternal dispute in the academic literature as to whether business returns are explained by exposure to macroeconomic variables or merely by firm characteristics. Until the new century the factor-based explanation for risk premium was the favourite, but after the seminal work by Daniel and Titman (1997), the characteristics-based explanation has become a great competitor of the traditional outlook.

- Some researchers have observed fading anomalies as a result of publication: once an anomaly is made public, agents invest in it, driving up prices and causing the anomaly to vanish. David McLean and Pontiff (2016) documents this impact in the United States, while Jacobs and Mülle (2020) finds that post-publication factor returns are sustained in other relevant markets. Herding may be destroying factor premia Krkoska and Schenk-Hoppé (2019), and the democratization of so-called smart-beta products (particularly the ETFs) that enable investors to actively invest in specific styles (value, low volatility, etc.) may speed the process up.

- Researchers have developed more sophisticated techniques to organize the so-called factor zoo and, more significantly, to detect false anomalies. Feng et al. (2020), for example, uses LASSO selection and Fama–MacBeth regressions to see if new factor models are worthwhile. They calculate the benefit of adding one new factor to a set of preset factors, demonstrating that many of the factors described in papers published in the 2010 decade do not provide much extra value.

- There is no such thing as a flawless approach, but the sheer volume of contributions in the field emphasizes the importance of robustness. The notion that factors are likely to change over time is a key obstacle for short-term strategies. We refer for instance to Cooper and Maio (2019).

- As we have seen in the Section 5.4 about price forecasting, the difficulty to test consistently, using traditional approaches, the EMH, leaves a huge space to alternative techniques.

- In the case of MVO, as Cochrane (2011) points out, even though it is not a particularly useful guide to computation, classic one-period mean-variance analysis is a brilliantly useful characterization of an optimal portfolio, useful for final investors to understand and think hard about risk allocations. Even when investors are considering highly non-normal payoffs, traditional mean-variance analysis continues to dominate portfolio applications. Nevertheless, many researchers have tried to improve the suitability of this model from different perspectives.

In conclusion, traditional financial economics has no perfect answers to all these pitfalls and challenges described. On the other hand, ML techniques have found an excellent breeding ground to develop all its potentialities. As Cerniglia and Fabozzi (2020) points out, financial theory, market behavior, ever-increasing data sources, and computational innovation are all required for good forecasting and pricing. By putting together the most comprehensive toolbox, you can create realistic computational models. This goal may be achieved using both financial econometrics and ML techniques. According to these authors, "ML tools provide the ability to make more accurate predictions by accommodating nonlinearities in data, understanding complex interaction among variables, and allowing the use of large, unstructured datasets. The tools of financial econometrics remain critical in answering questions related to inference among the variables describing economic

relationships in finance; when properly applied, their role has not diminished with the introduction of ML".

As we have reviewed, and according to Song et al. (2017), ML algorithms are commonly employed for financial market forecast and trading strategies. There are three different sorts of applications. The first sort of application forecasts asset prices or returns in the future. Generally, SVR and NN algorithms are used in this type of strategies. The drawback with this strategy is that it has a high error rate owing to the difficulty in predicting future asset values based on erratic financial market data. The second type uses classification algorithms to anticipate price movement directions, such as SVM and DTs. These approaches generally have significant forecast accuracy, but this does not always imply high profitability. For example, a model can anticipate small gains properly but massive losses wrongly, resulting in a substantial downside risk. Rule-based optimization is the third type. Its goal is to find the best trading indicator and parameter combinations (for example, technical indicators, fundamental indicators, and macroeconomic indicators). Optimization algorithms that have been explored include Gaussian Process (GP) and RL.

To sum up, the trend indicates that ML algorithms clearly outperform traditional econometrics approaches. However, the landscape is extremely diverse in terms of data and applied methods, which suggests a lack of common benchmarks, methodologies, and frameworks. Both classical and modern ML methods have been applied successfully. The high number of applications of LSTMs models is especially remarkable where time-series come into play.

### 6.2. Discipline Focus

After providing an overview for the field as a whole, we will increase our focus to discuss each one of the reviewed disciplines (including the extra algorithmic trading) from the ML perspective.

- Value/Factor investing: The landscape is not specifically dominated by any particular technique. PCA is successful in most works as a pre-processing technique whilst other classical ML methods like RFs, SVMs, or shallow NNs are present in almost all the analyzed works. RNNs does not have much presence in this discipline. The paradigms are mainly Supervised Learning and Unsupervised Learning.
- Portfolio Management: In this discipline, we observed a trend of favoring RNNs architectures to model long-term dependencies of financial time-series data. In particular, most of the reviewed methods make use of variations of LSTMs (usually combining them with other techniques like MVO). RL methods appear in this discipline coupled with RNNs in the form of LSTMs in the most recent works. The dominant paradigms are Supervised Learning (SL) and Reinforcement Learning (RL).
- Price Forecasting: The most heterogeneous discipline where all sorts of ML methods have been applied to either refine the output of other algorithms, to generate predictions on its own or even as a technique to process alternative data sources. A few works make use of social media, financial news, and sentiment analysis to increase prediction accuracy. There is no dominant paradigm in this discipline.
- Algorithmic Trading: In this case, most reviewed papers make use of SL to train architectures more typically suited to target other domains. For instance, CNNs (which are common in image processing scenarios) are applied to specially pre-processed financial data with success. Oftentimes, they are also coupled with RNNs techniques to model time dependencies, usually applying LSTMs.

