



Machine learning in empirical asset pricing

Alois Weigand¹

Published online: 26 February 2019
© Swiss Society for Financial Market Research 2019

Abstract

The tremendous speedup in computing in recent years, the low data storage costs of today, the availability of “big data” as well as the broad range of free open-source software, have created a renaissance in the application of machine learning techniques in science. However, this new wave of research is not limited to computer science or software engineering anymore. Among others, machine learning tools are now used in financial problem settings as well. Therefore, this paper mentions a specific definition of machine learning in an asset pricing context and elaborates on the usefulness of machine learning in this context. Most importantly, the literature review gives the reader a theoretical overview of the most recent academic studies in empirical asset pricing that employ machine learning techniques. Overall, the paper concludes that machine learning can offer benefits for future research. However, researchers should be critical about these methodologies as machine learning has its pitfalls and is relatively new to asset pricing.

Keywords Machine learning · Big data · Empirical asset pricing

JEL Classifications G12 · G13

1 Introduction

The concept of machine learning is not new. Its origination goes back to the 1950s when, for example, Turing (1950) and Samuel (1957) conducted pioneering machine learning research with simple algorithms. However, these insights and additional promising findings over the next decades hardly found any applications outside the field of computer science and software engineering due to constraints on computing and data. These hurdles disappeared due to technological changes in the twenty-first century. For instance, Arnott et al. (2018) highlight the speed of these

✉ Alois Weigand
alois.weigand@unisg.ch

¹ Swiss Institute of Banking and Finance (s/bf), University of St. Gallen, Unterer Graben 21, 9000 St. Gallen, Switzerland

transformations with the comparison of the Cray 2 the fastest supercomputer in the world in the late 1980s and early 1990s and the iPhone Xs. At the time, the Cray 2 was able to perform 1.9 billion operations per second, weighed 5500 pounds and cost over US\$ 30 million adjusted for inflation. Today's iPhone Xs is capable of 5 trillion operations per second, weighs six ounces and costs US\$ 1449. Moreover, the authors highlight the difference of the costs in data storage. A gigabyte of storage cost US\$ 10,000 in 1990, whereas it costs only about a penny today. Arnott et al. (2018) further mention the availability of "big data" and the broad range of open-source application software as triggers for the new wave of research that uses machine learning techniques in various scientific fields. Inter alia, machine learning is now applied in financial research as well.

Most recently, Gu et al. (2018) want to justify the growing role of machine learning in financial research by synthesizing it with modern empirical asset pricing research. More precisely, the authors use the research agenda of understanding the dynamics of market equity risk premiums as the basis for a comparative analysis of methods in the machine learning repertoire. These methods include generalized linear models, dimension reduction techniques, boosted regression trees, and random forests. The researchers find that machine learning tools improve the description of expected return behavior relative to traditional forecasting methods. In particular, penalization, dimension reduction as well as nonlinearities improve the results. In addition, Gu et al. (2018) identify trees and neural nets as the best performing methods and point out that all machine learning approaches agree on a small set of dominant predictive signals, which include variations on momentum, liquidity, and volatility. Moreover, the gap between machine learning approaches and traditional models widens further when the authors predict portfolio returns. These results show, according to the authors, that improved risk premium measurement through machine learning can simplify the investigation of economic mechanisms of asset pricing and that machine learning is a promising approach for new financial technologies.

To define the true role of machine learning in empirical asset pricing, this paper takes the study of Gu et al. (2018) as a starting point and defines machine learning in this specific setting in a first step. In a second step, it is shown why empirical asset pricing is a particularly attractive field for machine learning applications and what machine learning is not able to do in this context. In a third step, the results of several asset pricing studies that employ machine learning are presented and discussed in a detailed literature review. Implications for future research are outlined as well. Hence, this paper contributes to the sparse literature around empirical asset pricing via machine learning. It does so by condensing theoretical insights and empirical results of recent academic publications as well as working papers to give the reader a condensed overview of the topic. Therefore, this literature review does not dig into the technical aspects of machine learning.¹

¹ The author refers the interested reader to Gu et al. (2018) who provide a detailed description of machine learning tools for empirical asset pricing. These explanations start from scratch and cover the statistical model specification as well as programming guidelines of different methods. Furthermore, Dey (2016) reviews and explains different machine learning algorithms in detail.

