



A novel validation framework to enhance deep learning models in time-series forecasting

Ioannis E. Livieris¹ · Stavros Stavroyiannis² · Emmanuel Pintelas¹ · Panagiotis Pintelas¹

Received: 15 April 2020 / Accepted: 24 June 2020
© Springer-Verlag London Ltd., part of Springer Nature 2020

Abstract

Time-series analysis and forecasting is generally considered as one of the most challenging problems in data mining. During the last decade, powerful deep learning methodologies have been efficiently applied for time-series forecasting; however, they cannot guarantee the development of reliable prediction models. In this work, we introduce a novel framework for supporting deep learning in enhancing accurate, efficient and reliable time-series models. The major novelty of our proposed methodology is that it ensures a time-series to be “suitable” for fitting a deep learning model by performing a series of transformations in order to satisfy the stationarity property. The enforcement of stationarity is performed by the application of Augmented Dickey–Fuller test and transformations based on first differences or returns, without the loss of any embedded information. The reliability of the deep learning model’s predictions is guaranteed by rejecting the hypothesis of autocorrelation in the model’s errors, which is demonstrated by autocorrelation function plots and Ljung–Box Q test. Our numerical experiments were performed utilizing time-series from three real-world application domains (financial market, energy sector, cryptocurrency area), which incorporate most of global research interest. The performance of all forecasting models was compared on both problems of forecasting time-series price (regression) and time-series directional movements (classification). Additionally, the reliability of the models’ forecasts was evaluated by examining the existence of autocorrelation in the errors. Our numerical experiments indicate that our proposed methodology considerably improves the forecasting performance of a deep learning model, in terms of efficiency, accuracy and reliability.

Keywords Time-series · Deep learning · Spatiotemporal data · Reliability

1 Introduction

Time-series are encountered in a large variety of real-world applications, ranging from finance [10, 18] and commodities [8, 20] to healthcare [9, 29] and pollution management [4, 14]. Time-series data consist of discrete data points, obtained at successive equally spaced points in time. The main properties and characteristics of time-series data are responsible for distinguishing them from other types of data. More specifically, they frequently contain much noise, exhibit high volatility as well as extremal directional movements and possess a tendency for possible reversing these movements in the near-term future. Due to these significant characteristics, time-series forecasting is generally considered as one of the most challenging problems in data mining. As a result, the analysis of time-series data

✉ Ioannis E. Livieris
livieris@upatras.gr
Stavros Stavroyiannis
s.stavroyiannis@teipel.gr
Emmanuel Pintelas
ece6835@upnet.gr
Panagiotis Pintelas
pintelas@math.upatras.gr

¹ Department of Mathematics, University of Patras, 265-00 Patras, Greece

² Department of Accounting and Finance, University of the Peloponnese, 241-00 Antikalamos, Greece

41 has been an active subject of research for decades
42 [6, 7, 23, 30].

43 In the literature, the problem of time-series price and
44 movement forecasting has been comprehensively studied
45 for decades and numerous rewarding approaches have been
46 proposed. Traditional time-series methods such as ARIMA
47 (Auto-Regressive Integrated Moving Average) [5, 6] and
48 its variations as well as the more elaborated Machine
49 Learning methods [1, 3] probably constitute the most
50 famous and widely utilized methods for time-series pre-
51 diction. Nevertheless, these methods frequently do not
52 possess the ability to accurately model such complex data
53 and be successfully effective, since they cannot depict the
54 stochastic nature and high volatility of time-series.

55 During the last decade, the rapid advances in artificial
56 intelligence, as well as the vigorous developments in deep
57 learning techniques, attracted wide attention of scientific
58 and industrial communities for the development of efficient
59 and robust time-series forecasting models. Probably the
60 most popular and widely utilized deep learning methodol-
61 ogy is the development of an ANN-type network utilizing
62 convolutional and long short-term memory (LSTM) layers,
63 along with the classical dense layers. The former are uti-
64 lized to filter out the noise of the input data [11], and the
65 latter are tailored to efficiently capture complex temporal
66 dependencies and sequence pattern information by
67 exploiting their special architecture design [2]. Along this
68 line, researchers [18–20, 31] paid special attention to
69 exploit the advantages and benefits of both mentioned deep
70 learning techniques, proposing forecasting models utilizing
71 both convolutional and LSTM layers.

72 Recently, Pintelas et al. [24, 25] evaluated the perfor-
73 mance of several deep learning models for price and
74 movement forecasting of major cryptocurrencies. The
75 major novelty in their work was the application of a series
76 of tests for examining the prediction efficiency but mostly
77 the reliability of the models. In other words, they examined
78 whether the models have properly fitted the time-series
79 data and exploited all the available mined information,
80 during the training process. Based on their experimental
81 analysis, the authors stated that even the powerful deep
82 learning methodologies cannot guarantee the development
83 of reliable forecasting models. Additionally, they con-
84 cluded that new more sophisticated algorithmic methods
85 should be considered for the development of an accurate
86 and reliable prediction model.

87 In this work, we propose a novel framework for the
88 development of efficient and reliable deep learning models
89 which constitutes the main contribution. The novelty of our
90 proposed methodology is that it guarantees considerable
91 improvement in the deep learning model's forecasting
92 performance in terms of reliability and accuracy, regarding
93 any utilized time-series. More analytically, the proposed

methodology ensures the “*suitability*” of a time-series for
fitting a deep learning model by performing a proper
transformation, in order to satisfy the stationarity property
and therefore no autocorrelation in the model's errors. The
stationarity is secured by the application of Augmented
Dickey–Fuller test and transformations based on first dif-
ferences or returns. It is worth mentioning that an adequate
deep learning model trained with the transformed series
presents no autocorrelation in the errors and a big
improvement of the forecasting performance is expected,
compared with the same model trained with the original
non-transformed series. We conducted a detailed and
comprehensive experimental analysis on time-series from
three real-world application domains which incorporate
most of global research interest, that is, financial stock
market, energy sector and the novel cryptocurrency area.
All prediction models were evaluated on both problems of
forecasting time-series price (regression) and also for the
prediction of time-series directional movements (classifi-
cation). Furthermore, the reliability of the models' fore-
casts was evaluated by examining the existence of
autocorrelation of the errors using the autocorrelation
function plot and the Ljung–Box Q test.

The remainder of this paper is organized as follows:
Sect. 2 presents a brief review of state-of-the-art deep
learning-based models for time-series forecasting. Sec-
tion 3 presents a comprehensive description of the problem
of reliability introduced in deep learning models for time-
series forecasting. Section 4 presents our proposed frame-
work, providing special attention to its theoretical advan-
tages and benefits. Section 6 presents our experimental
methodology including the data preparation and prepro-
cessing as well as the detailed experimental analysis,
regarding the evaluation of proposed methodology. Sec-
tion 7 summarizes our findings and discusses the experi-
mental results. Finally, Sect. 8 presents the conclusions and
some future directions.

2 Related work

Time-series forecasting is generally considered as one of
the most challenging and significantly complex research
areas. The complexity of time-series' internal structure is
caused by the variety of factors which have a deep influ-
ence on the series and on the volatility of these factors
[5, 6, 30]. During the last years, the significant develop-
ments in computer science as well as rapid advances in
research lead to the exponential generation of temporal and
sequential data [13]. Therefore, time-series infiltrated
almost every task and assignment, requiring a human
cognitive process. Recently, Fawaz et al. [11] provided an
excellent review, presenting a comprehensive overview of

144 the application of deep learning approaches in various
145 time-series domains. More specifically, they presented in
146 detail the process of mining time-series data using deep
147 learning methods for discovering new insights, and how
148 those insights impact the process of decision making. In the
149 rest of our research, we focus on three time-series appli-
150 cation domains which incorporate most of global research
151 interest, that is, financial stock market, energy sector and
152 the novel cryptocurrency area.

153 Liu et al. [18] developed a CNN–LSTM framework for
154 modeling and analyzing stock markets' quantitative
155 selection and timing strategy. The convolutional-based
156 framework is used for determining the selection of the
157 quantitative stock strategy and subsequently, the LSTM-
158 based framework is utilized for performing the quantitative
159 timing strategy in order to improve the amounts of profits.
160 The stock time-series data used in their research range from
161 January 1, 2007, to December 31, 2017. Their experiments
162 demonstrated that their proposed CNN–LSTM framework
163 could be efficiently applied for defining a quantitative
164 strategy and achieving better profits than the classical
165 Benchmark index.

166 Fischer and Krauss [12] focused on developing an effi-
167 cient forecasting model based on deep learning techniques
168 and unveil the sources of stock profitability. More specifi-
169 cally, the utilized LSTM networks for predicting the
170 directional movements for Standard & Poor's 500
171 (S&P500) constituent stocks. The utilized data contained
172 prices of S&P500 constituents from Thomson Reuters from
173 December 1989 to September 2015. The experiments
174 reported that LSTM exhibited a Sharpe Ratio of 5.8 prior to
175 transaction costs and daily returns of 46% per day. Addi-
176 tionally, LSTM networks outperformed computationally
177 efficient classification methods such as deep neural net-
178 work, random forest and logistic regression. Finally, the
179 authors developed a rule-based decision support making
180 system which focused on selecting winning and losing
181 stocks, exploiting the LSTM predictions.

182 Zhao et al. [32] proposed a deep learning ensemble
183 forecasting model to address the problem of forecasting oil
184 prices. Their proposed model is based on Bootstrapping
185 aggregation (Bagging) ensemble strategy which exploits
186 the predictions of advanced deep learning base models,
187 called Stacked Denoising Auto-Encoders (SDAE). The
188 data utilized in their study contained monthly prices cov-
189 ering a period from January 1986 to May 2016, concerning
190 198 exogenous factors such as cost of crude oil imports,
191 refiner values of crude oil products, information on rigs and
192 development wells drilled, oil product consumption, crude
193 oil production as well as macroeconomic and financial
194 indicators. Their experimental analysis showed the fore-
195 casting superiority of the proposed model, which was sta-
196 tistically proved by three nonparametric tests.

197 Cen and Wang [8] aimed at forecasting the volatility
198 behaviors of crude oil prices for increasing the prediction
199 accuracy of oil market price. The authors considered a
200 methodology based on a transfer learning approach in order
201 to extend the size of training set. Their research contained
202 daily data of West Texas Intermediate covering a time
203 period from January 31, 2005, to December 5, 2016, and
204 daily data of Brent oil covering a range from January 31,
205 2006, to October 17, 2017, concerning oil factors such as
206 opening price, closing price, lowest price and highest price.
207 Their proposed methodology was evaluated by comparing
208 the performance of a classical LSTM model trained with
209 the initial data and with the data transfer approach. Their
210 experiments showed that their methodology improved the
211 performance of the LSTM model, and thus, the authors
212 stated that the prediction was able to catch most fluctua-
213 tions of crude oil prices.

214 Nakano et al. [21] considered improving the traditional
215 “buy-and-hold” strategy presenting a new methodology
216 which exploits the predictions of advanced machine
217 learning models on Bitcoin's high-frequency technical
218 trading. More specifically, they designed ANN-based
219 models to extract the useful trading signals from technical
220 indicators calculated from the time-series return data at
221 time intervals of 15 min. The utilized data in this research
222 concerned historical returns and technical indicators,
223 ranging from December 2017 to January 2018, during
224 which Bitcoin suffers from substantial volatility and a
225 significant number of negative returns. Their preliminary
226 experimental results reported that the utilization of various
227 technical indicators could prevent over-fitting and consid-
228 erably enhance trading performance.

229 Ji et al. [16] studied the prediction performance on
230 Bitcoin price of various deep learning forecasting models
231 such as deep neural networks, convolutional neural net-
232 works, LSTM networks, deep residual networks and their
233 combinations. In their study, they utilized Bitcoin data
234 from 2590 days (from November 29, 2011, to December
235 31, 2018), containing 29 features. The authors performed a
236 comprehensive experimental procedure, considering the
237 problems of predicting the next's day Bitcoin price, and
238 whether or not the next day price will increase or decrease.
239 Their numerical experiments demonstrated that deep neural
240 networks reported the best performance for price move-
241 ment, while the LSTM models exhibited the best perfor-
242 mance for forecasting Bitcoins' price, slightly
243 outperforming the rest prediction models.