### 6.3. Challenges and Future Research

We have identified the following main challenges and opportunities:

- Standard Datasets: The whole field is characterized by a lack of curated datasets to be reused by the community. Although some datasets are built upon the portfolio database of Fama and French (1993), most of them deviate from this standard. Furthermore, even those which reuse that database end up diverging in terms of the final data

available for investigation. Therefore, creating a standard database (complete and broad enough) to be reused by the research community is a need for further works.

- Reproducibility: No common methodology or framework for method training and benchmarking has been established. This hurts reproducibility since most of the analyzed methods are difficult, if not impossible, to compare against each other (unless reimplemented specifically for each scenario). In addition, almost no paper includes codes or data to be accessed by other researchers. Establishing a reproducibility framework for asset management ML research is a high-impact workstream for improving the quality of life and pace of the research community.
- Multimodal Data: Most methods are focused on analyzing numerical financial data to generate predictions. Analyzing alternative sources of information like news, social media, sentiment, and user-generated content can provide useful cues for financial decisions. Few works make use of those data sources at the moment. The challenge of combining all those multimodal sources and multiple architectures might unlock new levels of prediction accuracy.
- Heterogeneous Architectures: Arguably due to the state of immaturity in which financial ML sits nowadays (with regard to other more established synergies like image processing or NLP), no clear architectures for processing financial data have been established yet. There is a broad range of papers that spawn new models, and few that build upon solid groundwork to improve them. Finding the common patterns and unifying those diverse architectures could have a beneficial effect to the community for broad adoption in industry (in a similar way as other networks, such as UNet or ResNet, have done for image processing by becoming the de facto standard for many applications).
- Algorithmic trading: This application field is characterized by a very interesting trade-off. From an academic perspective, this area is relatively disconnected from the theoretical background about asset pricing and value investing, which has had a central role in financial economics during the last five decades, and it has been summarized in Section 3. However, in return, and precisely because of this characteristics—exclusive dependence on price data—it is the financial discipline that can maximize the contributions of ML applications. We can find a future challenge in the possibility of combining both issues, deepening trading algorithms with a higher relevance of financial fundamentals.

## 7. Conclusions

To the best of our knowledge, this is the first review paper in the literature which focuses on asset management using ML. In comparison with other papers, which are either broader (tackling finance as a whole) or narrower (portfolio management), this paper is devoted to an intermediate area which is of great interest for both academics and practitioners, but has not yet been reviewed and structured adequately. We formulated the asset management problem and broke it down into disciplines, providing the reader with enough background knowledge to either get familiar with the area or acquire a standard terminology. Furthermore, we also provided a background on ML for the more finance-oriented audience. We covered the contemporary literature of datasets and methods, creating a comprehensive review of 12 sources of data and more than 50 techniques. We also discussed the criteria that have been applied throughout them to train and measure their performance. In the end, we discussed the reviewed methods and datasets and provided useful insight in the shape of future research directions and open challenges in the field.

**Author Contributions:** Conceptualization, P.M.M.-F., A.G.-G., J.S.B.-S. and M.A.P.; methodology, P.M.M.-F. and A.G.-G.; software, P.M.M.-F. and A.G.-G.; validation, P.M.M.-F. and A.G.-G.; formal analysis, P.M.M.-F. and A.G.-G.; investigation, P.M.M.-F. and A.G.-G.; resources, P.M.M.-F. and A.G.-G.; data curation, P.M.M.-F. and A.G.-G.; writing—original draft preparation, P.M.M.-F. and A.G.-G.; writing—review and editing, P.M.M.-F. and A.G.-G.; visualization, P.M.M.-F. and A.G.-G.;

## Notes

[1]   The term "stochastic discount factor" is used because $m$ generalizes standard discount factor ideas. If there is no uncertainty, we may use the conventional present value formula to describe prices

$$p_t = \frac{1}{1 + r_d} x_{t+1}$$

where $r_d$ is the risk free rate, the return of a discount bond with a unique and riskless payoff of 1 dollar in the period $t + 1$.

[2]   An investor's first order conditions give the basic consumption-based model, in which the pricing kernel or SDF can be expressed as:

$$m_{t+1} = \delta \frac{u'(c_{t+1})}{u'(c_t)}$$

where $c_t$ denotes the level of consumption in period $t$, $u$ the utility function and $\delta$ the elasticity of intertemporal substitution of consumption.