2 Defining machine learning

In general, machine learning is a method used to teach machines how to handle data more efficiently while performing a certain task (Dey 2016). However, using a general definition of machine learning is not always suitable. The description should be tailored to the specific field in which these techniques are used in. Therefore, Gu et al. (2018) outline a context-specific definition of machine learning in empirical asset pricing:

- (1) A diverse collection of high-dimensional models for statistical prediction, combined with (2) so-called regularization methods for model selection and the mitigation of overfit, and (3) efficient algorithms for searching among a vast number of potential model specifications.

The three elements in the definition of Gu et al. (2018) describe the concept very well. The high-dimensional nature (1) of machine learning brings high flexibility in approximating the underlying asset returns. According to Mullainathan and Spiess (2017), this flexibility, however, goes in hand with a higher danger of overfitting. Therefore, machine learning approaches need to include refinements in its implementation that guard against overfit (2) as also mentioned by Cawley and Talbot (2010). Gu et al. (2018) as well as Hwang (2018) further emphasize the importance of the last element (3) as machine learning techniques should be designed to approximate an optimal specification with manageable computational cost.

3 Motivating machine learning

In machine learning, statistics as well as computer science are joining forces and are combining their strengths (Das and Behera 2017). Hence, machine learning holds considerable promise for financial applications if it is applied in the correct way (Arnott et al. 2018). Although there is a high danger of misapplying these techniques and to engage in data mining, machine learning may help in solving problems researchers always have faced in quantitative finance. In the case of predicting asset returns, three main problems, which are listed below, arise in the traditional asset pricing literature (Keim and Stambaugh 1986; Pesaran and Timmermann 1995; Torous and Valkanov 2000; Welch and Goyal 2008). Gu et al. (2018) also touch upon these difficulties and Kumar and Thenmozhi (2006) implicitly deal with these problems in their paper as well:

- Problem of Prediction
- Problem of Variable Selection
- Problem of Functional Form

The way machine learning may help to bypass these problems becomes evident when relating the difficulties to the elements (1 to 3) of the machine learning

definition of Gu et al. (2018) in the previous chapter. Firstly, the measurement of an asset's risk premium is a problem of prediction as the premium is a conditional expectation of a future realized excess return. A lot of machine learning tools are specialized in predictions as they offer a high degree of flexibility (1). Secondly, existing research proposes a variety of asset-level characteristics as well as macroeconomic predictors which may be used for modeling risk premiums (Green et al. 2013; Harvey and Liu 2016; Welch and Goyal 2008). Harvey et al. (2016) quantify this pool of possible variables by depicting the cumulative number of recently discovered predictors which increased sharply from 21 in 2003 to more than 240 in 2012. Moreover, these factors are often closely related and therefore highly correlated. In this regard, it is necessary that machine learning techniques help to optimize the degrees of freedom and to condense variation among predictors by emphasizing variable selection and by including dimension reduction tools (2). Thirdly, the ambiguity of the functional form of the relationship between the risk premium and the predictors is complicating the setting. The relationship in the model specification may not be linear and take a nonlinear form. In addition, interaction terms may be included in the correct model specification. Efficient algorithms (3) allow machine learning approaches to look for the optimal functional form at low computation costs. This specification search for a functional form may be very wide which allows to approximate complex nonlinear associations while avoiding overfit biases and false discovery.

Although machine learning may substantially improve the pricing of assets, Gu et al. (2018) and Arnott et al. (2018) also highlight what machine learning cannot do. These authors postulate that machine learning, unsupervised machine learning in particular, does not impose economic principles. Hence, machine learning cannot identify deep fundamental economic mechanisms or equilibria. These structures have to be added by economists, who also have to decide how the machine learning algorithms should work given this structure. This link between machine learning and equilibrium asset pricing is for example shown by Kelly et al. (2017) as well as Feng et al. (2017).

