

Machine learning in finance: forecasting and trading

Maxime LELIEVRE

Tom MERY

Matteo PEDUTO

Abstract—Over the last few years, the cryptocurrencies became increasingly popular among the investors despite their sometimes high volatility. Likewise, machine learning has become a popular tool for improving decision making in financial markets. This paper presents the implementations and performances' comparison of different machine learning models on several assets. The results suggest that the combination of models optimized on different assets of the same type reaches the best trading performances.

I. INTRODUCTION

Over the last few years, the cryptocurrencies became increasingly popular as an investment product and for a portfolio diversification strategy. In the period from February 2016 to December 2022, the estimated market capitalization of all cryptocurrencies passed from 27 to 800 billion USD and the 24 hour average trading volume of all cryptocurrencies already reached more than 100 billion USD in 2021. A still increasing body of literature focused on the pertinence of the efficient market hypothesis (EMH), as proposed by Fama (1970) [1], on the famous Bitcoin, see for example Urquhart (2016) [2], Nadarajah and Chu (2017) [3], Bariviera (2017) [4], Sensoy (2018) [5] and Wildi et al.(2019) [6]. In essence, the EMH postulates that efficient markets reflect all past, public or public and private information in market prices. Verification of the EMH is important for market participants as it implies that such information cannot be used to make persistent profits on trading on the market. In summary, recent research on the topic is inconclusive as to whether Bitcoin markets are efficient under the EMH or not. Wildi et al.(2019) results suggest that Bitcoin markets are becoming less rather than more efficient towards the sample end of their data (2019).

In this context, we propose in the continuation of Wildi et al.(2019) to extend their approach to other cryptocurrencies and to commodities and market indices to see whether positive trading performances can be achieved with machine learning models.

Four different machine learning methods are implemented and applied on three types of assets (cryptocurrencies, commodities and market indices) with three different assets for each. We further analyze whether the combination of several models trained on one data set gives better results than the results of the best model only. We determine also if the combination of the best models trained on several assets of the same type gives better results than the results of the best model trained exclusively on the asset. We test their performances with several trading performance metrics.

After a presentation of the used data, the pre-processing step and an explanation of some financial prerequisites, we start our analysis by unfolding the four machine learning methods implemented and then discuss the results.

II. DATA ANALYSIS

A. Data collection

The project thus focuses on nine different assets divided in three categories (cryptocurrencies, commodities and market indices). The same time period has been considered with respect to the category.

	Cryptocurrencies	Commodities	Market Indices
Assets 1	Bitcoin	Gold	S&P 500
Assets 2	Ethereum	Natural Gas	CAC40
Assets 3	Ripple	Oil	SMI
Period considered	2017-11-09 to 2022-12-13	2012-12-13 to 2022-12-09	2012-12-13 to 2022-12-12

TABLE I
THE NINE ASSETS DIVIDED BY TYPE

The collection of data has been performed on different web sites hereafter detailed with the assets' symbol used in the code.

Bitstamp [7]: Bitcoin (BTC-USD), Ethereum (ETH-USD), Ripple (XRP-USD)

Nasdaq [8]: Natural Gas (NYMEX-NG), Gold (LBMA-GOLD), Oil (OPEC-ORB), S&P500 (SP500)

Yahoo Finance [9]: CAC40 (CAC40), SMI (SMI)

In time series machine learning, the data cannot be split in any ways. In fact, it is important to maintain the chronological aspect. The train set must be observed before in time than the validation and test set. This characteristic makes complicated the cross validation process. This is why the data has been split chronologically. The first half for the training, the third quarter for the validation and the last quarter for the testing. The validation set is used to tune the hyper-parameters.

B. Data pre-processing

The different data sets extracted give several information about an assets' value on each day (open, close, adjusted close prices,...). The models are trained to forecast the log-returns of the open price of the assets as it is the most realistic one if one wants to take position based on the model's prediction. The open price reflects the price at which an asset first trades when the market opens. The log return at time t is computed as follow:

$$r_t = \log \left(\frac{P_t}{P_{t-1}} \right) \quad (1)$$

Going through the log return allows to transform the time series such that the resulting financial data become more stationary from a temporal perspective, meaning even if we shuffle the data order, we will still be able to properly train the

model and achieve successful test performance. It also allows to determine right away if the asset's price has increased or decreased over the day, by checking the sign of the log return, which is what matters in a trading perspective.