[3]   In the consumption-based model already described, it means that investors are risk neutral, i.e., $u(c)$ is linear or there is no variation in consumption, and we are in short time horizons where $\delta$ is close to one.

[4]   The alpha component is, according with the different factor models we have exposed, the independent term which is not associated with any factor of risk and, supposedly, can be associated with the skill of the investors to find extra returns in the securities they invest in.

## References

Abad, Pilar, Sonia Benito, and Carmen López. 2014. A comprehensive review of value at risk methodologies. *The Spanish Review of Financial Economics* 12: 15–32. [CrossRef]

Aboussalah, Amine Mohamed, and Chi Guhn Lee. 2020. Continuous control with stacked deep dynamic recurrent reinforcement learning for portfolio optimization. *Expert Systems with Applications* 140: 112891. [CrossRef]

Almahdi, Saud, and Steve Y. Yang. 2017. An adaptive portfolio trading system: A risk-return portfolio optimization using recurrent reinforcement learning with expected maximum drawdown. *Expert Systems with Applications* 87: 267–79. [CrossRef]

Arnott, Rob, Campbell R. Harvey, and Harry Markowitz. 2019. A backtesting protocol in the era of machine learning. *The Journal of Financial Data Science* 1: 64–74. [CrossRef]

Asness, Clifford, Andrea Frazzini, Ronen Israel, and Tobias Moskowitz. 2014. Fact, fiction, and momentum investing. *The Journal of Portfolio Management* 40: 75–92. [CrossRef]

Avramov, D., Si Cheng, and Lior Metzker. 2021. Machine Learning Versus Economic Restrictions: Evidence from Stock Return Predictability. Available online: https://ssrn.com/abstract=3450322 (accessed on 30 September 2021).

Aw, E. N., Joshua Jiang, and John Q. Jiang. 2019. Rise of the machines: Factor investing with artificial neural networks and the cross–section of expected stock returns. *The Journal of Investing* 29: 6–17. [CrossRef]

Ban, Gah Yi, Noureddine El Karoui, and Andrew E. B. Lim. 2018. Machine learning and portfolio optimization. *Management Science* 64: 1136–54. [CrossRef]

Belloni, A., V. Chernozhukov, and C. Hansen. 2014. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* 81: 608–50. [CrossRef]

Black, Fischer, Michael C. Jensen, and Myron Scholes. 1972. *The Capital Asset Pricing Model: Some Empirical Tests*. Westport: Praeger Publishers Inc., pp. 79–121.

Breeden, Douglas T. 1979. An intertemporal asset pricing model with stochastic consumption and investment opportunities. *Journal of Financial Economics* 7: 265–96. [CrossRef]

Breiman, Leo. 2001. Random forests. *Machine Learning* 45: 5–32. [CrossRef]

Bustos, Oscar, and Alexandra Pomares-Quimbaya. 2020. Stock market movement forecast: A systematic review. *Expert Systems with Applications* 156: 113464. [CrossRef]

Cakici, Nusret, and Adam Zaremba. 2021. Size, value, profitability, and investment effects in international stock returns: Are they really there? *The Journal of Investing* 30: 65–86. [CrossRef]

Carhart, M. M. 1997. On persistence in mutual fund performance. *The Journal of Finance* 52: 57–82. [CrossRef]

Cavalcante, Rodolfo C., Rodrigo C. Brasileiro, Victor L. F. Souza, Jarley P. Nobrega, and Adriano L. I. Oliveira. 2016. Computational intelligence and financial markets: A survey and future directions. *Expert Systems with Applications* 55: 194–211. [CrossRef]

Cerniglia, Joseph A., and Frank J. Fabozzi. 2020. Selecting computational models for asset management: Financial econometrics versus machine learning. Is there a conflict? *The Journal of Portfolio Management* 47: 107–18. [CrossRef]

Cervelló-Royo, Roberto, and Francisco Guijarro. 2020. Forecasting stock market trend: A comparison of machine learning algorithms. *Finance, Markets and Valuation* 6: 37–49. [CrossRef]

Chen, Andrew Y. 2019. The limits of p-hacking: A thought experiment. In *Finance and Economics Discussion Series*. Washington, DC: Divisions of Research & Statistics and Monetary Affairs, Federal Reserve Board.

Chen, Luyang, Markus Pelger, and Jason Zhu. 2020. Deep Learning in Asset Pricing. Available online: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3350138 (accessed on 4 April 2019).

Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv* arXiv:1406.1078.

Chordia, Tarun, Amit Goyal, and Jay Shanken. 2019. Cross-Sectional Asset Pricing with Individual Stocks: Betas Versus Characteristics. Available online: https://ssrn.com/abstract=2549578 (accessed on 1 November 2017).

