

Using Naïve Models to Improve US Dollar Exchange Rate Trend Prediction

Elia Yathie Matsumoto - EESP/FGV (elia.matsumoto@fgv.br)

Emilio Del-Moral-Hernandez - POLI/USP (emilio.delmoral@usp.br)

Claudia Emiko Yoshinaga - EAESP/FGV (claudia.yoshinaga@fgv.br)

Afonso de Campos Pinto - EESP/FGV (afonso.pinto@fgv.br)

Abstract—This paper extends previous research, which proposed a methodology based on the following hypothesis: dealing with the problem of predicting the next-day USD/BRL exchange rate daily trend, the existence of calendar effects allows us to improve trained voting-based ensemble models without model retraining. Despite the evidence of good results, in the present work, we propose adding naïve models to the originally proposed methodology because naïve models would also potentially benefit from the calendar effect becoming a benchmark to consider. The experiments confirmed that naïve models are not just challenging benchmarks but also models that can be included in the process to improve existing voting-ensemble models. On average, adding the naïve models to the original solution generated an increase higher than 100% in the value of the primary metric adopted for performance measurement. Constantly overwhelmed by more complex solutions, we can take these outcomes as a reminder not to neglect simplicity.

Keywords— *USD/BRL Exchange Rate, Behavioral Finance, Machine Learning, Ensemble Models, Naïve Model.*

I. INTRODUCTION

In recent years, we have witnessed a boost in research proposing new computational models based on increasingly larger and more complex architecture. Moreover, some huge models have delivered outstanding positive results in some areas, such as Image Recognition [1] and Natural Language Processing [2]. All of this may convey a notion opposed to the classic "Occam's Razor" principle, which essentially states that "the simpler is better" [3]. Also, the latter statement carries the sense of competition between simpler and more complex methods.

In this perspective, conversely, this present work proposes to exploit possible uses of simpler models to improve complex models instead of competing with them by improving our previous work [4], in which we applied the concept of technique combination to address a financial time series forecasting problem. Specifically, we focused on predicting the trend of the daily quotations for the US dollar to the Brazilian real exchange rate (USD/BRL). This kind of task is still challenging, as researchers all over the world continue struggling to get good results even when using larger models consistently [5][6][7]. The technical literature attributes this hardship to the complex nature of the economic and financial phenomena [8]. Thus, even many decades after the famous "The Meese-Rogoff puzzle" article publication [9] on pricing forecasting, several papers are still adopting random-walk models as a typical benchmark [10][11][12].

Our previous article [4] proposed "a method to improve existing voting-based ensemble models trained to predict the next-day USD/BRL exchange rate trend with no need for retraining or other costly

computational tasks.". Despite achieving promising results, some aspects of the performed experiments made us question the robustness of the proposed method during the pandemic and whether it could overcome the simplest random model possible, the naïve model, which takes the last observed value as the prediction for the next.

In the current paper, we describe what we did to answer these questions: (1) we repeated the original experiments, just using more recent data to include the pandemic period; (2) we added the comparison with the naïve models' outcomes and verified that this model category produced the best final result, surpassing all the others; (3) finally, we added the naïve models to the ensemble models and this final combination provided very satisfactory results.

This document has four more sections besides this introduction. In Section 2, we highlight the main points of the original paper related to the present work. In Section 3, we detail the methodology we used to compare the experiments. Section 4 describes the new experiments performed in this research and discusses their numerical results. The conclusions are presented in Section 5.

II. RELATED WORK HIGHLIGHTS

As already stated, this paper is related to previous research [4] focused on forecasting the USD/BRL trend, a subject well explored by the technical literature [13]. The novelty we brought was adding a Behavioral Finance element, the calendar effect, to construct a strategy to improve trained voting-based ensemble models without retraining. It only presumes the availability of the past predictions produced by the trained voting-based ensemble model and all the individual models that compose it.

Assuming the calendar effects affect investors' actions, the rationale is that the collective actions of investors constantly trying to raise their trading positions would induce periodic and deterministic patterns in the financial time-series movement. Accordingly, the forecasting models built to predict these financial time-series trends would also be prone to induce periodic and deterministic patterns in their outcomes, allowing us to identify the 'best model' among all available models based on their past performance. Hence, to improve the voting-based ensemble model outcomes, the study proposed replacing its original averaged prediction with the one produced by the model identified as 'best' for each trading day.

Since months have different numbers of days, the regular date numbering was not convenient to compare daily metrics on a month-by-month basis. Thus, to better evaluate the models' past performance, each trading day received a 'tag' based on its relative position to the month's first and last trading days. This 'tag' index, composed of 22 values, grouped the trading days into 22 groups, allowing us to compare the models' past performance more uniformly. The model composed of the outcomes of the 'best model' among all available models was named the 'Best Model per Tag' ensemble model, the BMT.

We applied the methodology in experiments using 15 years of USD/BRL observations, a total of 3,756 values from July 1, 2005, to June 30, 2020. The voting-based ensemble model (VOTING) was composed of five model categories: K-Nearest-Neighbor (KNN), Logistic Regression (LOGREG), MLP (Multi-Layered Perceptron Neural Networks), RF (Random Forest), and Support Vector Machine (SVM). Moreover, the binary random model (RND) based on sorting Bernoulli binary random variables was adopted as the minimum benchmark.