4 Reviewing the existing literature

Machine learning techniques are not yet widely spread in the literature of empirical asset pricing. Traditional methods still dominate the publications of the last years. For example, differences in expected returns across stocks are typically analyzed by cross-sectional regressions of future stock returns on a few lagged variables (Fama and French 2008; Lewellen 2015). Time-series forecasts of stock returns are usually done by time-series regressions of aggregated portfolio returns on a few macroeconomic variables (Kojien and Nieuwerburgh 2011). Nevertheless, machine learning tools are on the rise and a new wave of research employs these approaches. These papers are worth to have a detailed look at and allow understanding the usefulness of machine learning in asset pricing. In the following, the papers are grouped into sub-categories based on their field of application.

4.1 Pricing equities

A small number of studies examines the cross section as well as the time series of stock returns with machine learning tools. Harvey and Liu (2016) identify the best predictors by a bootstrap procedure that takes the possibility of time-series as well as cross-sectional dependence of factors into account. They find that the market factor is the dominant factor and a second important factor is linked to profitability. Nevertheless, the method of Harvey and Liu (2016) gives guidance for any further study that faces the problem of multiple testing, which will be a common obstacle in navigating the vast array of “big data”.

Another study in this field is the working paper of Giglio and Xiu (2017). The authors make use of a dimension reduction technique, a principal component analysis (PCA), to show that the risk premium of a factor may be discovered in a linear factor model regardless of the rotation of the other control factors as long as they together span the true factor space. This allows consistent and robust estimates of the risk premiums as it provides a systematic way to tackle the concern that a theoretical model is misspecified due to omitted variables (Giglio and Xiu 2017). The power of the authors’ six-factor model is demonstrated by the high correlation of predicted and realized average excess returns of differently sorted portfolios, which are sorted by size, momentum, industry, book-to-market value, beta, and variance among others. Hence, the PCA should be a widely used toolkit according to Giglio and Xiu (2017).

Similarly to Giglio and Xiu (2017), Kelly et al. (2017) test and estimate a factor return model by using a dimension reduction method. They use a method called Instrumented Principal Components Analysis (IPCA) which treats characteristics as instrumental variables for estimating dynamic loadings. This estimator is easy to work with as it reflects a standard PCA which additionally allows Kelly et al. (2017) to bring information beyond just returns into the estimation of factors and betas. At the stock level, the authors find that three IPCA factors explain the cross section of average returns significantly better than existing factor models. Moreover, Kelly et al. (2017) show that among the large collection of stock characteristics in the existing literature only seven are statistically significant at the 1% level in the IPCA model— market beta, size, earnings-to-price, book-to-market, assets-to-market, short-term reversal, and momentum. These characteristics further account for nearly 100% of the model accuracy. The authors postulate through machine learning that the key for a successful factor model is including information from stock characteristics into the estimation of factor loadings. Thus, machine learning allows Kelly et al. (2017) to develop a new research protocol for evaluating hypotheses about patterns in asset returns.

Kelly and Pruitt (2015) further show the advantages of regularization methods and data compression for pricing equities as the authors introduce the three-pass regression filter (TPRF). The methodology of this approach is a constrained least squares estimator and reduces to partial least squares as a special case. The authors’ findings suggest that the TPRF is suitable for forecasting in a many-predictor environment as the methodology efficiently identifies a subset of predictors that is useful for forecasting. According to Kelly and Pruitt (2015), the TPRF is also superior to

other alternatives in a variety of simulation specifications and in empirical applications using financial and macroeconomic data.

Moritz and Zimmermann (2016) apply a machine learning technique to analyze the cross section of stock returns as well. They use tree-based models—one of the most powerful machine learning tools according to Gu et al. (2018)—in the context of portfolio sorting and relate information in past returns to future returns. By comparing this model to a simple, linear Fama–MacBeth framework, Moritz and Zimmermann (2016) show that the linear framework does not exploit all relevant information in the data and that their machine learning approach is more powerful. Moreover, a trading strategy based on the authors' findings has an information ratio twice as high as the Fama–MacBeth framework which takes two-way interactions into account. This result leads to the conclusion of Moritz and Zimmermann (2016) that their model around tree-based conditional portfolio sorts may significantly speed up the process of scientific discoveries in the field.