Moreover, the models are using a determined number of lags. This latter represents the number of consecutive log-returns considered to predict the asset's log-return at time $t+1$. The number of lags is determined by applying the auto-correlation function (ACF) on the log returns, see Fig.1. The ACF defines how data points in a time series are related, on average, to the preceding data points. In other words, it measures the self-similarity of the signal over different delay times. This is performed because, under the assumption that the time-series is gaussian and stationary, a linear forecasting model should reach a better accuracy if a pattern can be identified.

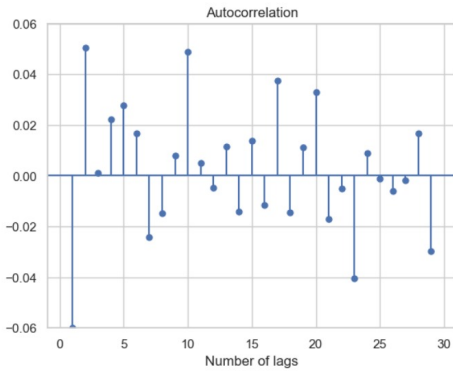


Fig. 1. Auto-correlation function applied on the log-returns of Bitcoin

Thus the data set is first transformed with the log returns before being split between the train, validation and test set. Finally, it is gathered in sequences of a chosen number of lags. For example, for 6 lags, one input is a sequence of 6 consecutive log returns and the goal is to forecast the following log return.

III. FINANCIAL PREREQUISITES

Before unfolding the implemented models, some financial basics must be explained to ensure the reader's understanding of the methodology used. Trading is the buying and selling of financial instruments in order to make a profit. These instruments range from a variety of assets that are assigned a financial value that goes up and down – and you can trade on the direction they take. Trading is opposed to investing, which suggests a buy-and-hold strategy and we will use it to evaluate our trading strategy. A good trade then consists of correctly predicting if the financial asset's value will increase or decrease and take the corresponding action. If one predicts it to increase, one will buy the asset to sell it later at a higher price. Conversely, if one predicts the value to decrease, one will sell the asset -it is called to short- and buy it later for a lower price.

We here give more details about the efficient market hypothesis (EMH) mentioned in the introduction. The EMH

postulates that efficient markets reflect all past information (weakform), public information (semi-strong form), or public and private information (strong form) in market prices. Some findings suggest that Bitcoin markets, while inefficient in their early days, transitioned into efficient markets recently. Others find support for the adaptive market hypothesis (AMH), an alternative theory that builds on evolutionary principles and assumes markets and market efficiency evolve over time. Verification of the EMH is important for market participants as it implies that such information cannot be used to make persistent profits on trading on the market.

The trading strategy followed during this analysis consists of selling or buying based on the predictions output by the implemented models in a horizon of one day. It is useless in financial prediction to measure a model performance through the value of the loss function used during the training (here MSE for neural nets). Thus, the performances of the models are assessed using three financial metrics hereafter explained.

A. Hit rate

The hit rate can be considered as the accuracy of the models. In trading, the hit rate is typically defined as the number of winning or profitable trades over a period of time for a trading strategy, divided by the total number of trades over the same period, and expressed as a percentage. It is determined by checking the sign of the predicted log return with respect to the sign of the target, expressed as a log return too. So if the predicted and the target log-returns have the same sign then it is considered as a profitable trade.

$$Hit\ Rate = \frac{\sum(sign(output) \cdot sign(target) >= 0)}{len(output)} \quad (2)$$

B. Annualized Sharpe ratio

The Sharpe ratio compares the return of an investment with its risk. Sharpe ratios above 1 are generally considered “good,” offering excess returns relative to volatility.