For models' performance comparison purposes, we adopted two metrics: (1) as the primary metric to assess the profit-generating potential of the predictions, the authors defined a metric (EARN), calculated according to Equation (1), to measure the earning of a theoretical USD/BRL trading strategy calculated using the predictions provided by the models; (2) to evaluate the predictions' correctness they adopted the Accuracy metric (ACC), the percentage of correct prediction.

$$EARN = \sum_{t=1}^n \hat{Y}_t * v_t \quad (1)$$

Where v_t is the USD/BRL variation on date t ; \hat{Y}_t is the prediction of the sign of v_t ; and n is the number of trading days. By multiplying the model prediction of the sign of v_t (\hat{Y}_t) by the USD/BRL variation observed on the date t (v_t), we obtain the outcome of trading US\$ 1.00 on the date $(t-1)$ to earn R\$ v_t on the date t . This value is positive if the prediction \hat{Y}_t is correct, and negative otherwise.

To summarize the outcomes of the original work, we display the average EARN and ACC values produced by each model built in the original experiments in TABLE I, henceforth named **Exp20**.

TABLE I: EARN and ACC averages ordered by EARN (descending order) in **Exp20**

Model category	EARN	ACC
BMT	1.832	0.531
VOTING	1.477	0.516
RF	1.166	0.518
MLP	1.156	0.507
KNN	0.788	0.521
SVM	0.473	0.510
LOGREG	0.193	0.510
RND ^a	-0.081	0.495

^a. Minimum benchmark.

In TABLE II, we display the EARN average values generated by all the models and the percentage of improvement provided by BMT over each of the other models' categories. On average, BMT generated EARN values 24% higher than VOTING, with 16% lower volatility. This outcome supported the central hypothesis of the research: the financial agents' collective actions affected by the trading calendar arguably induced deterministic patterns in the USD/BRL movement.

TABLE II: EARN average (in descending order) achieved in **Exp20** with the % of improvement provided by BMT

Model category	EARN	% of improvement provided by BMT
BMT	1.832	
VOTING	1.477	24.0%
MLP	1.166	57.1%
RF	1.156	58.5%
KNN	0.788	132.5%
LOGREG	0.473	287.3%
SVM	0.193	849.2%
RND ^b	-0.081	2361.7%

^b. Minimum benchmark.

Nonetheless, these positive results raised the suspicion that, under this hypothesis, a naïve model (that takes the previous value as the prediction of the next) would also benefit from the calendar effect and potentially produce reasonable results, making it a challenging benchmark to consider. We were also

concerned about the methodology's robustness during the pandemic. In the next section, we describe our methodology for investigating this guesswork.

III. METHODOLOGY

Motivated by the two questions stated at the end of the previous section, our purpose was to verify the robustness of the BMT during the pandemic and evaluate the naïve models (NAÏVE) performance. Thus, we carried out new experiments, working with a more recent timeframe and adding the NAÏVE outcomes to the process. Henceforth, we name:

- **Exp22**: the **Exp20** experiments reproduced with more recent data.
- **NExp22**: the **Exp22** experiments, including NAÏVE as a new benchmark and in the VOTING and BMT compositions.

For **Exp22**, we repeated the procedures of the original work building the VOTING ensembles composed of 5 model categories: KNN, LOGREG, MLP, SVM, and RF, and composing the BMT ensembles considering these 5 model categories plus VOTING. To verify the BMT robustness, we compared the EARN values generated by **Exp20** and **Exp22**, using RND as the minimum benchmark and ACC and EARN as performance metrics.

In **NExp22**, to verify the NAÏVE performance, we built the NAÏVE predictions equal to the previous output value and adopted them as the new minimum benchmark, keeping everything else identical to **Exp22**.

To verify the impact of the NAÏVE outcome inclusion in the process, we changed the VOTING and the BMT compositions by removing the model category that produced the outcomes with the highest volatilities in **Exp22** and replacing it with the NAÏVE. Lastly, we compared the **Exp22** and **NExp22** performance metrics. The following section details the experiments highlighted above and presents the numeric results.

IV. EXPERIMENTS DESCRIPTION & NUMERIC RESULTS

For comparison purposes, in this new set of experiments, **Exp22** and **NExp22**, we kept the same main components previously adopted in **Exp20**:

- The two performance metrics mentioned in Section 2 are the EARN metric, which measures the theoretical profit-generating potential provided by the models' outcomes, and the traditional Accuracy metric (ACC), which measures the models' overall performance.
- The five model categories listed in Section 2 (KNN, LOGREG, MLP, RF, and SVM) compose the VOTING.
- The RND category as the baseline benchmark.

A. Data Description and Preprocessing

We collected the data from the Institute for Applied Economic Research (IPEADATA)¹, extracting 3,998 daily raw observations from April 4, 2006, to March 31, 2022. The upper plot in Fig. 1 shows the daily USD/BRL (V_t), and the lower shows the curve of its variations ($v_t = V_t - V_{t-1}$).