There are also two studies that use shrinkage and selection methods to estimate stochastic discount factors and a nonlinear function for expected returns. These studies by Kozak et al. (2018) and Freyberger et al. (2018) highlight that the quest of traditional factor models to summarize the cross section of stock returns with a sparse set of characteristic-based factors is hopeless. Kozak et al. (2018) point out that there is simply not enough redundancy among the wide range of proposed predictors for such a simple model to adequately price the cross section. Therefore, a stochastic discount factor needs to load on a large number of characteristic-based factors. This process is facilitated by machine learning, which may be fruitful for future research on the economic interpretation of the stochastic discount factor (Kozak et al. 2018).

Deep learning networks are also applied to analyze the cross section of stock returns. For example, Messmer (2017) uses a deep feedforward neural network based on a large set of firm characteristics to predict the US cross section of stock returns. After the author applies a network optimization strategy, he finds that the generated long-short portfolios outperform linear benchmarks, which underscores the importance of nonlinear relationships between firm characteristics and expected returns. However, Messmer (2017) does not claim that deep learning is the best way to exploit these nonlinearities. In addition, Messmer (2017) studies the main drivers of expected returns which are the short-term reversal as well as the twelve-month momentum.

One recent paper by Brogaard and Zareei (2018) uses machine learning to reveal stock mispricing which is characterized by inconsistent cross-sectional or time-series patterns in an empirical asset pricing model. Overall, Brogaard and Zareei (2018) show that mispricing still exists, but it has decreased over time which implies that markets have recently become more efficient. The researchers also show that their algorithm is more successful in finding anomalies than the benchmark of a set of moving-average strategies. Hence, Brogaard and Zareei (2018) conclude that machine learning is informative in financial economics if it avoids p-hacking, data snooping or data-dredging.

In addition, Feng et al. (2018) construct deep learning dynamic factor models for predicting asset returns. More specifically, the authors jointly estimate hidden factors and regression coefficients by stochastic gradient descent and, thus, provide an

alternative to dynamic factor modeling. Their findings underline the adaptiveness of machine learning tools by the flexibility of the return predictors in the model while maintaining a high out-of-sample R^2 . Nevertheless, Feng et al. (2018) mention the caveats of their approach by highlighting the difficulties in interpreting the model as well as in performing causal inference from large datasets.

4.2 Pricing derivatives

In addition to the pricing of equities, machine learning is used for pricing derivatives. In the simulation of Black–Scholes option prices, Hutchinson et al. (1994) employ neural networks to show that this machine learning tool is a powerful alternative when the underlying price dynamics are unknown or the no-arbitrage condition cannot be solved analytically. With a neural net, they are able to recover the Black–Scholes formula from a 2-year training set of daily option prices and use the resulting network formula successfully to both price and delta-hedge options. In a comparison with traditional methods (ordinary least squares, radial basis function networks, multilayer perceptron networks, and projection pursuit), the outperformance of the machine learning tools is persistent in terms of out-of-sample R^2 . Overall, Hutchinson et al. (1994) are cautiously optimistic when reflecting on their results as well as the future importance of their approach. However, the authors see their work as a reference point for future research and highlight promising directions how to improve their machine learning approach. This initial work of Hutchinson et al. (1994) is revisited and confirmed by Culkin and Das (2017) who program a feedforward neural net to reproduce the Black–Scholes option pricing formula. Moreover, the nonparametric neural network of Hutchinson et al. (1994) is tested and confirmed by Amilon (2003) with a more extensive dataset which allows to better capture the relationships between the derivative and the underlying.

Yao et al. (2000) also aim to forecast option prices with neural networks. However, their approach is slightly different as the authors use back-propagation neural networks. Moreover, Yao et al. (2000) find that different results in terms of accuracy are achieved by grouping the data differently. Overall, the authors highlight that for volatile markets a neural net option pricing model performs better than the traditional Black–Scholes model in terms of normalized mean-squared error in the case of in-the-money and out-of-the money data. However, the Black–Scholes model is still good for pricing at-the-money options. Therefore, Yao et al. (2000) conclude that it may be too early to postulate that the forecasting power of neural networks is better than the performance of conventional models. Nevertheless, the advantages of neural nets are obvious (Yao et al. 2000).