244 Pintelas et al. [24, 25] performed a comprehensive
245 research, evaluating advanced deep learning models for
246 forecasting the prices and directional movements of major
247 cryptocurrencies. Furthermore, the authors conducted a
248 detailed discussion, concerning if deep learning models can
249 be trusted as reliable predictors and if the cryptocurrencies

prices follow a random walk process. Their experimental analysis presented that even the state-of-the-art deep learning models were unable to create reliable forecasting models. Moreover, the authors stated that a few hidden patterns in cryptocurrency prices may probably exist, although these prices seem to follow almost a random walk process.

Summarizing, most approaches proposed in the literature attempt to exploit deep learning techniques for extracting useful knowledge from time-series data, aiming at obtaining better performance compared to the already existing proposed models. In this research, we propose a different approach and introduce a novel framework for the development of efficient and reliable deep learning models. The novelty of the proposed methodology is that it guarantees the forecasting reliability of the model's predictions, independent of the time-series data and selected deep learning model. It is worth noticing that none of the mentioned research approaches considered to improve both the accuracy and reliability of a deep learning model by exploiting the information provided by the characteristics of the time-series as well as the error of prediction model.

3 Reliability evaluation on times-series forecasting

3.1 Significance of autocorrelation in model's forecasting reliability

Let y_1, y_2, \dots, y_n be the observations of a time-series. A nonlinear regression model of order m is defined by

$$y_t = f(x_t, \theta) + \epsilon_t, \quad (1)$$

where $x_t = (y_{t-1}, y_{t-2}, \dots, y_{t-m}) \in \mathbb{R}^m$ consists of m values of y_t , θ is the parameter vector and ϵ_t is the white noise residual. After the model structure has been defined, function $f(\cdot)$ can be determined by sophisticated machine learning or deep learning methods.

After a prediction model has been successfully fit, it is significant to evaluate and assess how well the model is able to capture patterns. The most commonly utilized metrics for evaluating the regression performance of a forecasting model are mean absolute error (MAE) and root mean square error (RMSE).

Nevertheless, both regression evaluation metrics help determine how close the predicted values are to the actual ones, they do not evaluate whether the model properly fits the time-series data, while the residuals are usually dedicated to evaluate this. In other words, provided that function f is appropriately estimated, the prediction model's residuals

$$\hat{\epsilon}_t = y_t - \hat{y}_t, \quad (2)$$

are identically distributed and asymptotically independent, where \hat{y}_t is the predicted value.

It is worth noticing that in case the assumption of no-autocorrelation in the residuals is violated in a forecasting model, implies that its predictions may be inefficient, since the model has not exploited all the available mined information during the training process. In other words, the dependence between the residuals indicates that the model has not properly fitted the time-series data and there exists significant information left over which should be taken into account.

Two significant tools for testing the existence autocorrelation of the residuals are the auto-correlation function (ACF) plot and the Ljung–Box Q test for residual autocorrelation [6]. More analytically, ACF is obtained from the linear correlation of each residual $\hat{\epsilon}_t$ to the others in different lags, $\hat{\epsilon}_{t-1}, \hat{\epsilon}_{t-2}, \dots$ and illustrates the intensity of the temporal autocorrelation, while Ljung–Box Q test is a “portmanteau” test which assesses the null hypothesis H_0 that “a series of residuals exhibits no autocorrelation for a fixed number of lags L , against the alternative H_1 that “some autocorrelation coefficient is nonzero.” More specifically, the Ljung–Box Q test statistic is defined by

$$Q = n(n+2) \sum_{k=1}^L \frac{\rho_k^2}{n-k}, \quad (3)$$

where ρ_k are autocorrelation coefficients at lag- k , defined by

$$\rho_k = \frac{\sum_{i=1}^{n-k} (y_i - \bar{y})(y_{i+k} - \bar{y})}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (4)$$

with $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Under H_0 the statistic Q asymptotically follows a $\chi_{(L)}^2$ distribution. The null hypothesis H_0 is rejected and state that the model exhibits autocorrelation if

$$Q > \chi_{1-\alpha, L}^2 \quad (5)$$

where the critical value of the Chi-square distribution is defined for significance level α , or critical level $p = 1 - \alpha$, known as p value.

3.2 Strict stationarity and weak stationarity

Time-series exhibit a variety of properties which appear so often that are called stylized facts which include autocorrelation, long memory, fractal and multi-fractal properties. The major drawback when dealing with prices or values (levels) of such series is, from the stochastic processes point of view, that they follow a random walk process. The autocorrelation coefficients ρ_k , with $k > 1$ are statistically

340 significant for a large number of lags L and the first-order
 341 autocorrelation coefficient ρ_1 is equal to one [27]. Such
 342 series are also named *unit root time-series* or *integrated of*
 343 *order one* and are denoted by $I(1)$.

344 Under these conditions, modeling the levels of such
 345 series is inefficient because the residuals of the models
 346 display autocorrelation setting the whole structure of sta-
 347 tistical significance under question. In order to study effi-
 348 ciently these series, they have to be stationary, which is a
 349 highly significant for ensuring the development of a reli-
 350 able prediction model.

351 Suppose that $F_y(y_{t_1+\tau}, \dots, y_{t_n+\tau})$ is the cumulative dis-
 352 tribution function of the unconditional joint distribution of
 353 $\{y_t\}$ at times $t_1+\tau, \dots, t_n+\tau$, then the stochastic process $\{y_t\}$
 354 is strictly stationary if

$$F_y(y_{t_1+\tau}, \dots, y_{t_n+\tau}) = F_y(y_{t_1}, \dots, y_{t_n}), \quad (6)$$

356 for all $\tau, t_1, \dots, t_n \in \mathbb{R}$ and $n \in \mathbb{N}$. Nevertheless, in time-
 357 series the strong form of stationarity is relaxed leading to
 358 the *weak-stationarity* or *covariance stationarity* [6].
 359 Therefore, a stochastic process is covariance stationary if
 360 the mean is constant, the second moment is finite, and the
 361 covariance function depends only on the difference
 362 between t_1 and t_2 and needs to be indexed by only one
 363 variable, i.e.,

$$cov_{yy}(t_1, t_2) = cov_{yy}(t_1 - t_2, 0). \quad (7)$$

365 where cov_{yy} is the auto-covariance of series y_t . Summa-
 366 rizing, stationarity means that the statistical properties of a
 367 stochastic process which generates a time-series are con-
 368 stant over time. Stationary processes are easier to analyze,
 369 model, and investigate, and it has been a common
 370 assumption of many practices involving statistical infer-
 371 ence, modeling and forecasting.

372 Having identified the problem, a solution for stationarity
 373 comes from the partial autocorrelation function, where the
 374 lag- k coefficient $\phi_{k,k}$ is given by the following formula

$$\begin{cases} \phi_{k,k} &= \frac{\rho_k - \sum_{j=1}^{k-1} \phi_{k-1,j} \rho_{k-j}}{1 - \sum_{j=1}^{k-1} \phi_{k-1,j} \rho_{k-j}}, \\ \phi_{k,j} &= \phi_{k-1,j} - \phi_{k,k} \phi_{k-1,k-j}, \end{cases} \quad (8)$$

376 for $k > 1$ and $\phi_{1,1} = \rho_1$ [26]. Clearly, in case the series
 377 exhibits a unit root, that is $\rho_1 = 1$, it immediate follows
 378 that the first-order partial autocorrelation coefficient $\phi_{1,1}$
 379 will be one. The significance of the partial autocorrelation
 380 function is that if only the first coefficient is statistically
 381 significant and the rest are not, which is the usual case in
 382 most time-series of scientific interest, then this is a guide
 383 that the initial series should be differenced by using the first
 384 differences of the series, namely

$$\Delta_t = y_t - y_{t-1}. \quad (9)$$

386 Therefore, taking the first difference of the levels of the
 387 series results in stationarity and these series are named
 388 *integrated of order zero* and are denoted by $I(0)$.

389 However, when dealing with time-series there might be
 390 an overlapping of a variety of non-stationarities, including
 391 unit-roots, structural breaks, level shifts, seasonal cycles, or
 392 a changing variance. Notice that the typical transformation
 393 when the series is $I(1)$ (non-stationary) is to take the first
 394 differences of the series and transform it to a series $I(0)$
 395 (stationary), while if the series incorporate structural breaks
 396 or a changing variance, i.e., due to crises, a nonlinear Box-
 397 Cox transformation [22] is the appropriate available option.
 398 A Box-Cox transformation is a way to transform non-
 399 normal dependent variables into a normal shape, since
 400 normality is a critical assumption for many statistical
 401 techniques. The one-parameter Box-Cox transformation is
 402 defined as

$$y_t = \begin{cases} \frac{y_t^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0; \\ \ln y_t, & \text{if } \lambda = 0. \end{cases} \quad (10)$$

404 where common nonzero Box-Cox transformations are for
 405 $\lambda = -3, -2, -0.5, 0, 0.5, 1$ and 2 . The large majority of
 406 time-series follow the rule $\lambda = 0$; therefore, the stationarity
 407 of these series is achieved via the returns that is the first
 408 logarithmic differences,

$$r_t = \ln y_t - \ln y_{t-1} \approx \frac{y_t - y_{t-1}}{y_{t-1}}, \quad (11)$$

410 the last expression being the percentage change or *returns*.

4 Research methodology and proposed framework

413 In this section, we introduce our proposed methodology for
 414 considerably improving the performance of a deep learning
 415 model for time-series forecasting in terms of accuracy and
 416 reliability. based on the well-established econometric theo-
 417 ry and time-series analysis with respect to stationarity and
 418 non-stationarity properties.

420 Revisiting the problem, when applying a machine
 421 learning or a deep learning model to time-series for fore-
 422 casting, the levels of the series are not-stationary, meaning
 423 that they possess unit roots and some order of integration¹.

¹ Non-stationary time-series which can be transformed in this way are called series integrated of order d . Usually, the order of integration is either $I(0)$ or $I(1)$; it's extremely rare to see values for d that are 2 or more in real-world applications [7]. Additionally, all series in this research are $I(1)$.

Author Proof

424 Notice that the identification of a unit root in a time-series
425 can be easily performed via the Augmented Dickey–Fuller
426 (ADF) test [6, 23]. The testing procedure is applied to the
427 model,

$$\Delta y_t = \alpha + \beta t + \gamma y_{t-1} + \sum_{i=1}^{k-1} \delta_i \Delta y_{t-i} + \epsilon_t, \quad (12)$$

429 where α is a constant, β is the coefficient of time trend, and
430 $\gamma = (\rho_1 - 1)$ where ρ_1 denotes the first-order autocorrela-
431 tion coefficient. It is worth mentioning that k is the lag
432 order of the autoregressive process chosen so that no serial
433 correlation exists in the residuals ϵ_t , ensuring that the test is
434 efficient and reliable. If $\alpha = 0$ and $\beta = 0$, then we have a
435 random walk stochastic process, while if $\alpha \neq 0$ and $\beta = 0$,
436 we have a random walk stochastic process with drift. The
437 unit root test is carried out testing the statistical signifi-
438 cance under the null hypothesis $H_0 : \{\gamma = 0 \text{ that is } \rho = 1\}$
439 against the alternative hypothesis $H_1 : \{\gamma < 0 \text{ that is } \rho < 1\}$.

440 The solution depending on the nature of the time-series
441 is to iteratively take the first differences (9) or the returns
442 (11) until the series is made stationary, implying that the
443 first-order autocorrelation coefficient ρ_1 is less than one. It
444 is worth noticing that the series transformation based on
445 first differences or returns implies that the autocorrelation
446 in the residuals of the model is removed. This indicates that
447 the prediction model is able to explain the data much
448 better, since it captures all possible nonlinearities, ensuring
449 the efficiency and effectiveness of the model.

450 Table 1 presents the pseudo-code of our proposed
451 framework. Initially, the time-series data are imported
452 (Step 1). Then, the ADF test is applied to examine whether
453 the levels of the series are non-stationary, meaning that
454 they possess a unit root (Step 2). In case the series is non-
455 stationary the transformation based on first differences or
456 returns is iteratively applied on the training data until the
457 new transformed series is stationary (Steps 4–7). Subse-
458 quently, the new transformed time-series data are used for
459 training the prediction model (Step 8).