Cochrane, John H. 2000. *Asset Pricing*. Princeton: Princeton University Press.

Cochrane, John H. 2011. Discount rates. NBER Working Paper No. w16972. Available online: https://ssrn.com/abstract=1820084 (accessed on 1 April 2011).

Cooper, Ilan, and Paulo F. Maio. 2019. New evidence on conditional factor models. *Journal of Financial and Quantitative Analysis* 54: 1975–2016. [CrossRef]

Cortes, Corinna, and Vladimir Vapnik. 1995. Support-vector networks. *Machine Learning* 20: 273–97. [CrossRef]

Cox, John C., and Stephen A. Ross. 1976. The valuation of oprions for alternative stochastic processes. *Journal of Financial Economics* 3: 145–66. [CrossRef]

Creswell, Antonia, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A. Bharath. 2018. Generative adversarial networks: An overview. *IEEE Signal Processing Magazine* 35: 53–65. [CrossRef]

Daniel, Kent, and Sheridan Titman. 1997. Evidence on the characteristics of cross sectional variation in stock returns. *Journal of Finance* 52: 1–33. [CrossRef]

David McLean, R., and Jeffrey Pontiff. 2016. Does academic research destroy stock return predictability? *Journal of Finance* 71: 5–32. [CrossRef]

DeMiguel, V., L. Garlappi, and R. Uppal. 2009. Optimal versus naive diversification: How inefficient is the 1/n portfolio strategy? *The Review of Financial Studies* 22: 1915–53. [CrossRef]

Ding, Xiao, Yue Zhang, Ting Liu, and Junwen Duan. 2015. Deep learning for event-driven stock prediction. Paper presented at Twenty-Fourth International Joint Conference on Artificial Intelligence, Buenos Aires, Argentina, July 25–31.

Durairaj, M., and B. H. Krishna Mohan. 2019. A review of two decades of deep learning hybrids for financial time series prediction. *International Journal on Emerging Technologies* 10: 324–331.

Elith, Jane, John R. Leathwick, and Trevor Hastie. 2008. A working guide to boosted regression trees. *Journal of Animal Ecology* 77: 802–13. [CrossRef]

Emerson, S., R. Kennedy, L. O'Shea, and J. O'Brien. 2019. Trends and applications of machine learning in quantitative finance. Paper presented at 8th International Conference on Economics and Finance Research (ICEFR 2019), Lyon, France, June 18–21.

Fama, Eugene, and Kenneth French. 1992. The cross section of expected stock returns. *The Journal of Finance* XLVII: 427–65. [CrossRef]

Fama, Eugene, and Kenneth French. 1993. Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics* 33: 3–56. [CrossRef]

Fama, Eugene F., and Kenneth R. French. 2015. A five-factor asset pricing model. *Journal of Financial Economics* 116: 1–22. [CrossRef]

Fama, Eugene F., and Kenneth R. French. 2018. Choosing factors. *Journal of Financial Economics* 128: 234–52. [CrossRef]

Fama, Eugene F., and James D. Macbeth. 1973. Risk, return, and equilibrium: Empirical tests. *Journal of Political Economy* 81: 607–36. [CrossRef]

Feng, Guanhao, Stefano Giglio, and D. Xiu. 2020. Taming the factor zoo: A test of new factors. *The Journal of Finance* 75: 1327–70. [CrossRef]

Feng, Guanhao, and Jingyu He. 2019. Factor Investing: Hierarchical Ensemble Learning. Available online: https://ssrn.com/abstract=3326617 (accessed on 31 January 2019).

Feng, Guanhao, Nicholas G. Polson, and J. Xu. 2018a. Deep Factor Alpha. Available online: https://www.arxiv-vanity.com/papers/1805.01104/ (accessed on 1 March 2018).

Feng, Guanhao, Nicholas G. Polson, and Jianeng Xu. 2018b. Deep Learning in Characteristics-Sorted Factor Models. Available online: https://ssrn.com/abstract=3243683 (accessed on 3 May 2018).

Feng, Guanhao, Nicholas G. Polson, and Jianeng Xu. 2019. Deep Learning in Asset Pricing. Available online: https://www.semanticscholar.org/paper/Deep-Learning-in-Asset-Pricing%E2%88%97-Feng-Kong/d0404ccdd0598f5ac6abee0ae97741323190aaf2 (accessed on 15 March 2019).

Freyberger, Joachim, Andreas Neuhierl, and M. Weber. 2020. Dissecting characteristics nonparametrically. *Review of Financial Studies* 33: 2326–77. [CrossRef]

Giglio, Stefano, and Dacheng Xiu. 2019. Asset Pricing with Omitted Factors. Available online: https://ssrn.com/abstract=2865922 (accessed on 14 September 2019).