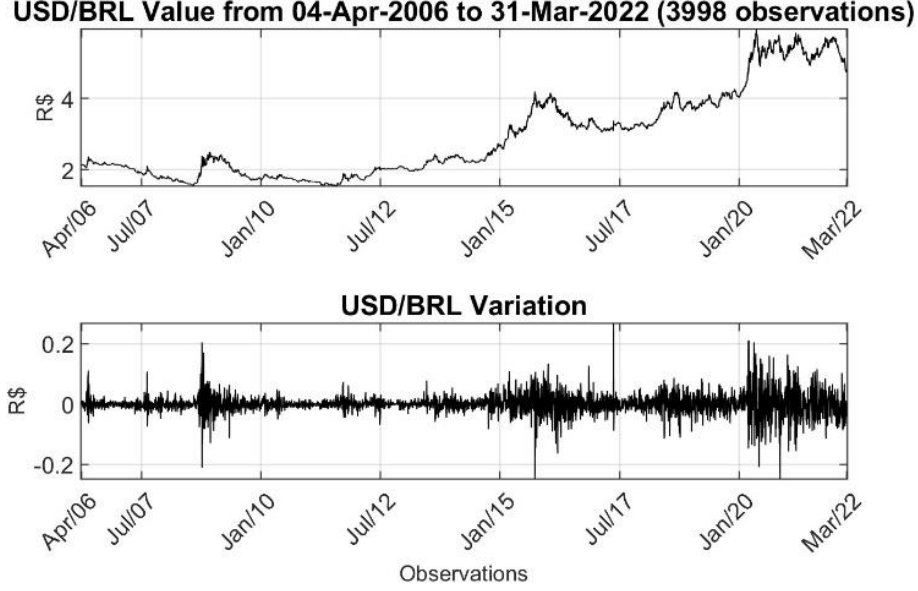


Fig. 1. (a) The USD/BRL daily values (upper); (b) The USD/BRL daily variation (lower).

TABLE III: USB/BRL variation (v_t) basic statistics

Basic Statistic	Value
Number of Observations	3.938
Mean	0.000
Standard Deviation	0.029
Minimum	-0.147
Maximum	0.155
Median	0.000
Volatility	61.775

We preprocessed the raw USD/BRL (V_t) variables by executing the following procedures:

- (1) **Trading days:** we considered only the trading days for the experiments, discarding weekends and holidays according to the Brazilian holiday calendar.
- (2) **Outliers:** as outlier identification criteria, we arbitrated a threshold value equal to 3.5 and discarded all the observations with USD/BRL log-variation values ($l_t = \log(V_t) - \log(V_{t-1})$) with Z-score values greater than it. In total, we discarded 59 observations.
- (3) **Stationarity:** we statistically tested and confirmed that the USD/BRL log-variation time-series (l_t) is stationary in mean using the Augmented Dickey-Fuller test with a 5% significance level and with the number of lagged difference terms equal to 252 (the number of trading days in one year).
- (4) **Date Tagging:** to better handle the difference in the trading days per month, we labeled the observations according to their relative position, using the first and last trading day of each month as references, using 22 tags defined as follows:

¹ Retrieved on April 10th, 2023, from <http://ipeadata.gov.br>.

- FIRST: First trading day of the month.
- LAST: last trading day of the month.
- F+N: trading day occurring n-days after the F trading day of the month, with N assuming integer values from 1 to 10.
- L-N: trading day occurring n-days before the L trading day of the month, with N assuming integer values from 1 to 10.

The basic statistics values of the cleaned-up USD/BRL daily variation time-series (v_t) are displayed in TABLE III.

B. Working Datasets

In this section, we describe how we composed the Working Datasets (WD), assuming all models have the same autoregressive formulation and use only time-series past values with no exogenous variables.

We defined the **output variable**, Y_t , as the sign of the USD/BRL variation, v_t , coding [-1] for negative variation and [+1] for positive variation. To compose the **input variables**, we arbitrated using the USD/BRL log-difference values (l_t), and worked with an observation window size equal to 10 (roughly half a month).

Hence, each observation in the dataset ended up with the final structure:

- Output variable, Y_t : the sign v_t on date t .
- Input Variables, X_t : (l_{t-1}, \dots, l_{t-10}), ten autoregressive values of the l_t Log-variation values.
- Tag label, Tag_t : the TAG of the date t .
- Return value, Ret_t : the v_t value on date t .

We divided the complete dataset into 132 sliding (or rolling) windows, defining 132 WDs. Each WD is composed of a training dataset with 60 months (5 years) of observations and a test dataset with a month of observations.