460 In contrast, in case the series is stationary, then the
461 levels of time-series are used for training the prediction
462 model (Step 10). Subsequently, the prediction model's
463 errors on the training data are used for further examination
464 and testing. Notice that a model trained with a series which
465 does not possess a unit root and has not been differenced,
466 may exhibit autocorrelation in the training errors. In other
467 words, although for any reasonable model the predicted
468 values will be close to the real values, the existence of large
469 autocorrelation coefficients characterizes the model as
470 inefficient [28]. Therefore, the residuals of the training data
471 are examined for autocorrelation by simply performing
472 ACF plots and/or Ljung–Box Q test (Step 11). In case the
473 residuals possess autocorrelation, the proposed

474 transformation is applied on the training data and the model
475 is re-trained using the new transformed series (Steps
476 13–14). Notice that if the levels of the series are stationary
477 and the residuals on the training set exhibit no autocorre-
478 lation, then there is no need to transform the series, since
479 this will lead to the dangerous phenomenon of over-dif-
480 ferencing the series. In other words, over-differencing leads
481 the whole process to be “non-invertible” and lacks an
482 infinite-order autoregressive representation. Figure 1 pre-
483 sents an overview of the proposed architecture in the form
484 of a flowchart.

485 Finally, it is worth mentioning that in case the model is
486 trained with a transformed series based on first differences
487 or returns, the reverse transformation is used in the pre-
488 dictions of the model to obtain the prediction for the levels
489 of the original time-series.

5 Data

491 In our research, we utilized three benchmark datasets from
492 the popular real-world application domains: finance, com-
493 modity and cryptocurrency, in order to demonstrate the
494 efficiency of our proposed methodology.

495 From finance domain, we utilized data from January 1,
496 2013, to December 31, 2019, of Standard & Poor's 500
497 (S&P500) prices in USD from <http://finance.yahoo.com>
498 Web site. The data were divided into training set consisting
499 of daily Brent prices from January 1, 2015, to December
500 31, 2018 (4 years), and a testing set consisting of daily
501 prices from January 1, 2018, to December 31, 2019 (1
502 year).

503 From commodity domain, we utilized daily data from
504 January 1, 2013, to December 31, 2019, of Brent prices in
505 USD from <https://www.eia.gov/> Web site. The data were
506 divided into training set consisting of daily Brent prices
507 from January 1, 2015, to December 31, 2018 (4 years), and
508 a testing set consisting of daily prices from January 1,
509 2018, to December 31, 2019 (1 year).

510 From cryptocurrency domain, we utilized daily data
511 from January 1, 2015, to December 31, 2019, of Bitcoin
512 (BTC) cryptocurrency in USD from <https://coinmarketcap.com>
513 Web site. The data were divided into training set
514 consisting of daily BTC prices from January 1, 2015, to
515 June 30, 2018 (3.5 years), and a testing set consisting of
516 daily prices from July 1, 2018, to December 31, 2019 (1/2
517 year).

518 All time-series data contained no missing values, while
519 the outlier prices were not removed in order not to destroy
520 the dynamics of each series, even if these prices are the
521 result of exceptional events.

522 Table 2 summarizes the descriptive statistics for the
523 training and testing set of each dataset, including the

Table 1 Proposed framework to enhance deep learning in time-series forecasting

```

Step 1. Import time-series data.
Step 2. Apply the ADF unit root test.
Step 3. If (Time-series possess a unit root) then
Step 4.   repeat
Step 5.     Apply the transformation based on first differences (9) or returns (11).
Step 6.     Apply the ADF unit root test.
Step 7.   until (Time-series is stationary)
Step 8.   Train the forecasting model using the transformed time-series.
Step 9. else
Step 10.  Train the forecasting model using the original time-series.
Step 11.  Calculate the residuals on the training set.
Step 12.  If (residuals possess autocorrelation) then
Step 13.   Apply the transformation based on first differences (9) or returns (11).
Step 14.   Re-train the forecasting model using the transformed time-series.
Step 15.  end if
Step 16.end if
    
```

/ Time-series is not I(0) */*
/ non-stationary */*

/ Time-series is I(0) */*
/ stationary */*

/ Training process completed */*

524 measures: Minimum, Maximum, Mean, Standard Devia- 553
 525 tion (Std. Dev.), Median, Skewness and Kurtosis, for pre- 554
 526 senting the nature of the distribution. Additionally, Table 3 555
 527 presents the increase and decrease cases and the corre- 556
 528 sponding percentages in S&P500, Brent and BTC datasets. 557

529 **6 Experimental analysis**

530 In this section, we apply our proposed methodology to the 561
 531 S&P500, Brent and BTC time-series to identify whether or 562
 532 not the training data are stationary, utilizing the ADF unit 563
 533 root test. 564

534 Table 4 presents the results of the ADF unit root test for 565
 535 the training data of all series under consideration, i.e., 566
 536 S&P500, Brent and BTC, performed on the level of the 567
 537 original series. By taking into consideration, the *t*-statistics 568
 538 (*t*-stat.) and the associated *p* values, we conclude that the 569
 539 null hypothesis H_0 : “*the levels possess a unit root and are* 570
 540 *non-stationary*” is accepted for S&P500, Brent and BTC 571
 541 series. 572

542 In the sequel, we perform the ADF test to the trans- 573
 543 formed time-series based on both first differences (first 574
 544 differenced series) and returns (returns series), to examine 575
 545 if the unit root has been removed, according to our pro- 576
 546 posed framework. 577

547 Table 5 presents the results of the ADF unit root test for 578
 548 the training data of all transformed time-series. Notice that 579
 549 (*) denotes statistical significance at the 5% critical level. 580
 550 Clearly, performing either the first differences or the 581
 551 returns transformation clearly solves the unit root problem, 582
 552 since all *p* values are practically zero and therefore the null

hypothesis H_0 is rejected and all transformed series are 553
 indeed stationary. 554

Thus, both transformed series are “suitable” for fitting 555
 a deep learning model which will present no autocorrela- 556
 tion in the errors, and a big improvement of the forecasting 557
 performance is expected, compared with the same model 558
 trained with the original non-transformed series. 559

560 **6.1 Numerical experiments**

561 In the sequel, we present a comprehensive experimental 562
 563 analysis, to evaluate the efficiency and reliability of our 564
 565 proposed methodology. More specifically, we compare the 566
 567 performance of two efficient deep learning forecasting 568
 569 models trained with the levels of the time-series (Time- 570
 571 series) and with the two transformed series based on first 572
 573 differences (first differenced series) and returns (returns 574
 575 series). 576

577 Under exhaustive experimentation (utilizing different 578
 579 number of the CNN and LSTM layers, different number of 580
 581 units in the LSTM layers, different number of filter in CNN 581
 582 layers), the selected models were an LSTM model which 582
 consist of a LSTM layer of 50 units and an output layer of 583
 one neuron (Fig. 2) and a CNN–LSTM which consists of 584
 two convolutional layers of 16 and 32 filters of size (2,) 585
 with the same padding, followed by a LSTM layer of 50 586
 units and an output layer of one neuron (Fig. 3). Notice that 587
 both models were trained with Adaptive Moment Estima- 588
 tion [17] with a batch size equal to 128, using a mean- 589
 squared loss function. The implementation code was 590
 written in Python 3.4 using Keras library [15] on a laptop 591
 (Intel(R) Core(TM) i7-6700HQ CPU 2.6GHz and 16GB 592

Author Proof

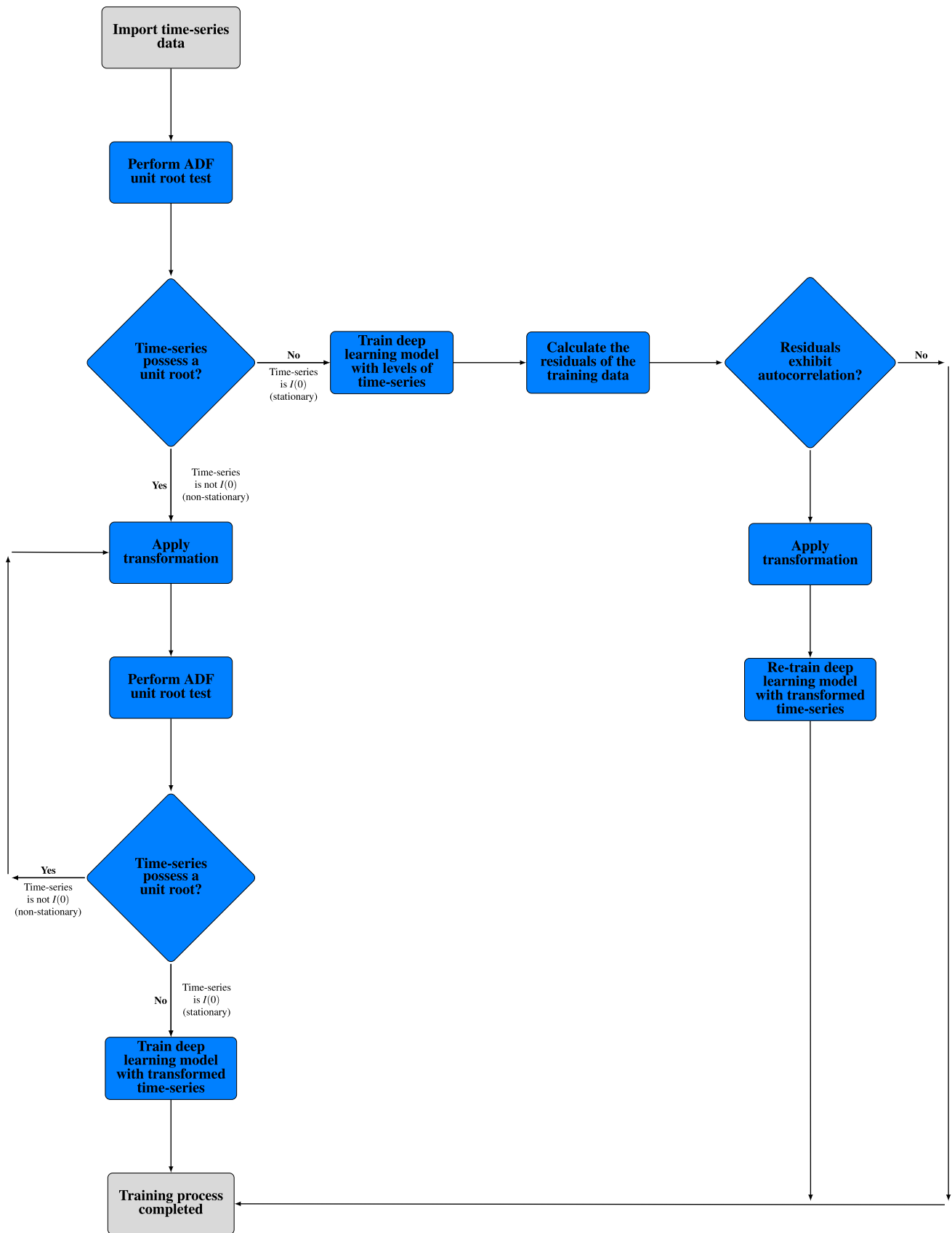


Fig. 1 Flowchart

RAM) running Windows 10.0 operating system. Each forecasting model was trained with the traditional time-series, the first differenced series and returns series, utilizing four different values of window size m , i.e., $m = 4, 6, 9$ and 12 . Finally, in order to reduce exponential trend and homogenize the variability and stability of the patterns, the traditional time-series data were transformed utilizing a natural logarithm (\ln) and the model's predicted value is used to predict the price on the following day. Moreover, we recall that in case the models were trained using the first difference or returns transformation to series, the reverse transformation is utilized and to predict the price on the following day.