Goetzmann, W. N., and A. Kumar. 2008. Equity portfolio diversification. *Review of Finance* 12: 433–63. [CrossRef]

Gogas, P., Theofilos Papadimitriou, and Dimitrios Karagkiozis. 2018. The Fama 3 and Fama 5 Factor Models under a Machine Learning Framework. Publisher=Rimini Centre for Economic Analysis. Available online: https://ideas.repec.org/p/rim/rimwps/18-05.html (accessed on 1 May 2018).

Green, Jeremiah, John R. M. Hand, and F. Zhang. 2016. The characteristics that provide independent information about average u.s. monthly stock returns. *Review of Financial Studies* 30: 4389–36. [CrossRef]

Gu, Shihao, Bryan Kelly, and Dacheng Xiu. 2020. Empirical asset pricing via machine learning. *The Review of Financial Studies* 33: 2223–73. [CrossRef]

Gu, Shihao, Bryan T. Kelly, and Dacheng Xiu. 2021. Autoencoder asset pricing models. *Journal of Econometrics* 222: 429–50. [CrossRef]

Harvey, C., and Y. Liu. 2019. Lucky Factors. Available online: https://ideas.repec.org/a/eee/jfinec/v141y2021i2p413-435.html (accessed on 8 April 2021).

Harvey, Campbell R. 2017. Presidential address: The scientific outlook in financial economics. *Journal of Finance* 72: 1399–40. [CrossRef]

He, Zhiguo, Bryan T. Kelly, and Asaf Manela. 2016. Intermediary asset pricing: New evidence from many asset classes. *Journal of Financial Economics* 126: 1–35. [CrossRef]

Heaton, J. B., N. G. Polson, and J. H. Witte. 2016. Deep Portfolio Theory. Available online: https://arxiv.org/abs/1605.07230 (accessed on 23 May 2016).

Heaton, J. B., N. G. Polson, and J. H. Witte. 2017. Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry* 33: 3–12. [CrossRef]

Henrique, Bruno Miranda, Vinicius Amorim Sobreiro, and Herbert Kimura. 2019. Literature review: Machine learning techniques applied to financial market prediction. *Expert Systems with Applications* 124: 226–51. [CrossRef]

Hester, Todd, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan, John Quan, Andrew Sendonaris, Ian Osband, and et al. 2018. Deep q-learning from demonstrations. Paper presented at AAAI Conference on Artificial Intelligence, New Orleans, IL, USA, February 2–7, vol. 32.

Hochreiter, Sepp, and Jürgen Schmidhuber. 1997. Lstm can solve hard long time lag problems. *Advances in Neural Information Processing Systems* 9: 473–479.

Hou, Kewei, Chen Xue, and Lu Zhang. 2017. Replicating Anomalies. Available online: https://ssrn.com/abstract=2961979 (accessed on 12 June 2017).

Hou, K., C. Xue, and L. Zhang. 2015. Digesting anomalies: An investment approach. *Review of Financial Studies* 28: 650–705. [CrossRef]

Huang, Jian, Junyi Chai, and Stella Cho. 2020. Deep learning in finance and banking: A literature review and classification. *Frontiers of Business Research in China* 14: art. n. 13. [CrossRef]

Huck, Nicolas. 2019. Large data sets and machine learning: Applications to statistical arbitrage. *European Journal of Operational Research* 278: 330–42. [CrossRef]

Huotari, Tommi, Jyrki Savolainen, and Mikael Collan. 2020. Deep reinforcement learning agent for s&p 500 stock selection. *Axioms* 9: 130.

Jacobs, Heiko, and Sebastian Mülle. 2020. Anomalies across the globe: Once public, no longer existent? *Journal of Financial Economics* 135: 213–30. [CrossRef]

Jain, Prayut, and Shashi Jain. 2019. Can machine learning-based portfolios outperform traditional risk-based portfolios? The need to account for covariance misspecification. *Risks* 7: 74. [CrossRef]

Jegadeesh, N. 1990. Evidence of predictable behavior of security returns. *The Journal of Finance* 45: 881–98. [CrossRef]

Jiang, Weiwei. 2021. Applications of deep learning in stock market prediction: Recent progress. *Expert Systems with Applications* 184: 115537. [CrossRef]

Jiang, Zhengyao, Dixing Xu, and Jinjun Liang. 2017. A deep reinforcement learning framework for the financial portfolio management problem. *arXiv* arXiv:1706.10059.

Kamley, Sachin, Shailesh Jaloree, and R. S. Thakur. 2016. Performance forecasting of share market using machine learning techniques: A review. *International Journal of Electrical and Computer Engineering* 6: 3196–204.

Kelly, Bryan T., Seth Pruitt, and Yinan Su. 2018. Characteristics Are Covariances: A Unified Model of Risk and Return. Available online: https://ssrn.com/abstract=3032013 (accessed on 15 October 2018).