TABLE IV: Samples of the Training and the Test WD defined by sliding windows divisions (date and size)

WD ID	Training		
	Initial date	Final date	# Obs.
1	May.16.2006	Mar.31.2011	1083
2	Jun.08.2006	Apr.29.2011	1098
...
131	Mar.01.2017	Jan.31.2022	1128
132	Apr.03.2017	Feb.25.2022	1119
WD ID	Test		
	Initial date	Final date	# Obs.
1	Apr.01.2011	Apr.29.2011	19
2	May.02.2011	May.31.2011	22
...
131	Feb.01.2022	Feb.25.2022	14
132	Mar.02.2022	Mar.31.2022	22

We set the initial timeframe from May 16th, 2006, to March 31st, 2011. Shifting the previous time frame a month ahead defined the following sliding windows. Then, we trained all the models using 60 months of observations to predict the observations of the following month. Table IV displays a sample of the sliding window divisions and the size of the WDs. Due to space constraints, we listed only 4 of 132 Training and Test WDs.

C. Exp22: Reproducing Exp20 with more recent Data.

For each of the 132 WDs defined as described in the previous section, we repeated the procedures specified in **Exp20** training five prediction model categories (KNN, LOGREG, MLP, RF, and SVM) according to the best practices acknowledged by the Machine Learning community [14].

These five model categories were used to compose the VOTING, using the criteria of majority voting. To be referenced as the minimum benchmark, we built the RND using a Bernoulli binary random variables generator function to produce the predictions.

As noted in Section 2, we used the tag labels to separate the models' daily outcomes into 22 tag groups to obtain the BMT predictions. We then compared the ACC values generated by the models in the past per each of these 22 tags.

We compared the five individual models' past ACC values and the VOTING to compose the BMT. In the event of a tie, we arbitrated the selection of the VOTING outcome. As past, we arbitrated the timeframe of the previous 12 WDs (predictions for the previous year), never including them in the model comparison process of the following WDs, preventing look-ahead bias problems.

Hence, the initial period corresponding to the first 12 WDs was not included in the overall evaluation. For this reason, we named 'Test timeframe' the period from the 13th to the 132nd, when we have the BMT predictions. TABLE V shows a summary of this division of timeframes.

To compare the ACC values of the eight model categories (BMT, VOTING, KNN, LOGREG, MLP, RF, SVM, and RND), we concatenated the predictions obtained using the testing datasets from the **Test Timeframe**:

- 120 WDs (from the 13th to the 132nd WD).
- 2,372 aligned predictions, covering a period of 120 months (10 years) from April 2nd, 2012, to March 31st, 2022.

To consistently compare the outcomes covering the whole **Test Timeframe** period, we calculated the ACC and the EARN values for 1,500 one-day sliding windows with 873 daily observations each. Each of the one-day sliding windows could be interpreted as the earnings path of a financial agent using the models' prediction to trade USD/BRL for 873 trading days (roughly equivalent to 41 months). TABLE VI shows a sample of the sliding window division.

TABLE VII displays the average EARN and ACC values produced by the eight model categories considering the 1,500 sliding windows. We listed the models sorted in descending order of the EARN average value.

TABLE V: Timeframe division considered for performance evaluation

Timeframe	WD	Initial/Final Dates	# of Obs. (# years)
Total	1 st to 132 nd	Apr.01.2011 / Mar.31.2022	2623 (11 years)
Initial "past"	1 st to 12 th	Apr.01.2011 / Mar.30.2012	251 (1 year)
Test	13 th to 132 nd	Apr.02.2012 / Mar.31.2022	2372 (10 years)

TABLE VI: Sliding window division sample

Sliding Window	Initial date	Final Date
#1	Apr.02.2012	Sep.18.2015
#2	Apr.03.2012	Sep.21.2015
...
#1499	Oct.16.2018	Mar.30.2022
#1500	Oct.17.2018	Mar.31.2022

TABLE VII: EARN and ACC averages ordered by EARN (descending order)

Model category	EARN	ACC
BMT	1.301 (G)	0.532 (G)
VOTING	1.203 (G)	0.523 (L)
RF	1.143 (G)	0.526 (G)
LOGREG	0.956 (G)	0.520 (G)
KNN	0.889 (E)	0.503 (L)
MLP	0.842 (G)	0.511 (L)
SVM	0.704 (G)	0.516 (G)
RND ^c	0.166	0.504

^c Minimum benchmark.

TABLE VIII: EARN average (in descending order) achieved in Exp22 with the % of improvement provided by BMT

Model category	EARN	% of improvement provided by BMT
BMT	1.301	
VOTING	1.203	8.2%
RF	1.143	13.8%
LOGREG	0.956	36.1%
KNN	0.889	46.3%
MLP	0.842	54.5%
SVM	0.704	84.9%
RND ^d	0.166	684.5%

^d Minimum benchmark.

We verified the model performance difference by applying the statistical test paired t-test with a 5% significance. We coded the test paired t-test outcomes besides the metric values as follows:

- (G) the model metric average is statistically greater than the next model on the list.
- (E) the model metric average is statistically equal to the next model on the list.
- (L) the model metric average is statistically lesser than the next model on the list.

TABLE II and TABLE VIII, respectively, display the percentage of improvement provided by BMT in **Exp20** and **Exp22**. As expected, the lower values in TABLE VIII indicate a drop in the models' performance when the pandemic period data is included. Despite this, BMT still generated the highest metric value.