The regression performance was evaluated utilizing the metrics: mean absolute error (MAE) and root mean square error (RMSE), while for the binary classification problem of predicting whether the price would increase or decrease on the following day, four performance metrics were used: Accuracy (Acc), F_1 -score (F_1), Sensitivity (Sen), Specificity (Spe), Positive Predicted Values (PPV) and Negative Predictive Values (NPV) which are defined by

$$Acc = \frac{TP+TN}{TP+FP+FN+FP}, \tag{13}$$

$$F_1 = \frac{2TP}{2TP + FP + FN}, \tag{14}$$

$$Spe = \frac{TP}{TP+FN}, \tag{15}$$

$$Spe = \frac{TN}{TN+FP}, \tag{16}$$

$$PPV = \frac{TP}{TP+FP}, \tag{17}$$

$$NPV = \frac{TN}{TN+FN} \tag{18}$$

where TP stands for the number of prices which were correctly identified to be increased, TN stands for the number of prices which were correctly identified to have a decreased, FP (type I error) stands for the number of prices which were misidentified to be increased, and FN (type II

Table 3 The number of up and down movements of S&P 500, Brent and BTC datasets

| Data | Decrease | % | Increase | % |
|-------------------|----------|-------|----------|-------|
| S&P500 | | | | |
| Training set | 815 | 54.01 | 694 | 45.99 |
| Testing set | 149 | 59.36 | 102 | 40.64 |
| Brent | | | | |
| Training set | 738 | 47.19 | 826 | 52.81 |
| Testing set | 140 | 53.64 | 121 | 46.36 |
| BTC | | | | |
| Training set | 900 | 55.87 | 711 | 44.13 |
| Testing set | 102 | 47.44 | 113 | 52.56 |

Table 4 ADF unit root test of all series under consideration

| Series | S&P500 | Brent | BTC |
|--------------------|----------|----------|----------|
| Time-series | | | |
| t stat. | - 2.8469 | - 1.1599 | - 2.7497 |
| p value | 0.1806 | 0.9170 | 0.2166 |

error) stands for the number of prices which misidentified to be decreased.

Moreover, we included area under curve (AUC) metric in our analysis which constitutes one of the most significant classification metrics and it is presented using the receiver

Table 5 ADF unit root test of all transformed times-series based on first differences and returns

| Series | S&P500 | Brent | BTC |
|---------------------------------|----------|----------|----------|
| First differenced series | | | |
| t stat. | - 39.146 | - 38.516 | - 7.3402 |
| p value | 0.0000* | 0.0000* | 0.0000* |
| Returns series | | | |
| t stat. | - 39.663 | - 38.259 | - 39.893 |
| p value | 0.0000* | 0.0000* | 0.0000* |

Table 2 Descriptive statistics for S&P500, Brent and BTC prices

| Data | Minimum | Maximum | Mean | SD | Median | Skewness | Kurtosis |
|-------------------|---------|----------|---------|---------|---------|----------|----------|
| S&P500 | | | | | | | |
| Training set | 1457.15 | 2930.75 | 2153.78 | 367.52 | 2087.90 | 0.31 | - 0.75 |
| Testing set | 2447.89 | 3240.02 | 2912.09 | 149.32 | 2918.65 | - 0.26 | 0.17 |
| Brent | | | | | | | |
| Training set | 26.01 | 118.90 | 71.46 | 25.84 | 63.27 | 0.43 | - 1.28 |
| Testing set | 50.57 | 74.94 | 64.31 | 4.44 | 63.99 | 0.16 | - 0.30 |
| BTC | | | | | | | |
| Training set | 178.10 | 19497.40 | 3261.48 | 3675.02 | 1152.36 | 1.40 | 1.84 |
| Testing set | 6640.52 | 13016.23 | 9221.26 | 1466.76 | 9244.97 | 0.27 | - 0.93 |

Author Proof

625 operating characteristic (ROC) curve. Notice that ROC
626 curve is created by plotting the true positive rate (Sensi-
627 tivity) against the false positive rate (Specificity) at various
628 threshold settings.

629 6.1.1 S&P500

630 Tables 6 and 7 present the performance comparison of
631 both LSTM and CNN–LSTM forecasting models, respec-
632 tively, for S&P500 dataset. The LSTM model improved its
633 average performance, in terms of MAE and RMSE scores
634 by 30.57%–43.88% and 18.73%–37.31%, respectively,
635 when trained with the first differenced series, while the
636 CNN–LSTM model considerably improved its MAE
637 average performance by 11.35%–45.34% and its RMSE
638 average performance by 4.82%–34.59%, in the same situ-
639 ation. Furthermore, the LSTM model reduced its average
640 MAE and RMSE scores by 23.76%–49.40% and 11.51%–
641 45.45% in case it was trained with the returns of S&P500
642 prices, while the CNN–LSTM model reduced its average
643 MAE performance by 21.10%–48.98% and its RMSE
644 average performance by 11.89%–38.96%. Summarizing,
645 we conclude that the regression performance of both
646 LSTM and CNN–LSTM forecasting models was consid-
647 erably improved, utilizing the first differenced and returns
648 series, instead of the traditional S&P500 time-series.

649 Furthermore, the classification performance of both
650 prediction models was also improved utilizing our pro-
651 posed methodology. More specifically, both LSTM and
652 CNN–LSTM models were biased in case they were trained
653 with the traditional time-series. In contrast, the trade-off
654 between sensitivity and specificity as well as between

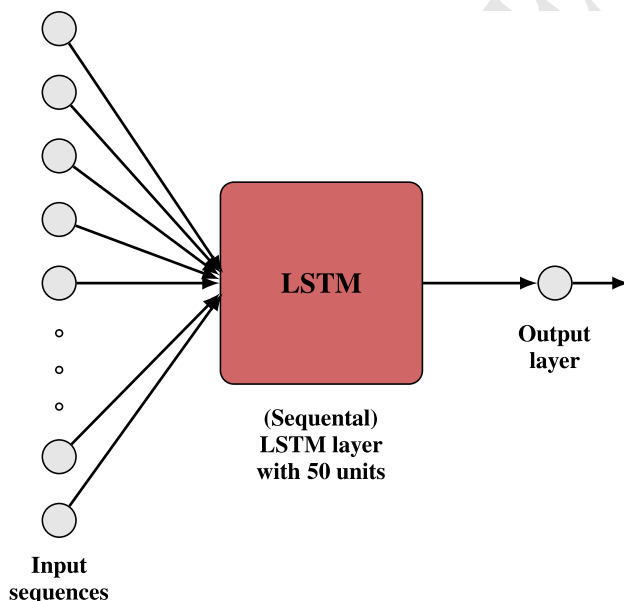


Fig. 2 LSTM forecasting model architecture

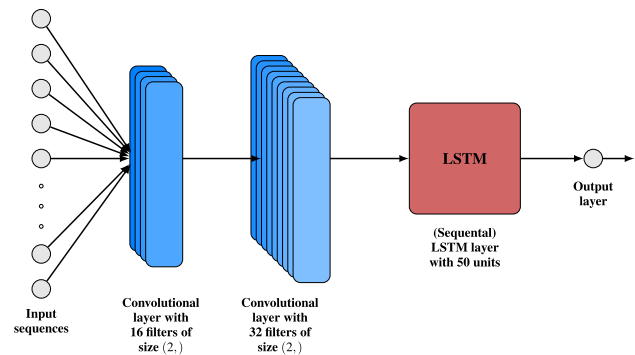


Fig. 3 CNN–LSTM forecasting model architecture

655 positive and negative predictive values of both models was
656 considerably increased in case the models were trained
657 with the first differenced.

658 It is worth noticing that both LSTM and CNN–LSTM
659 models exhibited the highest classification performance in
660 case they were trained with the first differenced series and
661 the best regression performance when they were trained
662 with the returns series. Moreover, the LSTM model trained
663 with the first differenced time-series reported the best
664 classification performance for all values of window size
665 m and the best regression performance for $m = 6, 9$ and 12.
666 The CNN–LSTM model trained with the first differenced
667 and returns series reported the best performance relative to
668 classification and regression accuracy, respectively.

669 6.1.2 Brent

670 Tables 8 and 9 present the performance comparison for
671 Brent forecasting problem of LSTM and CNN–LSTM,
672 respectively. The LSTM and CNN–LSTM models
673 improved their average MAE score by 6.44–29.66% and
674 5.88–17.52%, respectively, in case they were trained with
675 the first differenced series, instead of the traditional series.
676 Furthermore, their average RMSE score was reduced by
677 5.73–25.46% and 2.88–10.33% in the same situation.
678 However, the regression performance of both forecasting
679 models worsens, in case they were trained using the returns
680 series.

681 Regarding the classification performance, both LSTM
682 and CNN–LSTM models were biased in case they were
683 trained with the traditional series. On the other hand,
684 LSTM considerably improved its classification perfor-
685 mance utilizing either first differenced or returns series as
686 training data in terms of trade-off between sensitivity and
687 specificity as well as the trade-off between positive and
688 negative predicted values. Additionally, the CNN–LSTM
689 significantly improved its classification performance using
690 first differenced series as training data, in terms of both
691 accuracy and trade-off between sensitivity and specificity.

692 It is also worth mentioning that both LSTM and CNN–
693 LSTM forecasting models improved their F_1 -score, in case
694 they were trained with the transformed series.

695 The performance of the forecasting models, relative to
696 the value of the window size m , both models improved
697 their performance as the value of window size increases.
698 Moreover, it is worth mentioning that the best overall
699 performance was reported by CNN–LSTM trained with the
700 first differenced series with $m = 12$.

701 6.1.3 Bitcoin

702 Tables 10 and 11 present the performance comparison of
703 LSTM and CNN–LSTM, respectively, relative to BTC
704 dataset. Similar conclusions can be drawn with the previ-
705 ous benchmarks. More specifically, the LSTM model
706 improved its MAE and RMSE average performance by
707 24.29–31.90% and 19.44–23.38%, respectively, when
708 trained with the first differenced series instead of the tra-
709 ditional series, while the CNN–LSTM model improved its
710 MAE average performance by 14.34–38.12% and its
711 RMSE average performance by 6.0–29.1%, in the same
712 situation. Moreover, the MAE and RMSE average perfor-
713 mance of the LSTM model was reduced by 21.14–36.59%
714 and 15.99–27.49% in case it was trained with the returns of
715 Bitcoin's prices, while the CNN–LSTM model improved
716 its MAE average performance by 8.57–39.57% and its
717 RMSE average performance by 2.56–29.52%, in the same
718 situation. Summarizing, we can easily conclude that the
719 regression performance of both forecasting models was
720 considerably improved, utilizing our proposed methodol-
721 ogy for data preparation.

722 Regarding the classification performance, our proposed
723 methodology increased the accuracy of both prediction

724 models. More analytically, the interpretation of Tables 10
725 and 11 reveals that LSTM and CNN–LSTM models were
726 biased when trained with the traditional time-series. In
727 contrast, the trade-off between sensitivity and specificity as
728 well as the trade-off between positive and negative pre-
729 dictive values of both forecasting models was considerably
730 improved, in case they were trained with the first differ-
731 enced series or the returns series. Finally, AUC and F_1 -
732 score of both models were improved in case they were
733 trained with the transformed series instead of the traditional
734 time-series.

735 The LSTM model trained with the first differenced
736 series exhibited the best classification performance for all
737 values of window size m and the best regression perfor-
738 mance for $m = 6, 9$ and 12. Moreover, the CNN–LSTM
739 model exhibited the lowest (best) MAE and RMSE scores,
740 in case it was trained with the returns series, while it pre-
741 sented the best overall classification performance, in case it
742 was trained with the first differenced series. Finally, it is
743 worth mentioning that both LSTM and CNN–LSTM
744 models exhibited the best regression performance when
745 they were trained with the returns series and the highest
746 classification performance, in case they were trained with
747 the first differenced series.