Khan, Wasiat, Mustansar Ali Ghazanfar, M. Azam, A. Karami, K. H. Alyoubi, and A. S. Alfakeeh. 2020. Stock market prediction using machine learning classifiers and social media, news. *Journal of Ambient Intelligence and Humanized Computing* 1: 1–24. [CrossRef]

Khan, Wasiat, Usman Malik, Mustansar Ali Ghazanfar, Muhammad Awais Azam, Khaled H. Alyoubi, and Ahmed S. Alfakeeh. 2020. Predicting stock market trends using machine learning algorithms via public sentiment and political situation analysis. *Soft Computing* 24: 11019–43. [CrossRef]

Konstantinov, Gueorgui S., Andreas Chorus, and Jonas Rebmann. 2020. A network and machine learning approach to factor, asset, and blended allocation. *The Journal of Portfolio Management* 46: 54–71. [CrossRef]

Kozak, Serhiy, Stefan Nagel, and Shrihari Santosh. 2018. Interpreting factor models. *Journal of Finance* 73: 1183–223. [CrossRef]

Kozak, Serhiy, Stefan Nagel, and Shrihari Santosh. 2019. Shrinking the cross-section. *Journal of Financial Economics* 135: 271–92. [CrossRef]

Krauss, Christopher, Xuan Anh Do, and Nicolas Huck. 2017. Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500. *European Journal of Operational Research* 259: 689–702.

Krkoska, Eduard, and Klaus Reiner Schenk-Hoppé. 2019. Herding in smart-beta investment products. *Journal of Risk and Financial Management* 12: 47. [CrossRef]

Kumar, Indu, Kiran Dogra, Chetna Utreja, and Premlata Yadav. 2018. A comparative study of supervised machine learning algorithms for stock market trend prediction. Paper presented at International Conference on Inventive Communication and Computational Technologies, ICICCT 2018, Lalitpur, Nepal, July 20–22.

Lai, Tze Leung, and Haipeng Xing. 2008. *Statistical Models and Methods for Financial Markets*. New York: Springer.

LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* 521: 436–44. [CrossRef]

LeCun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86: 2278–324. [CrossRef]

Lee, Jinho, and Jaewoo Kang. 2020. Effectively training neural networks for stock index prediction: Predicting the s&p 500 index without using its index data. *PLoS ONE* 15: e0230635.

Lee, Tae Kyun, Joon Hyung Cho, Deuk Sin Kwon, and So Young Sohn. 2019. Global stock market investment strategies based on financial network indicators using machine learning techniques. *Expert Systems with Applications* 117: 228–42. [CrossRef]

Levenberg, Kenneth. 1944. A method for the solution of certain non-linear problems in least squares. *Quarterly of Applied Mathematics* 2: 164–68. [CrossRef]

Li, Bin, and Steven C. H. Hoi. 2014. Online portfolio selection: A survey. *ACM Computing Surveys* 46: 1–33. [CrossRef]

Li, Xiujun, Lihong Li, Jianfeng Gao, Xiaodong He, Jianshu Chen, Li Deng, and Ji He. 2015. Recurrent reinforcement learning: A hybrid approach. *arXiv* arXiv:1509.03044.

Lintner, John. 1965. The valuation of risk assets and the selection of risky investments in stock portfolios and capital budgets. *The Review of Economics and Statistics* 47: 13–37. [CrossRef]

Lloyd, Stuart. 1982. Least squares quantization in pcm. *IEEE Transactions on Information Theory* 28: 129–37. [CrossRef]

Lo, A. W., and A. C. MacKinlay. 1988. Stock market prices do not follow random walks: Evidence from a simple specification test. *Review of Financial Studies* 1: 41. [CrossRef]

Lu, Zhichen, Wen Long, Jiashuai Zhang, and Yingjie Tian. 2019. Factor integration based on neural networks for factor investing. Paper presented at Computational Science—ICCS 2019, 19th International Conference, Faro, Portugal, June 11–14; vol. 11538 LNCS, pp. 286–92.

López de Prado, M. 2016. Building diversified portfolios that outperform out of sample. *The Journal of Portfolio Management* 42: 59–69. [CrossRef]

Ma, Y., R. Han, and W. Wang. 2020. Prediction-based portfolio optimization models using deep neural networks. *IEEE Access* 8: 115393–405. [CrossRef]

Ma, Yilin, Ruizhu Han, and Weizhong Wang. 2021. Portfolio optimization with return prediction using deep learning and machine learning. *Expert Systems with Applications* 165: 113973. [CrossRef]

Malkiel, B. G., and E. F. Fama. 1970. Efficient capital markets: A review of theory and empirical work. *Journal of Finance* 25: 383–417. [CrossRef]

Markowitz, Harry. 1952. Portfolio selection. *The Journal of Finance* 7: 77–91.

Merton, Robert C. 1973. An intertemporal capital asset pricing model. *Econometrica* 41: 867–87. [CrossRef]

Messmer, Marcial. 2017. Deep Learning and the Cross-Section of Expected Returns. Available online: https://ssrn.com/abstract=3081555 (accessed on 2 December 2017).