D. *NExp22*: Adding the NAÏVE to *Exp22*

Using the same 132 WDs described in Section 4.2, we defined the NAÏVE outcome as $\hat{Y}_t = Y_{t-1}$, i.e., the prediction for the next USD/BRL variation signal is equal to the previous (or current) USD/BRL variation signal. Moreover, we adopted the NAÏVE outcomes as the new minimum benchmark.

As foreseen, on average, the EARN values generated by NAÏVE surpassed all the others, including the BMT's, as displayed in TABLE IX. Nevertheless, if we consider the model sequence ordered by ACC in descending order, NAÏVE would fall in the 5th position. This evidences that, on average, the NAÏVE predictions were correct for higher v_t values (USD/BRL variation) than the BMT, VOTING, RF, and LOGREG predictions. It also illustrates that only considering pure model performance metrics when dealing with trading applications may not be enough. TABLE X displays the percentage of improvement provided by the NAÏVE outcomes over the other models. Also, on average, the NAÏVE EARN volatility was 18% lower than that of the BMT.

TABLE IX: EARN and ACC averages ordered by EARN (descending order)

Model category	EARN	ACC
NAÏVE ^e	1.968 (G)	0.517 (L)
BMT	1.301 (G)	0.532 (G)
VOTING	1.203 (G)	0.523 (L)
RF	1.143 (G)	0.526 (G)
LOGREG	0.956 (G)	0.520 (G)
MLP	0.889 (L)	0.503 (L)
KNN	0.842 (G)	0.511 (L)
SVM	0.704 (G)	0.516 (G)
RND	0.166	0.504

^e. Minimum benchmark.

TABLE X: EARN average (in descending order) achieved in *NExp22* with the % of improvement provided by NAÏVE

Model category	EARN	% of improvement provided by NAÏVE
NAÏVE ^f	1.968	
BMT	1.301	51.2%
VOTING	1.203	63.6%
RF	1.143	72.1%
LOGREG	0.956	105.8%
KNN	0.889	121.3%
MLP	0.842	133.6%
SVM	0.704	179.6%
RND	0.166	1086.4%

^f. Minimum benchmark.

To observe the performance variation of the three main model categories (BMT, VOTING, and NAÏVE) over time, we plotted the curve of the ACC and the EARN values produced by the models in Fig. 3 and 4. Note that we are working with 1,500 one-day sliding windows, each representing the earning path of a financial agent using the models' prediction to trade USD/BRL for 873 days. In the graphs, we defined:

- The X-axis as the end date of the one-day sliding window.
- The Y-axis as the average value of the metric in the one-day sliding window.
- A horizontal dashed line to indicate minimum reference metric values:

- Y-axis equal to 0.5: below this line, the ACC average is below 0.5, indicating deficient performance.
- Y-axis equal to 0: below this line, the EARN average is negative, indicating trading loss.

For instance, the first point of the BMT curve in Fig. 2 is (Sep.18.2015, 0.5384). This information means that 0.5384 is the average ACC value achieved by BMT during the sliding window period ending on Sep.18.2015 (starting 873 working days earlier, on Apr.02.2012).

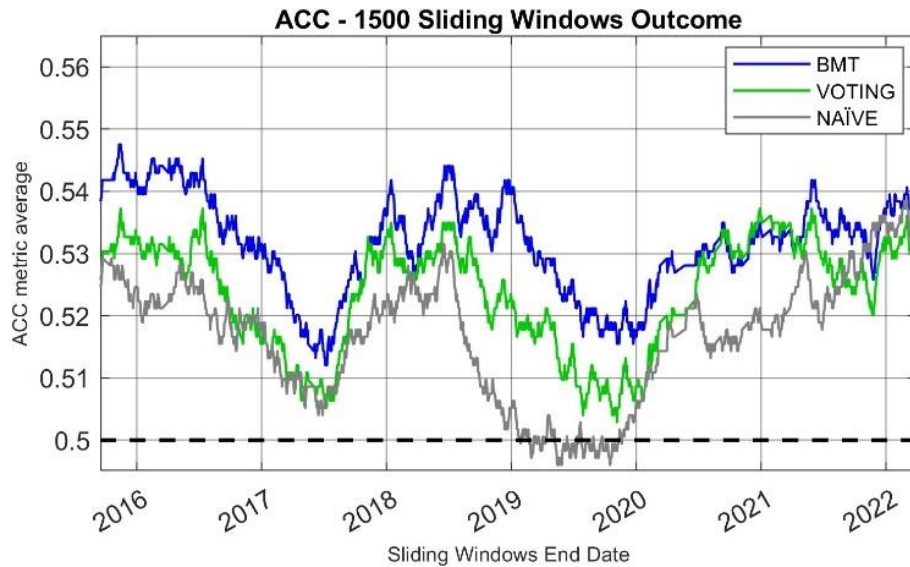


Fig. 2. ACC average values produced by BMT, VOTING, and NAÏVE.