748 6.2 Reliability evaluation of the forecasts

749 In the sequel, we evaluate the reliability of all forecasting
750 models by examining the existence of autocorrelation in
751 the residuals utilizing the Auto-Correlation Function (ACF)
752 plot and the Ljung–Box Q test for residual autocorrelation
753 [6]. In other words, we examine whether each trained
754 model has properly fitted the time-series by examining
755 whether the residuals are identically distributed and

Table 6 Performance comparison of the LSTM model for S&P500 dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F_1 | Sen | Spe | PPV | NPV |
|--------------------------|---------|-------|-------|---------|-------|-------|-------|-------|-------|-------|
| Time-series | 4 | 33.31 | 42.80 | 50.42 | 0.532 | 0.606 | 0.595 | 0.470 | 0.651 | 0.405 |
| First differenced series | | 18.69 | 26.83 | 47.87 | 0.467 | 0.546 | 0.530 | 0.404 | 0.565 | 0.410 |
| Returns series | | 16.85 | 23.35 | 46.59 | 0.435 | 0.570 | 0.597 | 0.274 | 0.546 | 0.347 |
| Time-series | 6 | 34.49 | 42.35 | 47.79 | 0.514 | 0.425 | 0.324 | 0.703 | 0.655 | 0.415 |
| First differenced series | | 20.10 | 29.38 | 48.11 | 0.466 | 0.554 | 0.545 | 0.388 | 0.566 | 0.417 |
| Returns series | | 19.31 | 31.56 | 49.00 | 0.460 | 0.590 | 0.617 | 0.304 | 0.565 | 0.352 |
| Time-series | 9 | 29.84 | 39.02 | 47.39 | 0.520 | 0.385 | 0.277 | 0.762 | 0.631 | 0.409 |
| First differenced series | | 20.58 | 29.27 | 50.36 | 0.497 | 0.559 | 0.531 | 0.463 | 0.593 | 0.451 |
| Returns series | | 18.94 | 29.40 | 50.07 | 0.486 | 0.574 | 0.565 | 0.406 | 0.583 | 0.409 |
| Time-series | 12 | 27.14 | 38.04 | 47.59 | 0.509 | 0.471 | 0.385 | 0.634 | 0.606 | 0.413 |
| First differenced series | | 18.84 | 25.43 | 49.40 | 0.489 | 0.547 | 0.514 | 0.465 | 0.585 | 0.435 |
| Returns series | | 20.69 | 30.12 | 47.12 | 0.466 | 0.524 | 0.493 | 0.439 | 0.562 | 0.402 |

Table 7 Performance comparison of the CNN–LSTM model for S&P500 dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F ₁ | Sen | Spe | PPV | NPV |
|--------------------------|---------|-------|-------|---------|-------|----------------|-------|-------|-------|-------|
| Time-series | 4 | 33.40 | 38.70 | 48.59 | 0.501 | 0.395 | 0.421 | 0.581 | 0.601 | 0.272 |
| First differenced series | | 18.26 | 25.31 | 50.20 | 0.478 | 0.551 | 0.529 | 0.426 | 0.575 | 0.411 |
| Returns series | | 17.04 | 23.62 | 46.02 | 0.426 | 0.568 | 0.601 | 0.251 | 0.540 | 0.330 |
| Time-series | 6 | 23.40 | 29.67 | 51.67 | 0.535 | 0.494 | 0.437 | 0.634 | 0.649 | 0.437 |
| First differenced series | | 20.75 | 28.24 | 51.41 | 0.516 | 0.554 | 0.507 | 0.525 | 0.610 | 0.451 |
| Returns series | | 18.47 | 26.14 | 48.35 | 0.466 | 0.557 | 0.550 | 0.383 | 0.566 | 0.398 |
| Time-series | 9 | 35.46 | 42.18 | 46.45 | 0.509 | 0.356 | 0.275 | 0.743 | 0.623 | 0.412 |
| First differenced series | | 22.78 | 31.94 | 53.41 | 0.458 | 0.513 | 0.480 | 0.436 | 0.554 | 0.403 |
| Returns series | | 20.61 | 29.82 | 52.13 | 0.508 | 0.561 | 0.527 | 0.488 | 0.601 | 0.414 |
| Time-series | 12 | 33.96 | 40.96 | 45.38 | 0.489 | 0.389 | 0.302 | 0.677 | 0.577 | 0.398 |
| First differenced series | | 23.97 | 33.03 | 52.08 | 0.459 | 0.519 | 0.487 | 0.465 | 0.557 | 0.395 |
| Returns series | | 23.97 | 33.03 | 51.20 | 0.459 | 0.519 | 0.487 | 0.432 | 0.557 | 0.391 |

Table 8 Performance comparison of the LSTM model for Brent dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F ₁ | Sen | Spe | PPV | NPV |
|--------------------------|---------|------|------|---------|-------|----------------|-------|-------|-------|-------|
| Time-series | 4 | 1.40 | 1.81 | 53.05 | 0.543 | 0.442 | 0.370 | 0.717 | 0.618 | 0.496 |
| First differenced series | | 0.99 | 1.35 | 48.42 | 0.485 | 0.499 | 0.479 | 0.490 | 0.522 | 0.467 |
| Returns series | | 1.10 | 1.92 | 51.24 | 0.506 | 0.506 | 0.478 | 0.533 | 0.541 | 0.480 |
| Time-series | 6 | 1.33 | 1.75 | 53.75 | 0.545 | 0.497 | 0.439 | 0.652 | 0.608 | 0.501 |
| First differenced series | | 1.25 | 1.86 | 50.19 | 0.504 | 0.503 | 0.472 | 0.537 | 0.541 | 0.497 |
| Returns series | | 1.31 | 2.18 | 50.27 | 0.504 | 0.509 | 0.484 | 0.525 | 0.541 | 0.488 |
| Time-series | 9 | 1.62 | 2.03 | 52.36 | 0.530 | 0.462 | 0.449 | 0.610 | 0.584 | 0.508 |
| First differenced series | | 1.37 | 1.86 | 51.43 | 0.514 | 0.534 | 0.519 | 0.508 | 0.550 | 0.498 |
| Returns series | | 1.21 | 1.84 | 50.19 | 0.500 | 0.528 | 0.525 | 0.475 | 0.535 | 0.486 |
| Time-series | 12 | 1.68 | 2.16 | 53.67 | 0.536 | 0.536 | 0.544 | 0.528 | 0.588 | 0.494 |
| First differenced series | | 1.45 | 1.94 | 51.66 | 0.515 | 0.543 | 0.537 | 0.493 | 0.551 | 0.499 |
| Returns series | | 1.28 | 2.03 | 51.66 | 0.513 | 0.552 | 0.557 | 0.470 | 0.550 | 0.498 |

756 asymptotically independent. We recall that the Ljung–
 757 Box Q test is a “portmanteau” test which assesses the null
 758 hypothesis H_0 that “a series of residuals exhibits no
 759 autocorrelation for a fixed number of lags L ,” against the
 760 alternative hypothesis H_1 that “some autocorrelation
 761 coefficient is nonzero.”

762 Tables 12 and 13 present the information of the sta-
 763 tistical analysis performed by Ljung–Box Q test for $L = 10$
 764 of LSTM and CNN–LSTM, respectively. Clearly, the null
 765 hypothesis H_0 of no autocorrelation in the residuals is
 766 accepted, in case the models were trained with the first
 767 differenced or returns series, relative to all benchmarks and
 768 window sizes. On the other hand, both prediction models
 769 reject the H_0 in case they were trained with the traditional
 770 time-series.

For completeness, we also present the ACF plots of
 LSTM and CNN–LSTM for S&P500, Brent and BTC
 datasets in order to illustrate the intensity of the temporal
 autocorrelation. In each ACF plot, the confident limits are
 denoted with blue dashed lines and are constructed
 assuming that the residuals follow a Gaussian probability
 distribution. Notice that the ACF plot of each model for
 S&P500, Brent and BTC datasets was calculated for
 $m = 9$, $m = 9$ and $m = 6$, respectively, for which the
 models exhibited the best performance.

Figures 4, 5 and 6 present the ACF plots of LSTM
 model for S&P500, Brent and BTC datasets, respectively.
 The ACF plots of the forecasting model trained with the
 traditional time-series violate the assumption of no auto-
 correlation in the residuals. More specifically, the signifi-
 cant spikes that occurred in several lags suggest the

Table 9 Performance comparison of the CNN–LSTM model for Brent dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F ₁ | Sen | Spe | PPV | NPV |
|--------------------------|---------|------|------|---------|-------|----------------|-------|-------|-------|-------|
| Time-series | 4 | 1.13 | 1.51 | 53.01 | 0.535 | 0.494 | 0.464 | 0.607 | 0.583 | 0.500 |
| First differenced series | | 0.98 | 1.35 | 50.77 | 0.508 | 0.522 | 0.502 | 0.514 | 0.545 | 0.491 |
| Returns series | | 1.17 | 1.61 | 49.10 | 0.468 | 0.507 | 0.508 | 0.428 | 0.508 | 0.457 |
| Time-series | 6 | 1.35 | 1.74 | 53.17 | 0.548 | 0.421 | 0.330 | 0.765 | 0.632 | 0.496 |
| First differenced series | | 1.17 | 1.64 | 51.47 | 0.521 | 0.542 | 0.509 | 0.521 | 0.550 | 0.505 |
| Returns series | | 1.67 | 2.39 | 52.16 | 0.510 | 0.546 | 0.490 | 0.558 | 0.564 | 0.495 |
| Time-series | 9 | 1.54 | 1.94 | 52.97 | 0.544 | 0.414 | 0.353 | 0.734 | 0.612 | 0.500 |
| First differenced series | | 1.27 | 1.80 | 53.74 | 0.548 | 0.541 | 0.521 | 0.535 | 0.565 | 0.511 |
| Returns series | | 2.07 | 2.88 | 51.35 | 0.506 | 0.549 | 0.511 | 0.517 | 0.550 | 0.497 |
| Time-series | 12 | 1.44 | 1.90 | 53.82 | 0.537 | 0.551 | 0.558 | 0.515 | 0.587 | 0.500 |
| First differenced series | | 1.35 | 1.85 | 54.57 | 0.545 | 0.564 | 0.549 | 0.542 | 0.581 | 0.515 |
| Returns series | | 1.87 | 2.63 | 54.44 | 0.538 | 0.565 | 0.518 | 0.575 | 0.585 | 0.497 |

Table 10 Performance comparison of the LSTM model for BTC dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F ₁ | Sen | Spe | PPV | NPV |
|--------------------------|---------|--------|--------|---------|-------|----------------|-------|-------|-------|-------|
| Time-series | 4 | 404.57 | 551.17 | 51.17 | 0.500 | 0.462 | 0.693 | 0.294 | 0.501 | 0.524 |
| First differenced series | | 275.52 | 428.92 | 52.54 | 0.524 | 0.499 | 0.501 | 0.547 | 0.500 | 0.549 |
| Returns series | | 256.53 | 399.66 | 51.08 | 0.514 | 0.529 | 0.580 | 0.448 | 0.499 | 0.542 |
| Time-series | 6 | 436.77 | 576.23 | 50.85 | 0.501 | 0.446 | 0.653 | 0.348 | 0.509 | 0.521 |
| First differenced series | | 297.88 | 459.85 | 53.66 | 0.538 | 0.532 | 0.558 | 0.517 | 0.521 | 0.565 |
| Returns series | | 302.62 | 469.09 | 50.01 | 0.491 | 0.480 | 0.497 | 0.484 | 0.486 | 0.515 |
| Time-series | 9 | 447.30 | 617.72 | 49.48 | 0.498 | 0.475 | 0.703 | 0.259 | 0.509 | 0.526 |
| First differenced series | | 338.67 | 497.64 | 53.15 | 0.529 | 0.491 | 0.488 | 0.571 | 0.526 | 0.555 |
| Returns series | | 352.73 | 518.96 | 50.86 | 0.510 | 0.505 | 0.533 | 0.487 | 0.509 | 0.529 |
| Time-series | 12 | 483.96 | 672.34 | 50.19 | 0.497 | 0.477 | 0.702 | 0.396 | 0.508 | 0.510 |
| First differenced series | | 348.62 | 515.16 | 51.08 | 0.514 | 0.521 | 0.569 | 0.458 | 0.509 | 0.541 |
| Returns series | | 363.98 | 546.14 | 50.70 | 0.510 | 0.518 | 0.565 | 0.455 | 0.506 | 0.535 |

787 model’s predictions may be inefficient. In contrast, all ACF
 788 plots of the LSTM trained with the first differences and
 789 returns of prices reveal that there is no autocorrelation in
 790 the residuals which suggests the reliability of the model
 791 and advocates the efficiency of its forecasts.

792 Figures 7, 8 and 9 show the ACF plots of CNN–LSTM
 793 model for S&P500, Brent and BTC datasets, respectively.
 794 Clearly, all ACF plots of CNN–LSTM model trained with
 795 the first differenced and returns series illustrate that there
 796 exists no autocorrelation in the residuals. This implies that
 797 the model is reliable, with respect to the efficiency of its
 798 forecasts. On the other hand, the significant spikes presented
 799 in Figs. 7a, 8a and 9a reveal the CNN–LSTM trained with
 800 the traditional time-series has not properly fitted the training
 801 data and exhibited unreliable predictions.