To calculate the Sharpe ratio, which is a measure that balances drift and volatility aspects, investors first subtract the risk-free rate R_f from the asset's rate of return R_a , often using U.S. Treasury bond yields as a proxy for the risk-free rate of return. In this analysis, it is set up to 0 for simplicity. Then, the expectation of the excess return is divided by the standard deviation of the asset's excess return σ_a . For the annualized sharpe ratio, it is almost the same but the numerator is multiplied by the square root of the number of periods in a year (365 for the cryptocurrencies and 252 for the commodities and market indices due to the fact that the markets of these two latest are open only during the working days).

$$Annualized\ Sharpe\ Ratio = \sqrt{nb_periods} \cdot \frac{E[R_a - R_f]}{\sigma_a} \quad (3)$$

C. Maximum Drawdown

The maximum drawdown (MDD) is the maximum observed loss from a peak to a trough of a portfolio, before a new peak is attained. MDD is an indicator used to assess the relative

riskiness of one's trading strategy versus another. For risk-averse investors, a low MDD is preferred as they would avoid trading with periods of high losses.

$$MDD = \frac{Trough\ Value - Peak\ Value}{Peak\ Value} \quad (4)$$

IV. MODELS AND METHODS

This section unfolds the models and methods used to analyze weather positive trading performance can be achieved with machine learning. The first part of the analysis optimized the four methods, namely Neural Networks (NN), Convolutional Neural Networks (CNN), Long Short Term Memory (LSTM) and Random Forest (RF), on every single asset. Optimization of machine learning models in financial trading is a complicated topic where a lot of research is still ongoing. In addition, the main discoveries remain secreted for obvious economical reasons. Thus, the goal of the optimization performed in this project was not to obtain the best absolute trading metrics performances. In fact, the analysis of the architecture of the methods -mentioned hereafter- was not the main focus of the tuning. The structure of the methods is then the same for each single asset and was determined with satisfying results for the trading metrics, and in order to have approximately the same number of parameters independently of the method (103 for NN, 109 for CNN, 91 for LSTM). The four methods implemented are explained hereafter. The first section tries to reproduce the analysis of Wildi et al.(2019) using the same architecture of a feed-forward neural networks but applied on a broader time period. The three others analyze weather positive trading performances can be achieved with other methods. These methods can be split in two categories: Neural nets that include Feed-forward neural networks (NN), Convolutional Neural Networks (CNN), Long-Short-Term-Memory (LSTM) and Random Forest which is a different kind of machine learning method.

A. Neural Networks

Recall that the goal is to forecast the log-return at time (t+1) of an asset based on a sequence of the previous log-returns of the asset. Therefore, for the following neural networks we try to minimize the Mean-Squared-Error (MSE) loss during training. Since the trading strategy is only based on the sign of the forecasted output, one could simply train a binary classifier with a logistic loss. But forecasting the log-returns allow to get more information and the trading strategy could then be adjusted in consequence. The following networks have been trained using PyTorch library with GPU capability on Google Colab. The implementation is available here [10]. Also, since the considered data sets are small, it has been chosen to train the networks using full-batch gradient descent using Rprop optimizer. Rprop, short for resilient backpropagation, tries to resolve the problem that gradients may vary widely in magnitudes. For full-batch learning this problem can be solved by only using the sign of the gradient. Finally, for the task of forecasting financial time-series, neural networks are highly subjects to convergence toward local optimum. Therefore,

each of the following neural nets are trained 100 times with different random initialization to form an ensemble of models (ensemble learning). The output of the ensemble is then taken as the median of the output of the 100 sub-models (majority vote). The following section describe the architecture of each type of sub-models.

1) *Feed-Forward*: The feedforward neural networks are composed of two hidden layers of dimensions six and three. Each layers are followed by a Relu-activation function.

2) *Convolutional*: To expand the analysis further, a convolutional neural network is implemented with two hidden layers before being flattened to end up with a fully connected neural network. There are two convolutional layers. The first one with 8 channels, a stride of 1, a kernel size of 3 and no-padding. The second one with 4 channels, a stride of 1, a kernel size of 1 and no-padding. The Relu-activation function is used after each convolutional layer and the weights.

3) *LSTM*: A Long Short Term Memory (LSTM) is implemented. Two LSTM are concatenated and each layer has two hidden features.

B. Random forest

A random forest regressor is also implemented using Scikit-Learn library. The random forest model is composed of 10000 estimators.