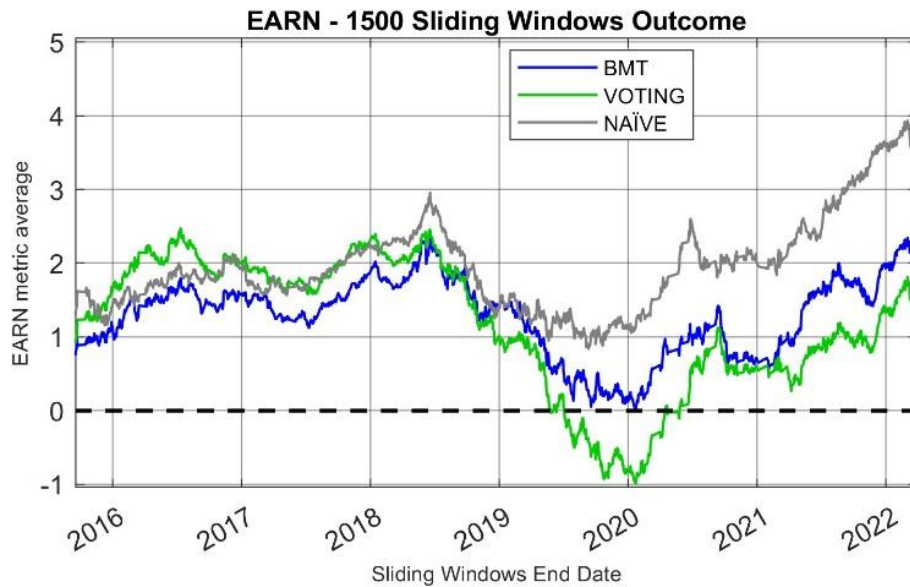


Fig. 3. EARN average values produced by BMT, VOTING, and NAÏVE.

Applying the same principle, the first point of the BMT curve in Fig. 3 is (Sep.18.2015, 0.7531). This information means that, on average, 0.7531 is the BRL(R\$) amount earned by a financial agent daily trading USB/BRL using the BMT predictions during the sliding window period ending on Sep.18.2015.

Fig. 3 shows that the NAÏVE EARN curve always stays above the BMT's and VOTING's even though, in Fig.2, the NAÏVE ACC curve positioning is just the opposite, showing, again, how essential it is to have metrics to measure the potential earnings associated with methodologies for trading applications.

As an attempt to take advantage of the good results produced by NAÏVE, and since neither VOTING nor BMT have any restrictions regarding the individual models' composition, we added the NAÏVE outcomes to the VOTING and the BMT compositions, in both cases, replacing the MLP's. We chose MLP because, on average, its outcome generated the EARN values with the highest volatilities.

TABLE XI shows the EARN and ACC average values produced by the ensemble models (VOTING and BMT) in **NExp22**, including the NAÏVE outcomes in their compositions. Once again, as anticipated, replacing MLP with NAÏVE improved the performance of both the ensemble models. TABLE XII displays the percentage of improvement provided by BMT for the set of experiments performed with more recent data and the inclusion of NAÏVE in the ensemble model compositions. Comparing the **Exp20** outcomes with that of **NExp22**, the percentage of improvement provided by BMT over VOTING increased by 51.3% (from 24.0% to 36.3%, values from TABLE II and XII, respectively).

TABLE XI: EARN and ACC averages ordered by EARN (descending order) – with NAÏVE included in the BMT and the VOTING composition

Model category	EARN	ACC
BMT	2.710 (G)	0.541 (G)
VOTING	1.989 (E)	0.525 (G)
NAÏVE ^g	1.968 (G)	0.517 (L)
RF	1.143 (G)	0.526 (G)
LOGREG	0.956 (G)	0.520 (G)
KNN	0.889 (E)	0.503 (L)
MLP	0.842 (G)	0.511 (L)
SVM	0.704 (G)	0.516 (G)
RND	0.166	0.504

^g Minimum benchmark.

TABLE XII: EARN average (in descending order) achieved in **NExp22** with the % of improvement provided by BMT – with NAÏVE included in the BMT and the VOTING composition

Model category	EARN	% of improvement provided by BMT
BMT	2.710	
VOTING	1.989	36.3%
NAÏVE ^g	1.968	37.7%
RF	1.143	137.1%
LOGREG	0.956	183.4%
KNN	0.889	204.7%
MLP	0.842	221.8%
SVM	0.704	285.1%
RND	0.166	1534.2

^h Minimum benchmark.

TABLE XIII: EARN statistics generated by BMT and VOTING in **Exp22** and **NExp22**, ordered by Average (in descending order)

Experiment / Model	EARN		
	<i>Avg.</i>	<i>Std.Dev.</i>	<i>Vol (Std.Dev./Avg.)</i>
NEXp2022 BMT	2.710	0.754	0.278
NEXp2022 VOTING	1.989	0.649	0.326
NEXp2022 NAÏVE	1.968	0.658	0.334
EXp2022 BMT	1.301	0.531	0.408
EXp2022 VOTING	1.203	0.899	0.747

To compare the **Exp22** and the **NExp22** outcomes, we displayed the EARN statistics generated by BMT and VOTING in the two sets of experiments in TABLE XIII. The positive impact of including the NAÏVE outcomes can also be confirmed by observing the graphs in Fig. 4 and 5, which respectively show the ACC and the EARN curves produced by NAÏVE, VOTING, and BMT with the inclusion of the NAÏVE predictions in their compositions. In Fig. 4, the VOTING and the BMT curves went up away from the NAÏVE curve and the horizontal dashed line, the minimum model performance indicator (ACC average equal to 0.5). Likewise, in Fig. 5, for all the 1,500 one-day sliding windows end dates, the VOTING and the BMT curves ended significantly above the horizontal dashed line (EARN average equal to 0.0).