7 Discussion

802
 803 In this section, we perform a discussion relative to the
 804 theoretical and experimental contribution of our research.

805 We presented a detailed theoretical background
 806 regarding the problem of time-series forecasting and the
 807 reliability of the forecasts of a prediction model. Since
 808 most time-series datasets are extremely noisy and chaotic
 809 by nature, the development of a reliable deep learning
 810 prediction models is considered a significantly challenging
 811 task. Moreover, the achievement of high accuracy or low
 812 RMSE score cannot be considered as a reliable metric since
 813 a model may just accidentally perform well on a specific
 814 time period, while on a new different period, it may exhibit
 815 a totally different and probably poor prediction
 816 performance.

Table 11 Performance comparison of the CNN–LSTM model for BTC dataset

| Series | Horizon | MAE | RMSE | Acc (%) | AUC | F ₁ | Sen | Spe | PPV | NPV |
|--------------------------|---------|--------|--------|---------|-------|----------------|-------|-------|-------|-------|
| Time-series | 4 | 315.41 | 448.46 | 49.58 | 0.489 | 0.512 | 0.693 | 0.277 | 0.488 | 0.516 |
| First differenced series | | 270.19 | 421.57 | 53.33 | 0.535 | 0.533 | 0.562 | 0.507 | 0.527 | 0.563 |
| Returns series | | 288.38 | 436.99 | 46.81 | 0.471 | 0.487 | 0.534 | 0.409 | 0.479 | 0.493 |
| Time-series | 6 | 366.81 | 514.74 | 52.21 | 0.513 | 0.476 | 0.489 | 0.552 | 0.518 | 0.539 |
| First differenced series | | 269.45 | 415.27 | 53.99 | 0.524 | 0.511 | 0.525 | 0.523 | 0.519 | 0.549 |
| Returns series | | 256.97 | 398.49 | 53.94 | 0.540 | 0.527 | 0.544 | 0.536 | 0.524 | 0.566 |
| Time-series | 9 | 399.61 | 534.20 | 51.41 | 0.499 | 0.531 | 0.584 | 0.383 | 0.500 | 0.536 |
| First differenced series | | 258.88 | 396.16 | 51.24 | 0.514 | 0.515 | 0.546 | 0.482 | 0.507 | 0.541 |
| Returns series | | 253.61 | 393.25 | 50.00 | 0.501 | 0.494 | 0.517 | 0.485 | 0.505 | 0.527 |
| Time-series | 12 | 417.20 | 553.32 | 51.40 | 0.510 | 0.510 | 0.592 | 0.399 | 0.511 | 0.523 |
| First differenced series | | 258.15 | 392.60 | 50.70 | 0.507 | 0.488 | 0.509 | 0.483 | 0.505 | 0.538 |
| Returns series | | 252.12 | 389.97 | 50.52 | 0.506 | 0.499 | 0.520 | 0.492 | 0.503 | 0.531 |

817 In this research, we demonstrated theoretically whether
 818 time-series data are “suitable” for fitting a deep learning
 819 model, which constitutes the main contribution of our
 820 research. In other words, we introduced a novel framework
 821 which can efficiently identify if a time-series is suitable for
 822 developing and training a deep learning model, which will
 823 perform reliable and stable prediction performances, inde-
 824 pendent of the characteristics of the series in any time period.

825 By the term “suitable,” we mean that the time-series
 826 data has successfully passed our proposed theoretical cri-
 827 teria and it can be used for training a prediction model. In
 828 contrast, if the series fails satisfying the requested criteria,
 829 then it is considered as “unsuitable” and every attempt for
 830 building a reliable prediction model will be probably in
 831 vain. Therefore, we provide a “starting point” for any
 832 attempt on developing any prediction framework for any
 833 time-series forecasting problem. This starting point is
 834 indeed the critical point in which every attempt and
 835 investment for building a model will result in a stable and
 836 reliable predictor or it will be totally wasted out in the case

837 that the utilized starting dataset was unsuitable. By iden-
 838 tifying the suitability of any time-series data, a “green
 839 light” for the machine learning developer is provided, in
 840 order to invest computational effort for building a fore-
 841 casting framework.

842 Furthermore, we established a novel and complete
 843 framework which provides a solution for any “unsuitable”
 844 identified time-series by performing a transformation based
 845 on first differences or returns and transform these series to
 846 “suitable.” Although these two techniques were well
 847 known as a rule of thumb for transformation and prepro-
 848 cessing for time-series data, it was not proved why, when
 849 and how these formulae work and if they actually can be
 850 successfully applied. Most approaches were relying on a
 851 “trial and error” logic something not appropriate and
 852 viable especially on cases when costly and time-consuming
 853 real-world projects aim to build accurate and reliable
 854 forecasting models. In this work, we proved that these
 855 formulae actually filtered these “unsuitable” data,

Table 12 Ljung–Box *Q* test for 10 lags with significance level $\alpha = 5\%$ (LSTM)

| Series | Horizon | <i>p</i> value | <i>H</i> ₀ | <i>p</i> value | <i>H</i> ₀ | <i>p</i> value | <i>H</i> ₀ |
|--------------------------|---------|----------------|-----------------------|----------------|-----------------------|----------------|-----------------------|
| Time-series | 4 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.170 | Accepted | 0.692 | Accepted | 0.764 | Accepted |
| Returns series | | 0.071 | Accepted | 0.370 | Accepted | 0.228 | Accepted |
| Time-series | 6 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.295 | Accepted | 0.271 | Accepted | 0.052 | Accepted |
| Returns series | | 0.486 | Accepted | 0.733 | Accepted | 0.383 | Accepted |
| Time-series | 9 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.286 | Accepted | 0.143 | Accepted | 0.235 | Accepted |
| Returns series | | 0.244 | Accepted | 0.113 | Accepted | 0.110 | Accepted |
| Time-series | 12 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.325 | Accepted | 0.790 | Accepted | 0.391 | Accepted |
| Returns series | | 0.423 | Accepted | 0.172 | Accepted | 0.400 | Accepted |

Table 13 Ljung–Box Q test for 10 lags with significance level $\alpha = 5\%$ (CNN–LSTM)

| Series | Horizon | p value | Null H_0 | p value | Null H_0 | p value | Null H_0 |
|--------------------------|---------|-----------|------------|-----------|------------|-----------|------------|
| Time-series | 4 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.102 | Accepted | 0.622 | Accepted | 0.834 | Accepted |
| Returns series | | 0.137 | Accepted | 0.344 | Accepted | 0.254 | Accepted |
| Time-series | 6 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.657 | Accepted | 0.416 | Accepted | 0.160 | Accepted |
| Returns series | | 0.180 | Accepted | 0.168 | Accepted | 0.064 | Accepted |
| Time-series | 9 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.105 | Accepted | 0.083 | Accepted | 0.560 | Accepted |
| Returns series | | 0.619 | Accepted | 0.382 | Accepted | 0.435 | Accepted |
| Time-series | 12 | 0.000 | Rejected | 0.000 | Rejected | 0.000 | Rejected |
| First differenced series | | 0.059 | Accepted | 0.060 | Accepted | 0.629 | Accepted |
| Returns series | | 0.056 | Accepted | 0.869 | Accepted | 0.217 | Accepted |

856 eliminating this costly and time consuming “trial and
857 error” approach.

858 It is worth mentioning that an attractive property of our
859 proposed framework is that it can be easily extended to
860 cover the wider scientific area of time-series forecasting
861 applications without the requirement of any extra modifi-
862 cations or additional constraints. In more detail, the pro-
863 posed framework performs an efficient preprocessing step
864 in order to exploit the internal representation of the times-
865 eries, through the utilization of statistic and econometric
866 test. Conclusively, we point out that our experimental
867 analysis indicated that although deep learning models
868 constitute a widely accepted and efficient choice for time-
869 series forecasting, our proposed framework provides a
870 significant boost in increasing the forecasting performance.

871 Nevertheless, an extensive research is under considera-
872 tion to identify which of these two methodologies can be a
873 priori efficiently applied depending on the characteristics
874 of each time-series in order to obtain better prediction
875 results. A possible approach could be the application of a
876 sophisticated preprocessing framework based on the
877 intrinsic time-series specific properties such as stationarity,
878 heteroskedasticity, seasonal cycles and changing variance,
879 for performing that a priori identification and the proper
880 time-series transformation methodology.

881 8 Conclusions and future research

882 Time-series forecasting and analysis is generally consid-
883 ered as one of the most challenging problems in data
884 mining. In the literature, most time-series forecasting
885 approaches attempt to exploit machine learning and deep
886 learning algorithms, aiming at obtaining better perfor-
887 mance compared to the already existing or proposed
888 models. Nevertheless, they cannot guarantee to develop
889 reliable forecasting models.

890 In this work, we propose a different approach and
891 introduce a novel methodology for the development of
892 efficient and reliable deep learning prediction models. The
893 major novelty of our proposed framework is that it guar-
894 antees the forecasting reliability of the deep learning
895 model’s predictions, independent of the used time-series
896 data. This is achieved by applying a series of transforma-
897 tions, which ensure that a time-series satisfies the station-
898 arity property and it is suitable for fitting a deep learning
899 model. In addition to the theoretical advantages of the
900 proposed framework, we provided empirical evidence
901 about its efficiency and robustness. More specifically, we
902 performed a series of numerical experiments using time-
903 series from three application domains, which attracted most
904 of research interest, namely financial stock market, energy
905 sector and cryptocurrency area. All compared models
906 where evaluated on both forecasting time-series price (re-
907 gression) and time-series directional movements (classifi-
908 cation) as well as on the reliability of their forecasts by
909 examining the existence of autocorrelation of the errors.
910 Our comprehensive experimental analysis illustrated that
911 our proposed methodology considerably improved the
912 forecasting performance of a deep learning model, in terms
913 of accuracy and reliability.

914 By taking into consideration that our proposed frame-
915 work can be easily exploit any deep learning model, a
916 prediction model exhibiting even better forecasting ability
917 could be developed through the exploitation of deep
918 learning techniques together with regularization method-
919 ologies or through additional optimized configuration of
920 the utilized models.

921 It is worth mentioning that the introduced framework
922 can be easily extended to cover the wider scientific area of
923 time series forecasting applications such as weather fore-
924 casting, earthquake prediction, heartbeat rate and so on,
925 without the requirement of any extra modifications or
926 additional constraints. Furthermore, one issue which we