4.3 Predicting default

Khandani et al. (2010) and Butaru et al. (2016) employ machine learning tools to model credit card delinquencies and defaults. The study of Khandani et al. (2010) uses radial basis functions, tree-based classifiers, and support-vector machines. Khandani et al. (2010) consider these methods suitable for the setting of consumer

credit-risk analytics, which deals with large sample sizes as well as the complexity of the possible relationships among consumer transactions and characteristics. Hence, the authors believe that machine learning forecasts are highly adaptive and are able to measure the dynamics of changing credit cycles as well as the levels of default rates. Their results support this view as their out-of-sample forecasts are highly correlated with realized delinquencies in the following six or twelve months, respectively.

In the study of Butaru et al. (2016), decision trees and random forests are used as the main methodology. Further, the publication uses an extensive credit card data set which allows these methods to show their full potential. The authors find that decision trees and random forests work better than logistic regressions in both out-of-sample and out-of-time forecasts of credit card delinquencies. Therefore, the paper of Butaru et al. (2016) gives an in-depth illustration of the potential benefits that “big data” and machine learning tools bring to risk managers, shareholders, regulators, consumers, and all stakeholders that are affected by unexpected losses and the reduction in consumer credit costs. This also holds for the findings of Khandani et al. (2010).

Another recent working paper written by Sirignano et al. (2018) establishes a deep learning model of multi-period mortgage delinquency, foreclosure, and prepayment risk. The authors’ empirical findings yield insight into the behavior of mortgage borrowers. Sirignano et al. (2018) find that the relationship between borrower behavior and risk factors is nonlinear, which casts doubt on linear models applied in prior research. Additionally, the authors find evidence that prepayments involve the strongest nonlinear effects among all events. These results have significant implications for risk management, investment management, and mortgage securities. Moreover, in an out-of-sample analysis, Sirignano et al. (2018) show that by addressing nonlinearities, the accuracy of loan- and pool-level risk forecasts, the investment performance of mortgage trading strategies, and the valuation and hedging of mortgage-backed securities improves significantly. Hence, according to Sirignano et al. (2018), machine learning unlocks great potential in the field of default prediction.

4.4 Appraising real estate

Machine learning tools are used in real estate price estimations as well. In this context, a research proceeding of Chiarazzo et al. (2014) models real estate sales prices based on an artificial neural network. The authors highlight the usefulness of machine learning in the complex system of real estate where motivations, tastes, and budget availability often do not follow rational behaviors. Furthermore, property prices are affected by a wide range of external parameters given by the location as well as the environment. Chiarazzo et al. (2014) show that the model is successful in incorporating these characteristics. Nevertheless, further improvements to the model have to be made according to the researchers. For example, Chiarazzo et al. (2014) apply clustering methods to improve the statistical performance of the artificial neural network to successfully capture specific characteristics of groups of properties in the future.

In contrast, the methodology of Barr et al. (2017) is more developed as they efficiently estimate a home price index for single houses based on a gradient boosted model. These granular indices can then be aggregated over geographies of any kind (for example on ZIP code level). Barr et al. (2017) point out that this approach has several strengths compared to the commonly reported indices like the median or repeat sales indices which are computed for bigger regions or metropolitan areas only as they do not produce reliable results at the property level or ZIP code level. Hence, the authors point out that the ability to produce accurate indices at varying levels of granularity allows scientists to take multiple perspectives and to answer a broader range of questions. For example, if researchers are not sure how a property behaved in the past compared to neighboring houses, it is not possible to predict how it will behave in the future (Barr et al. 2017). Machine learning opens the door to answer these questions.

4.5 Other fields of application

The pricing of bonds is another promising field for machine learning. A recent working paper by Bianchi et al. (2018) measures bond risk premiums within a regression-based context by employing a variety of machine learning methods. The authors find that machine learning tools, and neural networks in particular, may be very useful in improving the measurement of bond excess returns. More precisely, machine learning techniques achieve a higher out-of-sample predictive R^2 compared to traditional data compression and penalized regression techniques. Furthermore, empirical findings by Bianchi et al. (2018) show that macroeconomic information has substantial out-of-sample forecasting power for bond excess returns when complex nonlinear features are introduced via machine learning methodologies.