According to the outcomes obtained in **NExp22**, the predictions provided by the VOTING and BMT models, with the inclusion of the NAÏVE predictions, were able to consistently generate positive earnings for 10 years (the Test Timeframe period, detailed in TABLE VI), with BMT in the first place.

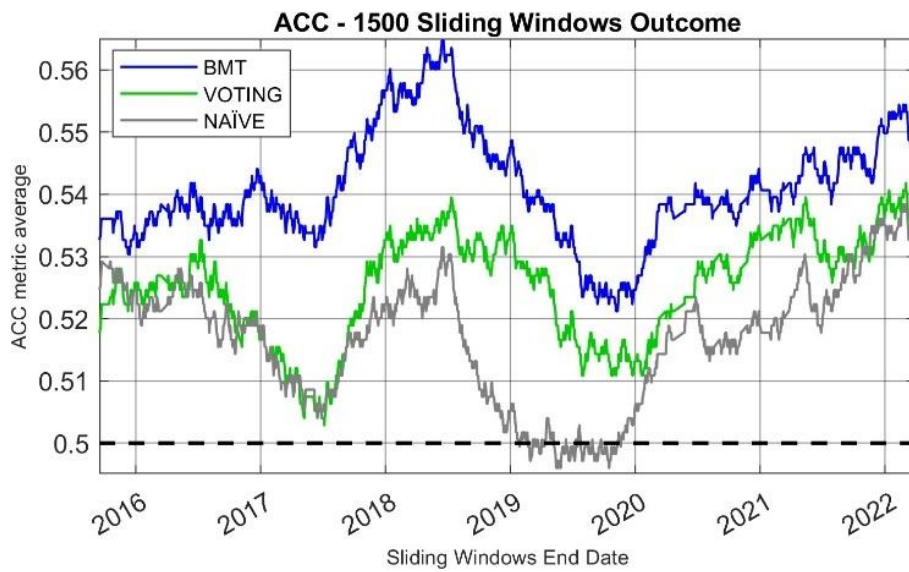


Fig. 4. ACC average values produced by BMT, VOTING, and NAÏVE with the new configuration.

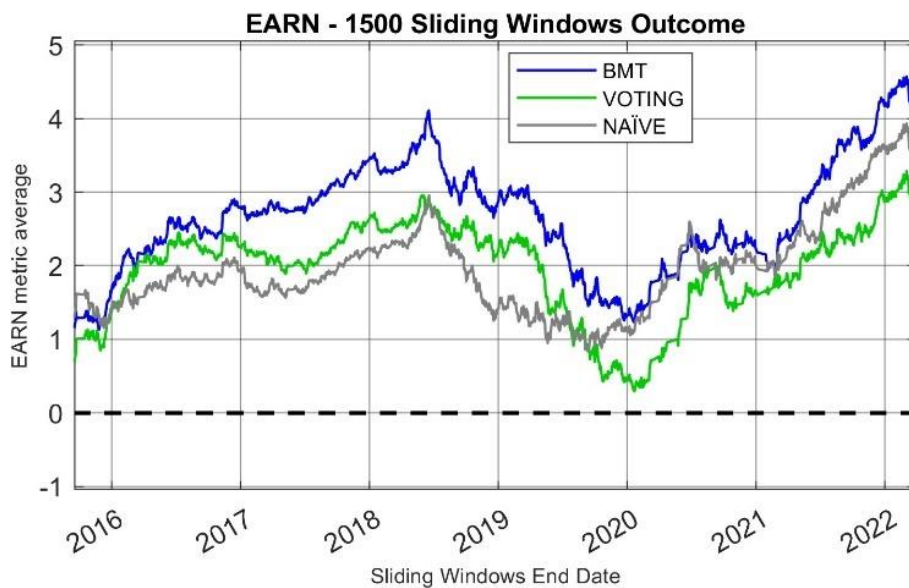


Fig. 5. EARN average values produced by BMT, VOTING, and NAÏVE with the new configuration.

V. CONCLUSIONS

Our prior research served as the starting point for this present paper. In this earlier study, we investigated the benefits of adding the Behavioral Finance Theory perspective to the Machine Learning framework applied to the problem of predicting the next-day USD/BRL trend by proposing a method to improve existing voting-based ensemble models (VOTING) trained to provide this prediction.