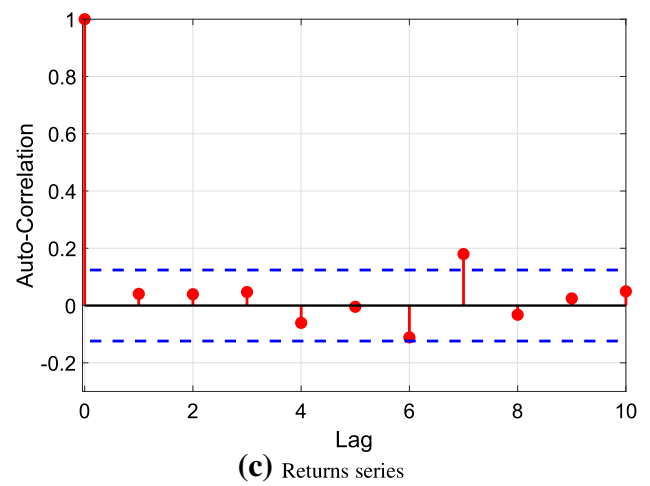
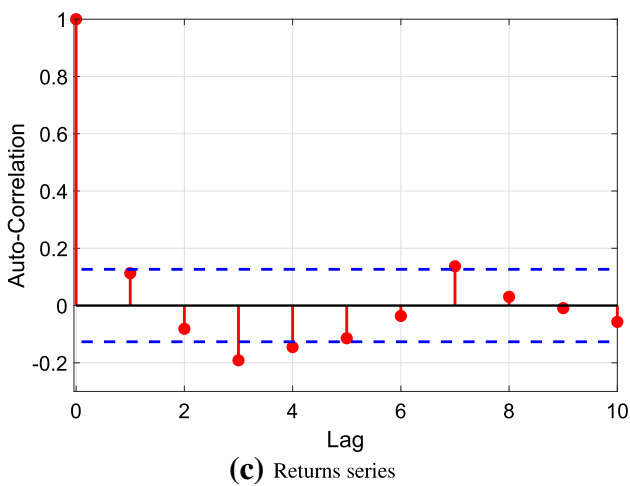
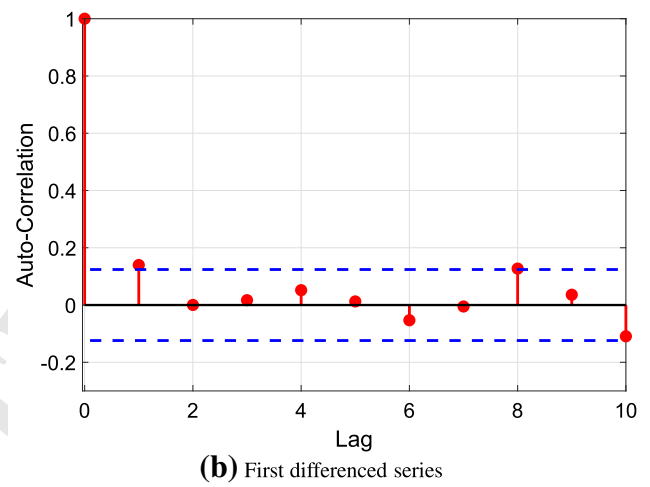
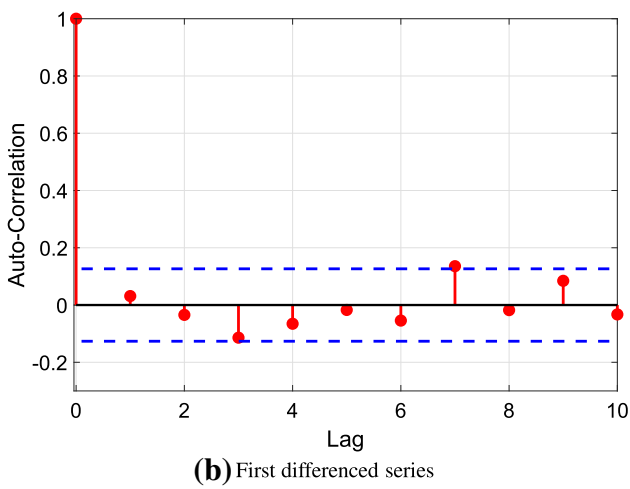
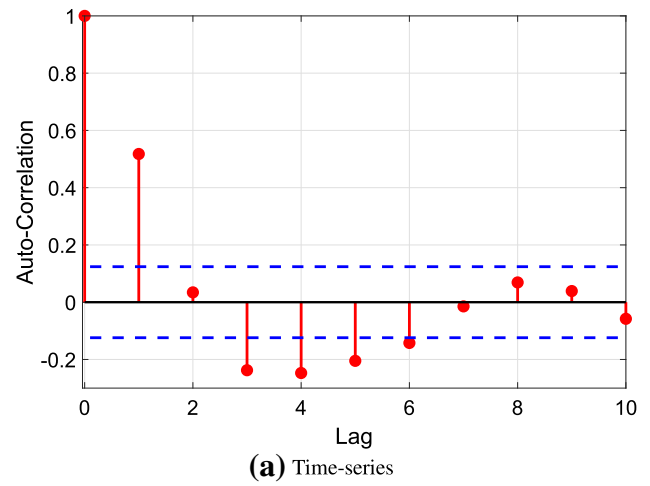
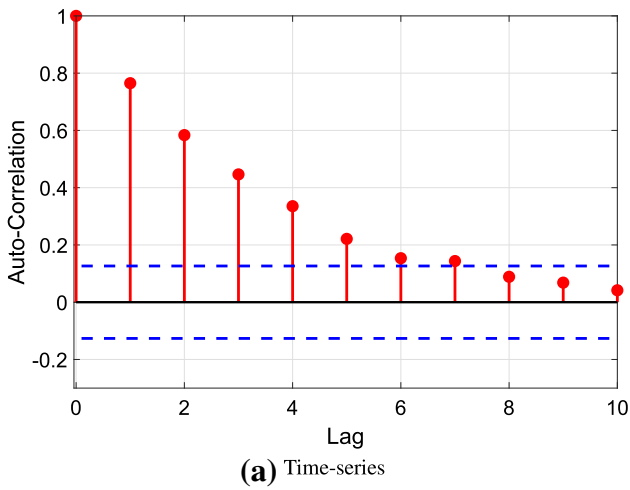


Fig. 4 Autocorrelation of residuals for S&P500 dataset of LSTM model

Fig. 5 Autocorrelation of residuals for Brent dataset of LSTM model

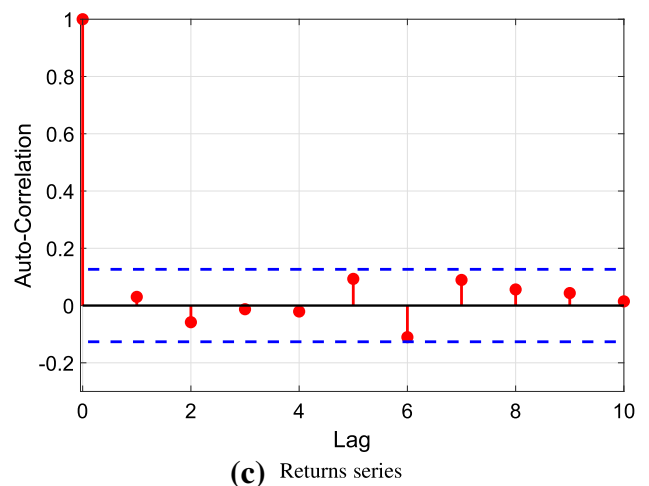
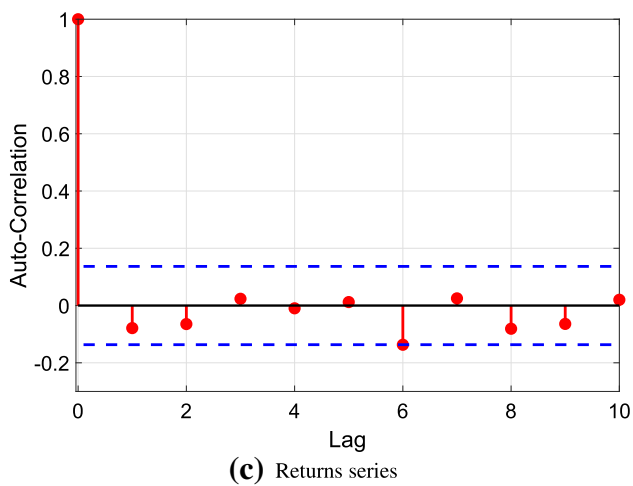
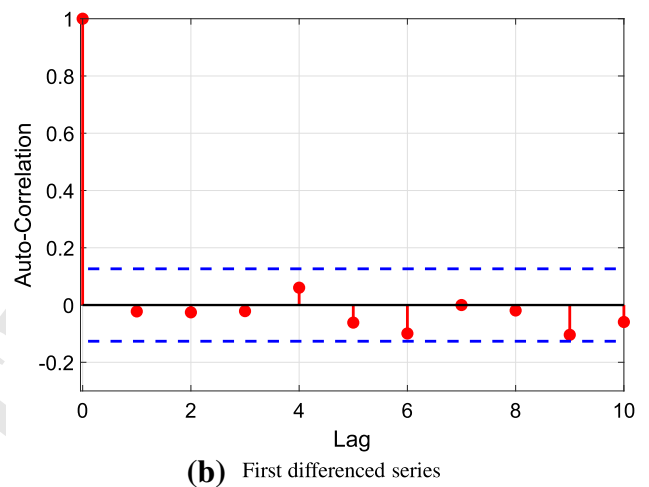
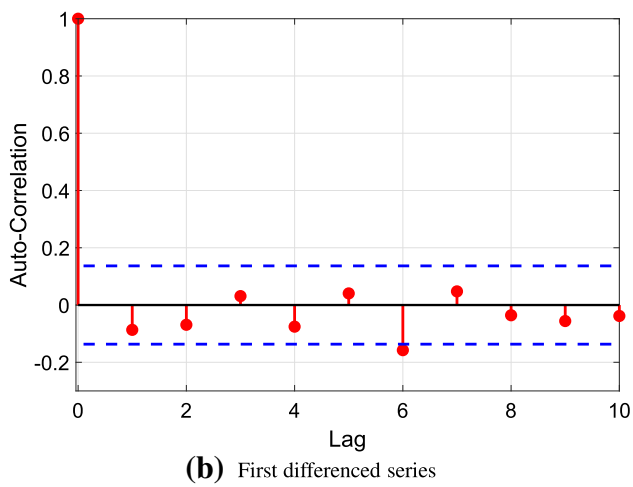
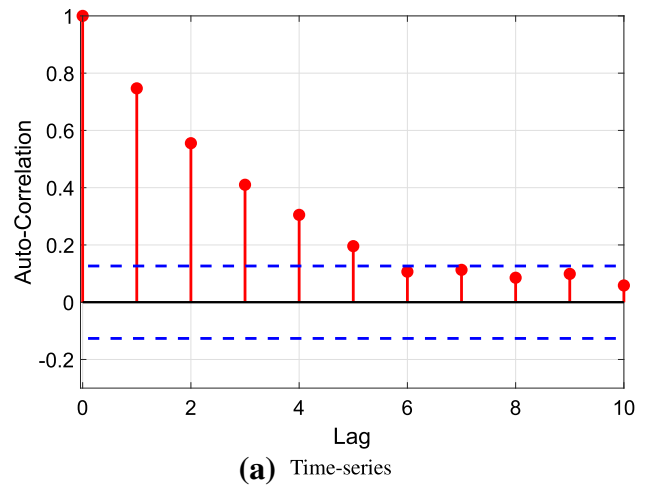
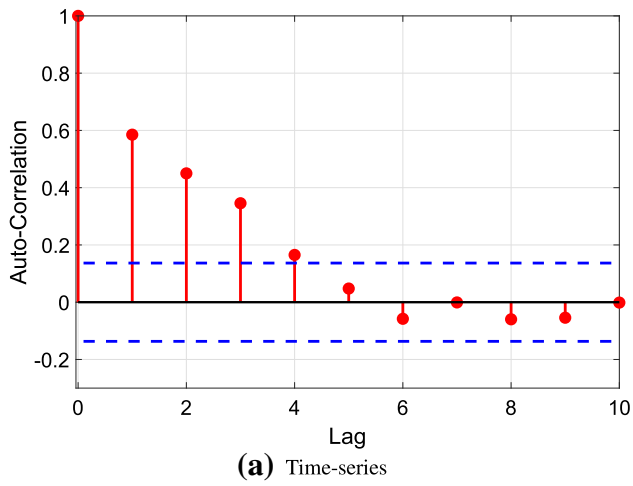