Minh, Dang Lien, Abolghasem Sadeghi-Niaraki, Huynh Duc Huy, Kyungbok Min, and Hyeonjoon Moon. 2018. Deep learning approach for short-term stock trends prediction based on two-stream gated recurrent unit network. *IEEE Access* 6: 55392–404. [CrossRef]

Moritz, Benjamin, and Tom Zimmermann. 2016. Tree-Based Conditional Portfolio Sorts: The Relation between Past and Future Stock Returns. Working Paper. Available online: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2740751 (accessed on 9 January 2022).

Nabipour, Mojtaba, Pooyan Nayyeri, Hamed Jabani, Shahab S., and Amir Mosavi. 2020. Predicting stock market trends using machine learning and deep learning algorithms via continuous and binary data; a comparative analysis. *IEEE Access* 8: 150199–212. [CrossRef]

Nikou, Mahla, Gholamreza Mansourfar, and Jamshid Bagherzadeh. 2019. Stock price prediction using deep learning algorithm and its comparison with machine learning algorithms. *Intelligent Systems in Accounting, Finance and Management* 26: 164–74. [CrossRef]

Niño, Jaime, Germán Hernández, Andrés Arévalo, and Diego León. 2018. Cnn with limit order book data for stock price prediction. Paper presented at Future Technologies Conference (FTC) 2018, Vancouver, BC, Canada, November 15–16.

Novy-Marx, Robert, and Mihail Velikov. 2016. A taxonomy of anomalies and their trading costs. *Review of Financial Studies* 29: 104–47. [CrossRef]

Nti, Isaac Kofi, Adebayo Felix Adekoya, and B. Weyori. 2019. A systematic review of fundamental and technical analysis of stock market predictions. *Artificial Intelligence Review* 53: 3007–57. [CrossRef]

Ozbayoglu, Ahmet Murat, Mehmet Ugur Gudelek, and Omer Berat Sezer. 2020. Deep learning for financial applications: A survey. *Applied Soft Computing Journal* 93: 106384. [CrossRef]

Paiva, Felipe Dias, Rodrigo Tomás Nogueira Cardoso, Gustavo Peixoto Hanaoka, and Wendel Moreira Duarte. 2019. Decision-making for financial trading: A fusion approach of machine learning and portfolio selection. *Expert Systems with Applications* 115: 635–55. [CrossRef]

Park, Hyungjun, Min Kyu Sim, and Dong Gu Choi. 2020. An intelligent financial portfolio trading strategy using deep q-learning. *Expert Systems with Applications* 158: 113573. [CrossRef]

Pastor, L., and R. F. Stambaugh. 2003. Liquidity risk and expected stock returns. *Journal of Political Economy* 111: 642–85. [CrossRef]

Pearson, Karl. 1901. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* 2: 559–72. [CrossRef]

Pitera, Marcin, and Thorsten Schmidt. 2018. Unbiased estimation of risk. *Journal of Banking and Finance* 91: 133–45. [CrossRef]

Raffinot, Thomas. 2017. Hierarchical clustering-based asset allocation. *Journal of Portfolio Management* 44: 89–99. [CrossRef]

Rasekhschaffe, Keywan Christian, and Robert C. Jones. 2019. Machine learning for stock selection. *Financial Analysts Journal* 75: 70–88. [CrossRef]

Rosenblatt, Frank. 1958. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review* 65: 386. [CrossRef]

Ross, Stephen A. 1976. The arbitrage theory of capital asset pricing. *Journal of Economic Theory* 13: 341–60. [CrossRef]

Ross, Stephen A. 1978. A simple approach to the valuation of risky streams. *The Journal of Business* 51: 453–75. [CrossRef]

Sezer, Omer Berat, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. 2020. Financial time series forecasting with deep learning: A systematic literature review: 2005–2019. *Applied Soft Computing* 90: 106181. [CrossRef]

Sezer, Omer Berat, and Ahmet Murat Ozbayoglu. 2018. Algorithmic financial trading with deep convolutional neural networks: Time series to image conversion approach. *Applied Soft Computing* 70: 525–38. [CrossRef]

Sezer, Omer Berat, Murat Ozbayoglu, and Erdogan Dogdu. 2017. A deep neural-network based stock trading system based on evolutionary optimized technical analysis parameters. *Procedia Computer Science* 114: 473–80. [CrossRef]

Sharpe, William F. 1964. American finance association capital asset prices: A theory of market equilibrium under conditions of risk. *The Journal of Finance* 19: 425–42.