C. Optimization

It was found that neural networks quickly overfits. To adress this problem, several regularization techniques have been experimented such that L1-regularization, L2-regularization, dropout or even early stopping by tracking the validation losses during training. None of them were conclusive, the only effect of these regularizations was to decrease the difference between train and validation losses by actually decreasing the performances on the train set. So finally, with our fixed architecture, the hyper-parameter that impacts the most the performances is the learning rate. Thus, the optimization is determined by finding the learning rate, via grid search for ten values evenly distributed in the log-interval [0.0001 ; 0.001], that gives the best hit rate on the validation set. It is, indeed, the most relevant metric in this project because the hit rate is often proportional to the sharpe ratio (it is usually not the case for strong variations). It is also important to notice that the hit rate oscillates depending on the number of training's epochs. A convergence, however, is noticed when 500 epochs are computed, see Fig.2. The training is then performed on this number of epochs to optimize the different methods. To speed-up the process, only 10 sub-models are assembled together. For the random forest, the tuned hyper-parameter is the number of maximum features, tuned via grid search again, taken during the bootstrapping step.

Once the best models found for each asset, the performance are computed on the test set and reported in section VI.

V. EXPERIMENTS

The goal here is to analyze weather the combination of several models trained on the same data set with different

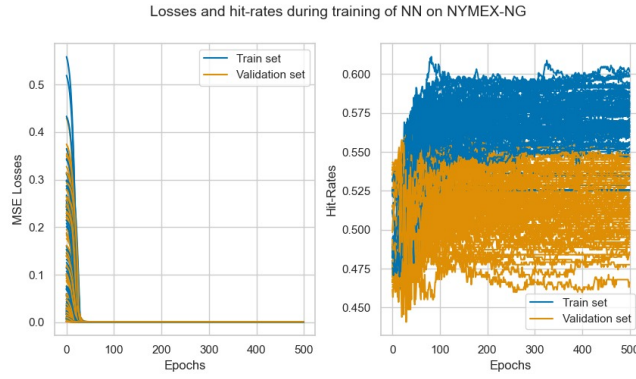


Fig. 2. Plot showing the convergence of the hit rate when computing 500 epochs

accuracy allows to make better predictions compared to the predictions of the best of the four models. One could think that each model learns something from the data in a possible different way that could contribute to improve the predictions when the results are combined.

In a second part, the analysis focused on all the nine assets to investigate whether the models' predictions are improved when one combines the best model trained on the assets of the same type. In other words, if what is learnt on a different asset of the same type (cryptocurrencies, commodities or market indices) can improve the predictions on one asset in particular.

The log return of the combination of the models on each data set is then measured by taking the mean, the median and absolute maximum of the log-returns of each model. The results reported in the Table II and III are the best values obtained between the latter three strategies.

VI. RESULTS

The following Table II and Table III resume the hit rates and Sharpe ratios of the best models for each asset as well as the comparison with the *Buy & Hold* strategy and two combinations of models. Combination 1 states for the combination of the 4 best models trained on the same asset. Combination 2 stands for the combination of the 3 best models trained on the asset of the same type. The results in bold represent the best model for each asset. The metrics of the combinations are highlighted in green when their performances are better than the best model alone. The *Buy & Hold* is in red when the models' and combination metrics performances could not reach higher values.

VII. DISCUSSION

The first result to notice is that positive trading performances can be achieved with one optimized model but it depends on the asset. The assets Oil and CAC40 do not show positive results when compared with the *Buy & Hold*. Overall, CNN and LSTM are the ones that perform the best. Random forest is also more accurate for the gold. The models are already satisfying because their architectures were not optimized. Concerning the combinations, the combination 1, which combines