The basic idea is that the calendar effects (a market anomaly acknowledged by the Behavioral Finance theory) affect investors' decisions, inducing deterministic patterns in the USD/BRL time-series and the outcomes of the models to forecast them. Therefore, by tracking these models' past performance, it would be possible to estimate their future performance and choose, among all available models, the 'best' one to provide the next prediction. Hence, this research proposed applying this strategy to improve a voting-based ensemble model and named it the Best Model per TAG ensemble model, BMT.

Despite the promising BMT results, we conjectured whether the calendar effect would also positively affect the predictions provided by naïve models (NAÏVE). We also anticipated that positive naïve models' outcomes could help to improve the outcomes generated by VOTING and BMT. Further, we wanted to verify the robustness of the BMT during the pandemic. For this paper, we thus collected more recent data, from April 2006 to March 2022, to include the pandemic and applied the BMT construction methodology by adding the NAÏVE's outcomes in the process.

We carried out a set of experiments and, as suspected, the NAÏVE's outcomes surpassed all the others, including the BMT's; we, therefore, included the former in the VOTING and the BMT compositions, and, as also expected, we got better results: comparing to the original experiments, the percentage of improvement provided by BMT over VOTING increased 51.3% (from 24.0% to 36.3%). Moreover, by comparing only the new experiments' outcomes, we got an average increase of 65.3% for VOTING and 108.3% for BMT.

Once again, our results supported the original research's initial hypothesis: it is possible to use calendar effects to improve trained voting-based ensemble models without individual model retraining. Furthermore, we verified that the naïve model is not just a challenging benchmark; it is a model category that should be included in the BMT composition process to improve the outcomes of existing voting-ensemble models. We took this outcome as a reminder not to neglect simplicity.

We plan to continue our research by refining ensemble model construction methodologies and exploring the applicability of these methods to predict other financial assets in different financial markets.

REFERENCES

- [1] Y. Li, "Research and Application of Deep Learning in Image Recognition," *2022 IEEE 2nd Int. Conf. Power, Electron. Comput. Appl. ICPECA 2022*, pp. 994–999, 2022.
- [2] T. Wu *et al.*, "A Brief Overview of ChatGPT: The History, Status Quo and Potential Future Development," *IEEE/CAA J. Autom. Sin.*, vol. 10, no. 5, pp. 1122–1136, 2023.
- [3] P. Domingos, "The Role of Occam's Razor in Knowledge Discovery," *Data Min. Knowl. Discov.*,

vol. 3, no. 4, pp. 409–425, 1999.

- [4] E. Y. Matsumoto, E. Del-Moral-Hernandez, C. E. Yoshinaga, and A. de Campos Pinto, “Forecasting US Dollar Exchange Rate Movement with Computational Models and Human Behavior,” *Expert Syst. Appl.*, vol. 194, p. 116521, 2022.
- [5] S. Khuntia and J. K. Pattanayak, “Evolving Efficiency of Exchange Rate Movement: An Evidence from Indian Foreign Exchange Market,” *Glob. Bus. Rev.*, vol. 21, no. 4, pp. 956–969, 2019.
- [6] M. Argotty-Erazo, A. Blázquez-Zaballos, C. A. Argoty-Eraso, L. L. Lorente-Leyva, N. N. Sánchez-Pozo, and D. H. Peluffo-Ordóñez, “A Novel Linear-Model-Based Methodology for Predicting the Directional Movement of the Euro-Dollar Exchange Rate,” *IEEE Access*, vol. 11, pp. 67249–67284, 2023.
- [7] A. Adekunle, “Fiscal deficit and exchange rate movement: Empirical evidence from Nigeria,” *Acta Univ. Danubius. Œconomica*, vol. 19, no. 2 SE-Business Administration and Business Economics, pp. 7–20, May 2023.
- [8] M. López de Prado, *Advances in financial machine learning*. John Wiley & Sons, 2018.
- [9] R. A. Meese and K. Rogoff, “Empirical exchange rate models of the Seventies: do they fit out of sample?,” *J. Int. Econ.*, vol. 14, pp. 3–24, 1983.
- [10] E. F. Marçal and E. H. Junior, “Is it Possible to Beat the Random Walk Model in Exchange Rate Forecasting? More Evidence for the Brazilian Case,” *Brazilian Rev. Financ.*, vol. 14, no. 1, 2016.
- [11] A. R. S. Parmezan, V. M. A. Souza, and G. E. A. P. A. Batista, “Evaluation of statistical and machine learning models for time series prediction: Identifying the state-of-the-art and the best conditions for the use of each model,” *Inf. Sci. (Ny)*, vol. 484, pp. 302–337, 2019.
- [12] Y. Zhang and S. Hamori, “The Predictability of the Exchange Rate When Combining Machine Learning and Fundamental Models,” *J. Risk Financ. Manag.*, vol. 13, no. 3, p. 48, 2020.
- [13] M. C. Zorzi and M. Rubaszek, “Exchange rate forecasting on a napkin,” *J. Int. Money Financ.*, vol. 104, p. 102168, Jun. 2020.
- [14] A. L. Sicsú, A. Samartini, and N. L. Barth, *Técnicas de machine learning*. Sao Paulo/Brazil: Editora Blucher, 2023.