International stock return predictability may also be investigated by machine learning tools. One recent study by Rapack and Zhou (2013) uses the least absolute shrinkage and selection operator to predict global equity market returns by lagged returns of individual countries. More specifically, the authors find that the machine learning approach outperforms the conventional prediction approach especially in non-US countries. Rapack and Zhou (2013) further confirm the superiority of machine learning tools over conventional methods as they highlight that forecasts can be significantly improved by accommodating model uncertainty and parameter instability as well as including economically motivated model restrictions, forecast combinations, diffusion indices, and regime shifts. These findings have important implications for the development of both asset pricing models and investment management strategies. However, the authors show that return predictability is linked to business-cycle fluctuations—a crucial insight for the definition of the training-sample period in future machine learning studies.

5 Conclusion

This literature review summarizes several scientific papers on machine learning in empirical asset pricing and related fields. These fields range from pricing equities, the prediction of global equity indices, derivative pricing, real estate appraisals, the measurement of bond risk premiums to the forecast of credit card as well as mortgage delinquencies. Different machine learning approaches are used and adapted to each of these specific settings to best suit the needs of the individual research goal. These machine learning approaches show potential in overcoming shortcomings of traditional methods in empirical asset pricing which are characterized by the general problems of prediction, variable selection as well as the definition of a suitable functional form. Moreover, machine learning effectively deals with the challenges of “big data”, which is important for future research as scientists are constantly trying to learn more of the large volumes of data that are available nowadays. Hence, the future of machine learning in asset pricing looks promising and may not just be a trend of the last decade. However, researchers have to use this new method with caution as machine learning in asset pricing is still at a very early stage and has its pitfalls. Among others, these potential shortcomings of machine learning include p-hacking, data snooping, or data-dredging and difficulties in interpreting the statistical models as well as problems in performing causal inferences from large datasets. Therefore, future research should consider existing research protocols and learn from mistakes previous studies have made.

Acknowledgements The author thanks Prof. Dr. Manuel Ammann as well as Prof. Dr. Markus Schmid for their constructive and insightful comments.

References

- Amilon, H.: A neural network versus Black–Scholes: a comparison of pricing and hedging performance. *J. Forecast.* **22**, 317–335 (2003)
- Arnott, R., Harvey, C., Markowitz, H.: A backtesting protocol in the era of machine learning. Working paper, Duke University (2018)
- Barr, J., Ellis, E., Kassab, A., Redfearn, C.: Home price index: a machine learning methodology. *Int. J. Semant. Comput.* **11**, 111–133 (2017)
- Bianchi, D., Büchner, M., Tamoni, A.: Bond risk premia with machine learning. Working paper, Warwick Business School (2018)
- Brogaard, J., Zareei, A.: Machine learning and the stock market. Working Paper, University of Utah (2018)
- Butaru, F., Chen, Q., Clark, B., Das, S., Lo, A., Siddique, A.: Risk and risk management in the credit card industry. *J. Bank. Finance* **72**, 218–239 (2016)
- Cawley, G., Talbot, N.: On over-fitting in model selection and subsequent selection bias in performance evaluation. *J. Mach. Learn. Res.* **11**, 2079–2107 (2010)
- Chiarazzo, V., Caggiani, L., Marinelli, M., Ottomanelli, M.: A neural network based model for real estate price estimation considering environmental quality of property location. *Transp. Res. Procedia* **3**, 810–817 (2014)
- Culkin, R., Das, S.: Machine learning in finance: the case of deep learning for option pricing. Working Paper, Santa Clara University (2017)