	<i>B&H</i>	<i>NN</i>	<i>CNN</i>	<i>LSTM</i>	<i>RF</i>	<i>CI</i>	<i>C2</i>
<i>Bitcoin</i>	0.473	0.457	0.446	0.479	0.459	0.451	0.505
<i>Ethereum</i>	0.490	0.457	0.490	0.481	0.475	0.488	0.497
<i>Ripple</i>	0.490	0.497	0.497	0.532	0.470	0.514	0.525
<i>Nat. gas</i>	0.516	0.518	0.519	0.494	0.498	0.505	0.502
<i>Gold</i>	0.493	0.502	0.511	0.477	0.512	0.509	0.514
<i>Oil</i>	0.562	0.492	0.508	0.527	0.522	0.525	0.536
<i>S&P500</i>	0.564	0.553	0.554	0.558	0.548	0.572	0.533
<i>CAC40</i>	0.548	0.520	0.526	0.528	0.498	0.504	0.520
<i>SMI</i>	0.535	0.507	0.511	0.506	0.488	0.517	0.546

TABLE II

HIT RATE COMPARISON FOR EACH ASSET

	<i>B&H</i>	<i>NN</i>	<i>CNN</i>	<i>LSTM</i>	<i>RF</i>	<i>CI</i>	<i>C2</i>
<i>Bitcoin</i>	-1.163	-0.981	-2.066	-0.140	-2.020	-1.856	0.413
<i>Ethereum</i>	-0.870	-0.560	-0.327	1.123	0.245	0.308	-0.404
<i>Ripple</i>	-0.947	-0.143	0.521	1.191	-1.453	0.531	1.497
<i>Nat. gas</i>	0.743	0.523	1.107	-0.308	0.352	0.577	0.615
<i>Gold</i>	0.034	-0.056	-0.960	-0.468	0.305	-0.390	0.596
<i>Oil</i>	0.686	0.389	0.652	0.994	1.012	0.810	1.049
<i>S&P500</i>	0.518	0.833	0.801	0.768	0.090	1.136	0.551
<i>CAC40</i>	0.597	-0.192	0.702	0.499	-0.674	0.169	0.685
<i>SMI</i>	0.219	-0.588	-0.321	-0.367	-0.805	-0.400	0.522

TABLE III

SHARPE RATIO COMPARISON FOR EACH ASSET

the best optimized models of one asset, performs worse than the best model alone, except for the S&P500. This can be explained by the fact that some model learnt nothing meaningful and thus worsen the overall performance of the combination. However, the combination 2 shows significant improvement compared to the best model alone. Indeed, for five assets out of nine, it is the one with the best performances when looking at the performance metrics. Here again it depends on the type of assets. For the cryptocurrencies and the commodities, the combination 2 often performs the best which might be explained by the difference of inter-dependence of the corresponding assets. The cryptocurrencies and the commodities have both a worldwide scale whereas the market indices depend on a region, USA for S&P500 and Europe for CAC40 and SMI.

Overall, we observe that the cryptocurrencies are the assets that have the best results from the machine learning models when we look at the difference of performance between the *Buy & Hold* and the models. This last result is consistent with the results found by Wildi et al.(2019) [6] and questions the pertinence of the EMH applied on cryptocurrencies.

VIII. SUMMARY

To conclude, the results obtained suggest positive trading performances for almost each asset based on the comparison with the *Buy & Hold* strategy. In addition, the combination of the best models on the assets of the same type allows an extra-gain in the performance. The assets of the same type being similar, each model can learn a different trend and when combining in the accurate way, the trading performance results increased.

ACKNOWLEDGEMENTS

The authors thank Marc Wildi for hosting our project in his lab and his helpful suggestions.

REFERENCES

- [1] Fama EF (1970) Efficient capital markets: A review of theory and empirical work. *The Journal of Finance* 25(2):383–417
- [2] Urquhart A (2016) The inefficiency of bitcoin. *Economics Letters* 148:80 – 82
- [3] Nadarajah S, Chu J (2017) On the inefficiency of bitcoin. *Economics Letters* 150:6 – 9
- [4] Bariviera AF (2017) The inefficiency of bitcoin revisited: A dynamic approach. *Economics Letters* 161:1 – 42
- [5] Sensoy A (2018) The inefficiency of bitcoin revisited: A high-frequency analysis with alternative currencies. *Finance Research Letters*
- [6] Bundi N, Wildi M (2019) Bitcoin and Market-(In)Efficiency: a Systematic Time Series Approach.
- [7] For Bitcoin, Ethereum and Ripple.
- [8] For Natural Gas
For gold
For oil
For S&P500
- [9] For CAC40 and SMI.
- [10] Link to the github.