Shen, Jingyi, and M. Omair Shafiq. 2020. Short-term stock market price trend prediction using a comprehensive deep learning system. *Journal of Big Data* 7: 66. [CrossRef]

Simonian, Joseph, Chenwei Wu, Daniel Itano, and Vyshaal Narayanam. 2019. A machine learning approach to risk factors: A case study using the fama–french–carhart model. *The Journal of Financial Data Science* 1: 32–44. [CrossRef]

Sirignano, Justin A., and R. Cont. 2019. Universal features of price formation in financial markets: Perspectives from deep learning. *Quantitative Finance* 19: 1449–59. [CrossRef]

Snow, Derek. 2020. Machine learning in asset management-part 1: Portfolio construction-trading strategies. *The Journal of Financial Data Science* 2: 10–23. [CrossRef]

Song, Qiang, Anqi Liu, and Steve Y. Yang. 2017. Stock portfolio selection using learning-to-rank algorithms with news sentiment. *Neurocomputing* 264: 20–28. [CrossRef]

Sugitomo, Seisuke, and Shotaro Minami. 2018. Fundamental factor models using machine learning. *Journal of Mathematical Finance* 8: 111–18. [CrossRef]

Sun, Chuanping. 2020. Dissecting the Factor Zoo: A Correlation-Robust Machine Learning Approach. Available online: https://ssrn.com/abstract=3263420 (accessed on 4 November 2020).

Sutton, Richard S., and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*. Cambridge: MIT Press.

Ta, Van Dai, Chuan Ming Liu, and Direselign Addis Tadesse. 2020. Portfolio optimization-based stock prediction using long-short term memory network in quantitative trading. *Applied Sciences* 10: 437. [CrossRef]

Tibshirani, Robert. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* 58: 267–88. [CrossRef]

Timmermann, Allan, and Clive W. J. Granger. 2004. Efficient market hypothesis and forecasting. *International Journal of forecasting* 20: 15–27. [CrossRef]

Tkáč, Michal, and Robert Verner. 2016. Artificial neural networks in business: Two decades of research. *Applied Soft Computing Journal* 38: 788–804.

Tobek, Ondrej, and Martin Hronec. 2020. Does it pay to follow anomalies research? machine learning approach with international evidence. *Journal of Financial Markets* 2: 100588. [CrossRef]

Tristan, Lim, and Ong Chin Sin. 2021. Portfolio management: A financial application of unsupervised shape-based clustering-driven machine learning method. *International Journal of Computing and Digital Systems* 10: 235–43.

Troiano, Luigi, Elena Mejuto Villa, and Vincenzo Loia. 2018. Replicating a trading strategy by means of lstm for financial industry applications. *IEEE Transactions on Industrial Informatics* 14: 3226–34. [CrossRef]

Tsantekidis, Avraam, Nikolaos Passalis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. 2017. Forecasting stock prices from thelimit order book using convolutional neural networks. Paper presented at 2017 IEEE 19th Conference on Business Informatics (CBI), Thessaloniki, Greece, July 24–27.

Tsantekidis, A., N. Passalis, A. Tefas, J. Kanniainen, M. Gabbouj, and A. Iosifidis. 2017. Using deep learning to detect price change indications in financial markets. Paper presented at 2017 25th European Signal Processing Conference (EUSIPCO), Kos, Greece, August 28.

Van Dijk, M. A. 2011. Is size dead? A review of the size effect in equity returns. *Journal of Banking and Finance* 35: 3263–74. [CrossRef]

Vo, Nhi N. Y., Xuezhong He, Shaowu Liu, and Guandong Xu. 2019. Deep learning for decision making and the optimization of socially responsible investments and portfolio. *Decision Support Systems* 124: 113097. [CrossRef]

Wang, Wuyu, Weizi Li, Ning Zhang, and Kecheng Liu. 2020. Portfolio formation with preselection using deep learning from long-term financial data. *Expert Systems with Applications* 143: 113042. [CrossRef]

Watkins, Christopher J. C. H., and Peter Dayan. 1992. Q-learning. *Machine Learning* 8: 279–92. [CrossRef]

Weigand, Alois. 2019. Machine learning in empirical asset pricing. *Financial Markets and Portfolio Management* 33: 93–104. [CrossRef]

Xing, Frank Z., Erik Cambria, and Roy E. Welsch. 2018. Natural language based financial forecasting: A survey. *Artificial Intelligence Review* 50: 49–73. [CrossRef]

Xue, Jingming, Qiang Liu, Miaomiao Li, Xinwang Liu, Yongkai Ye, Siqi Wang, and Jianping Yin. 2018. Incremental multiple kernel extreme learning machine and its application in robo-advisors. *Soft Computing* 22: 3507–17. [CrossRef]

Yun, Hyungbin, Minhyeok Lee, Yeong Seon Kang, and Junhee Seok. 2020. Portfolio management via two-stage deep learning with a joint cost. *Expert Systems with Applications* 143: 113041. [CrossRef]

Zhong, Xiao, and David Enke. 2019. Predicting the daily return direction of the stock market using hybrid machine learning algorithms. *Financial Innovation* 5: 1–20. [CrossRef]