- Das, K., Behera, R.: A survey on machine learning: concept, algorithms and applications. *Int. J. Innov. Res. Comput. Commun. Eng.* **5**, 1301–1309 (2017)
- Dey, A.: Machine learning algorithms: a review. *Int. J. Comput. Sci. Inf. Technol.* **7**, 1174–1179 (2016)
- Fama, E., French, K.: Dissecting anomalies. *J. Finance* **63**, 1653–1678 (2008)
- Feng, G., Giglio, S., Xiu, D.: Taming the factor zoo. Working Paper, City University of Hong Kong (2017)
- Feng, G., He, J., Polson, N.: Deep learning for predicting asset returns. Working Paper, City University of Hong Kong (2018)
- Freyberger, J., Neuhierl, A., Weber, M.: dissecting characteristics nonparametrically. Working paper, University of Wisconsin-Madison (2018)
- Giglio, S., Xiu, D.: Inference on risk premia in the presence of omitted factors. Working Paper, University of Chicago (2017)
- Green, J., Hand, J., Zhang, X.: The supraview of return predictive signals. *Rev. Account. Stud.* **18**, 692–730 (2013)
- Gu, S., Kelly, B., Xiu, D.: Empirical asset pricing via machine learning. Working Paper, University of Chicago (2018)
- Harvey, C., Liu, Y.: Lucky factors. Working Paper, Duke University (2016)
- Harvey, C., Liu, Y., Zhu, H.: and the cross-section of expected returns. *Rev. Financ. Stud.* **29**, 5–68 (2016)
- Hutchinson, J., Lo, A., Poggio, T.: A nonparametric approach to pricing and hedging derivative securities via learning networks. *J. Finance* **49**, 851–889 (1994)
- Hwang, T.: Computational power and the social impact of artificial intelligence. Working Paper, Cornell University (2018)
- Keim, D., Stambaugh, F.: Predicting returns in the stock and bond market. *J. Financ. Econ.* **17**, 357–90 (1986)
- Kelly, B., Pruitt, S.: The three-pass regression filter: a new approach to forecasting using many predictors. *J. Econom.* **186**, 294–316 (2015)
- Kelly, B., Pruitt, S., Su, Y.: Some characteristics are risk exposures, and the rest are irrelevant. Working Paper, University of Chicago (2017)
- Khandani, A., Kim, A., Lo, A.: Consumer credit-risk models via machine-learning algorithms. *J. Bank. Finance* **34**, 2767–2787 (2010)
- Koijen, R., Nieuwerburgh, S.V.: Predictability of returns and cash flows. *Ann. Rev. Financ. Econ.* **3**, 467–491 (2011)
- Kozak, S., Nagel, S., Santosh, S.: Shrinking the cross section. Working Paper, University of Michigan (2018)
- Kumar, M., and M. Thenmozhi. 2006. Forecasting Stock Index Movement: A Comparison of Support Machines and Random Forest. Working paper, Indian Institute of Technology
- Lewellen, J.: The cross-section of expected stock returns. *Crit. Finance Rev.* **4**, 1–44 (2015)
- Messmer, M.: Deep learning and the cross-section of expected returns. Working Paper, University of St. Gallen (2017)
- Moritz, B., Zimmermann, T.: Tree-based conditional portfolio sorts: the relation between past and future stock returns. Working Paper, Ludwig Maximilian University Munich (2016)
- Mullainathan, S., Spiess, J.: Machine learning: an applied econometric approach. *J. Econ. Perspect.* **31**, 87–106 (2017)
- Pesaran, H., Timmermann, A.: Predictability of stock returns: robustness and economic significance. *J. Finance* **50**, 1201–1228 (1995)
- Rapack, D., Zhou, G.: Forecasting financial variables. *Handb. Econ. Forecast.* **2A**, 327–383 (2013)
- Samuel, A.: Some studies in machine learning using the game of checkers. *IBM J. Res. Dev.* **3**, 535–554 (1957)
- Sirignano, J., Sadhwani, A., Giesecke, K.: Deep learning for mortgage risk. Working Paper, University of Illinois (2018)
- Torow, W., Valkanov, R.: Boundaries of predictability: noisy predictive regressions. Working Paper, University of California (2000)
- Turing, A.: Computing machinery and intelligence. *Mind* **49**, 433–460 (1950)
- Welch, I., Goyal, A.: A comprehensive look at the empirical performance of equity premium prediction. *Rev. Financ. Stud.* **21**, 1455–1508 (2008)
- Yao, J., Li, Y., Tan, C.: Option price forecasting using neural networks. *Omega* **28**, 455–466 (2000)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Alois Weigand is currently a Ph.D. candidate and research assistant at the School of Finance of the University of St.Gallen, Switzerland. His research interests are empirical asset pricing and alternative investments with a special focus on real estate.