Fig. 6 Autocorrelation of residuals for BTC dataset of LSTM model

Fig. 7 Autocorrelation of residuals for S&P500 dataset of CNN-LSTM model

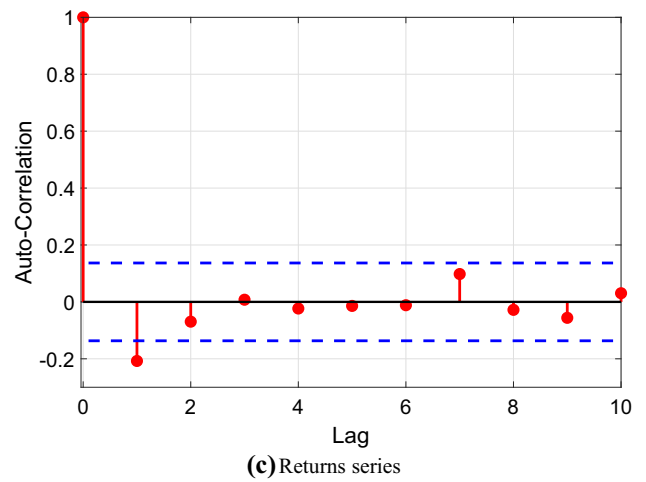
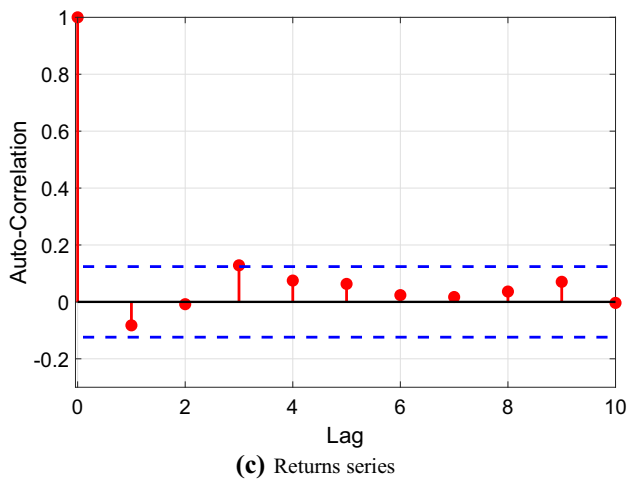
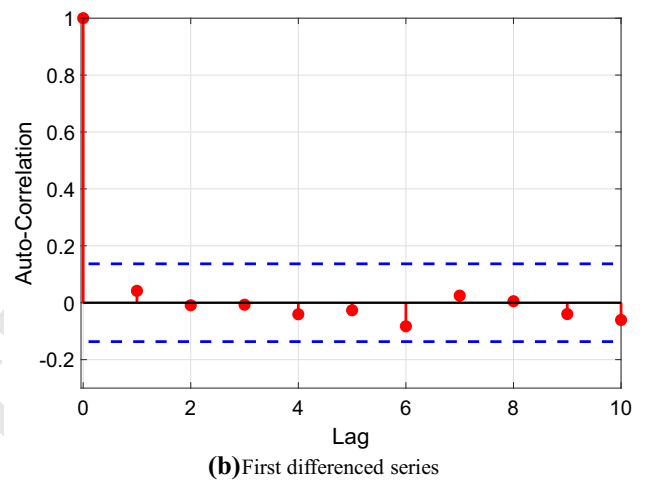
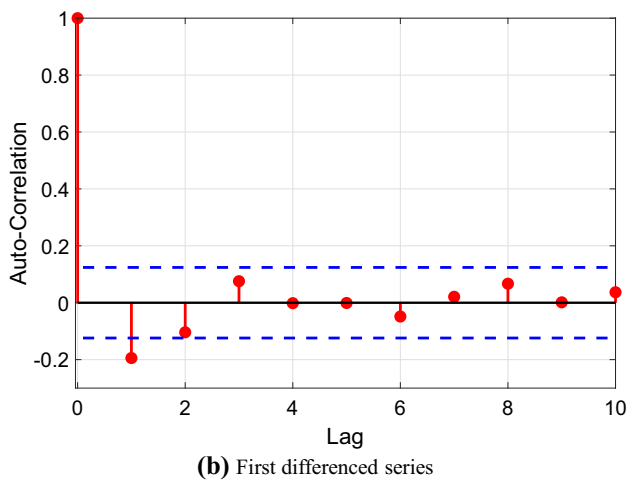
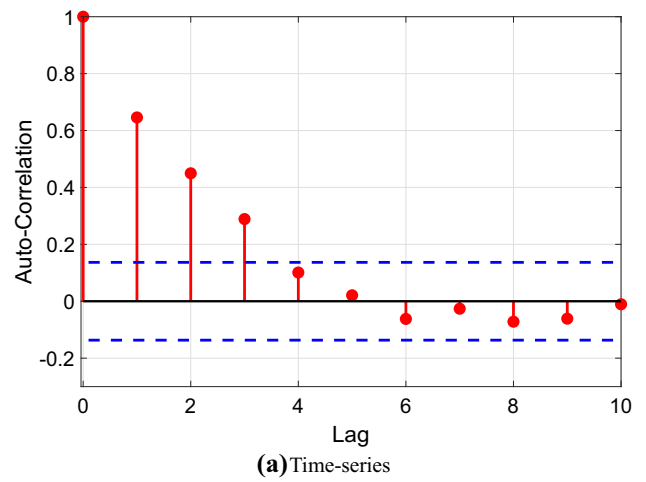
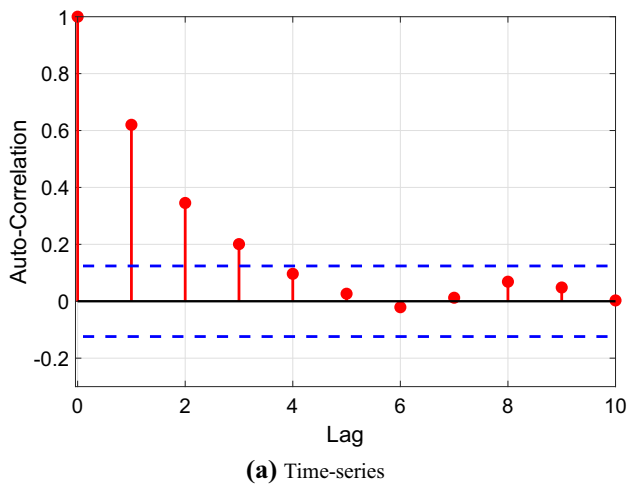


Fig. 8 Autocorrelation of residuals for Brent dataset of CNN-LSTM model

Fig. 9 Autocorrelation of residuals for BTC dataset of CNN-LSTM model

927 have not thoroughly investigated is the possibility of some
928 minor information loss due to non-stationarity in the time-
929 series and the imposition of the proposed transformations.
930 This is to be included and fully investigated in our future
931 research. Furthermore, we intend to verify that the pro-
932 posed framework works with any kind of regression
933 algorithm.

934 Another direction for future research is to enhance our
935 experimental framework with new performance metrics
936 based on profits and returns. Finally, an interesting idea is the
937 application of our proposed framework for the prediction of
938 anomaly detection in order to “catch” outliers or other rare
939 signals, which could indicate forecasting instability.
940
941

942 Compliance with ethical standards

943 **Conflicts of interest** The authors declared no potential conflicts of
944 interest with respect to the research, authorship and/or publication of
945 this article.

946 References

- 947 1. Ahmed NK, Atiya AF, Gayar NE, El-Shishiny H (2010) An
948 empirical comparison of machine learning models for time series
949 forecasting. *Econom Rev* 29(5–6):594–621
- 950 2. Bengio Y, Courville A, Vincent P (2013) Representation learn-
951 ing: a review and new perspectives. *IEEE Trans Pattern Anal*
952 *Mach Intell* 35(8):1798–1828
- 953 3. Bontempi G, Taieb SB, Le Borgne Y (2012) Machine learning
954 strategies for time series forecasting. In: *European business*
955 *intelligence summer school*. Springer, Berlin, pp 62–77
- 956 4. Bougoudis I, Demertzis K, Iliadis L (2016) HISYCOL a hybrid
957 computational intelligence system for combined machine learn-
958 ing: the case of air pollution modeling in Athens. *Neural Comput*
959 *Appl* 27(5):1191–1206
- 960 5. Box GEP, Jenkins GM, Reinsel GC, Ljung GM (2015) *Time*
961 *series analysis: forecasting and control*. Wiley, London
- 962 6. Brockwell PJ, Davis RA (2016) *Introduction to time series and*
963 *forecasting*. Springer, Berlin
- 964 7. Burke S, Hunter J (2005) *Modelling non-stationary economic*
965 *time series: a multivariate approach*. Springer, Berlin
- 966 8. Cen Z, Wang J (2019) Crude oil price prediction model with long
967 short term memory deep learning based on prior knowledge data
968 transfer. *Energy* 169:160–171
- 969 9. Chambon S, Galtier MN, Arnal PJ, Wainrib G, Gramfort A
970 (2018) A deep learning architecture for temporal sleep stage
971 classification using multivariate and multimodal time series.
972 *IEEE Trans Neural Syst Rehabil Eng* 26(4):758–769
- 973 10. Donate JP, Li X, Sánchez GG, de Miguel AS (2013) Time series
974 forecasting by evolving artificial neural networks with genetic
975 algorithms, differential evolution and estimation of distribution
976 algorithm. *Neural Comput Appl* 22(1):11–20
- 977 11. Fawaz HI, Forestier G, Weber J, Idoumghar L, Muller P (2019)
978 Deep learning for time series classification: a review. *Data Min*
979 *Knowl Disc* 33(4):917–963
- 980 12. Fischer T, Krauss C (2018) Deep learning with long short-term
981 memory networks for financial market predictions. *Eur J Oper*
982 *Res* 270(2):654–669

13. Fu T (2011) A review on time series data mining. *Eng Appl Artif*
983 *Intell* 24(1):164–181
14. Gocheva-Ilieva SG, Voynikova DS, Stoimenova MP, Ivanov AV,
984 Iliev IP (2019) Regression trees modeling of time series for air
985 pollution analysis and forecasting. *Neural Comput Appl*
986 *Appl* 31(12):9023–9039
15. Gulli A, Pal S (2017) *Deep learning with Keras*. Packt Publishing Ltd
987
16. Ji S, Kim J, Im H (2019) A comparative study of Bitcoin price
988 prediction using deep learning. *Mathematics* 7(10):898
17. Kingma DP, Ba J (2015) Adam: A method for stochastic opti-
989 mization. In: *2015 International conference on learning*
990 *representations*
18. Liu S, Zhang C, Ma J (2017) CNN-LSTM neural network model
991 for quantitative strategy analysis in stock markets. In: *Internat-*
992 *ional conference on neural information processing*. Springer,
993 Berlin, pp 198–206
19. Livieris IE, Pintelas E, Kiriakidou N, Stavroyiannis S (2020) An
994 advanced deep learning model for short-term forecasting U.S.
995 natural gas price and movement. In: *16th International conference*
996 *on artificial intelligence applications and innovations*
20. Livieris IE, Pintelas E, Pintelas P (2020) A CNN-LSTM model
997 for gold price time series forecasting. In: *Neural computing and*
998 *applications*
21. Nakano M, Takahashi A, Takahashi S (2018) Bitcoin technical
999 trading with artificial neural network. *Physica A* 510:587–609
22. Osborne J (2010) Improving your data transformations: applying
1000 the Box-Cox transformation. *Pract Assess Res Eval* 15(1):12
23. Pal A, Prakash PKS (2017) *Practical time series analysis: master*
1001 *time series data processing, visualization, and modeling using*
1002 *python*. Packt Publishing Ltd
24. Pintelas E, Livieris IE, Stavroyiannis S, Kotsilieris T, Pintelas P
1003 (2020) Fundamental research questions and proposals on pre-
1004 dicting cryptocurrency prices using DNNs. Technical Report
1005 TR20-01, University of Patras, [https://nemertes.lis.upatras.gr/](https://nemertes.lis.upatras.gr/jspui/bitstream/10889/13296/1/TR01-20.pdf)
1006 [jspui/bitstream/10889/13296/1/TR01-20.pdf](https://nemertes.lis.upatras.gr/jspui/bitstream/10889/13296/1/TR01-20.pdf)
25. Pintelas E, Livieris IE, Stavroyiannis S, Kotsilieris T, Pintelas P
1007 (2020) Investigating the problem of cryptocurrency price pre-
1008 diction: a deep learning approach. In: *16th International confer-*
1009 *ence on artificial intelligence applications and innovations*
26. Shaman P (2010) Generalized Levinson–Durbin sequences,
1010 binomial coefficients and autoregressive estimation. *J Multivar*
1011 *Anal* 101(5):1263–1273
27. Stavroyiannis S (2019) Can Bitcoin diversify significantly a
1012 portfolio? *Int J Econ Bus Res* 18(4):399–411
28. Tanaka K (2017) *Time series analysis: nonstationary and non-*
1013 *invertible distribution theory*, vol 4. Wiley, London
29. Urtnasan E, Park J, Lee K (2018) Automatic detection of sleep-
1014 disordered breathing events using recurrent neural networks from
1015 an electrocardiogram signal. In: *Neural computing and applica-*
1016 *tions*, pp 1–10
30. Weigend AS (2018) *Time series prediction: forecasting the future*
1017 *and understanding the past*. Routledge, London
31. Xingjian SHI, Chen Z, Wang H, Yeung D, Wong W, Woo W
1018 (2015) Convolutional LSTM network: a machine learning
1019 approach for precipitation nowcasting. In: *Advances in neural*
1020 *information processing systems*, pp 802–810
32. Zhao Y, Li J, Yu L (2017) A deep learning ensemble approach
1021 for crude oil price forecasting. *Energy Econ* 66:9–16

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