



Article Minimum Message Length in Hybrid ARMA and LSTM Model Forecasting

Zheng Fang ¹, David L. Dowe ^{1,*}, Shelton Peiris ² and Dedi Rosadi ³

- ¹ Department of Data Science and Artificial Intelligence, Monash University, Clayton, VIC 3800, Australia; zfan51@student.monash.edu
- ² School of Mathematics and Statistics, University of Sydney, Camperdown, NSW 2006, Australia; shelton.peiris@sydney.edu.au
- ³ Department of Statistics, Gadjah Mada University, Sleman, Yogyakarta 55500, Indonesia; dedirosadi@gadjahmada.edu
- Correspondence: david.dowe@monash.edu

Abstract: Modeling and analysis of time series are important in applications including economics, engineering, environmental science and social science. Selecting the best time series model with accurate parameters in forecasting is a challenging objective for scientists and academic researchers. Hybrid models combining neural networks and traditional Autoregressive Moving Average (ARMA) models are being used to improve the accuracy of modeling and forecasting time series. Most of the existing time series models are selected by information-theoretic approaches, such as AIC, BIC, and HQ. This paper revisits a model selection technique based on Minimum Message Length (MML) and investigates its use in hybrid time series analysis. MML is a Bayesian information-theoretic approach and has been used in selecting the best ARMA model. We utilize the long short-term memory (LSTM) approach to construct a hybrid ARMA-LSTM model and show that MML performs better than AIC, BIC, and HQ in selecting the model—both in the traditional ARMA models (without LSTM) and with hybrid ARMA-LSTM models. These results held on simulated data and both real-world datasets that we considered. We also develop a simple MML ARIMA model.

Keywords: long short-term memory; minimum message length; time series; neural network; deep learning; Bayesian statistics; probabilistic modeling

1. Introduction

Forecasting in time series is a difficult task due to the presence of trends and/or seasonal components. For example, economic time series data are highly impacted by seasonal factors and often show trends with long-run cycles. Such trends and seasonality are difficult to capture by the traditional Autoregressive Moving Average model (ARMA) [1]. The Bayesian Minimum Message Length (MML) principle [2], the Akaike Information Criterion (AIC) [3], Schwarz's Bayesian information criterion (BIC) [4] and Hannan–Quinn (HQ) [5] are often used in model selection for the ARMA model [6–8]. The models selected by MML87 [9] in ARMA time series have lower prediction errors than those from AIC, BIC, and HQ [10]. Schmidt previously showed that MML87 outperforms a variety of other (information-theoretic) approaches in ARMA time series modeling [11] (chapters 5 to 8). In this paper, we extended the traditional ARMA time series model to form the hybrid ARMA-LSTM by combining the neural network of long short-term memory (LSTM) in order to test the performance of MML in model selection. The results suggest that MML outperforms AIC, BIC, BIC, BIC, and HQ.

The ARIMA is used with integer differencing to achieve stationarity if the time series is not stationary. A time series with seasonal components can be modeled using the family of seasonal ARIMA (or SARIMA) models. On the other hand, this ARIMA family has been generated to include long memory time series using a suitable fractional order



Citation: Fang, Z.; Dowe, D.L.; Peiris, S.; Rosadi, D. Minimum Message Length in Hybrid ARMA and LSTM Model Forecasting. *Entropy* **2021**, *23*, 1601. https://doi.org/10.3390/ e23121601

Academic Editors: Eric Nalisnick and Dustin Tran

Received: 30 September 2021 Accepted: 25 November 2021 Published: 29 November 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). differencing in (0, 0.5) to form the family of autoregressive fractionally integrated moving average (ARFIMA) models. Nevertheless, the deep learning LSTM technique might be more suitable to capture the information that is less obvious in the time series, as it allows for a much more general class of models. Time series analysts require a lot of effort to discover the appropriate model in order to identify the dependency in time series data [12]. Historically, the ARMA model was introduced by Box and Jenkins in 1976 [13], and it is popular and widely used in the time series science community and provides accurate forecasts in both in-sample and out-of-sample data when the parameters are correctly estimated [14]. It is a hybrid (or mixture) of autoregressive (AR) and moving average (MA) processes, but the ARMA model can only be used in stationary time series [15].

In parallel, machine learning has seen the development of neural network models in computer science, ultimately influencing statistics. Similar to the families of the ARMA model, deep learning also has several variants, such as a deep neural network (DNN), a convolutional neural network (CNN), and a recurrent neural network (RNN). This report investigates a particular form of the RNN called long short-term memory (LSTM), which is typically used in time series [16]. In recent years, LSTM has been shown to work well in forecasting for data with complex time dependency, such as the stock market and energy consumption prediction [17]. In this paper, we select the best ARMA(p,q) model and then train the LSTM model for the residuals through the ARMA model. The time-step order used in LSTM is the parameter q in ARMA(p, q) determined by different information-theoretic criteria [18,19].

Our results show that MML compares favorably with the other information-theoretic approaches, including AIC, BIC, and HQ, when conducting ARMA-LSTM. Further, we compare the ARMA-LSTM selected by MML with the ARMA model selected by MML. These results also show that MML outperforms when compared to AIC [20], BIC [20], and HQ [21] in terms of selecting a model with lower prediction error, and this holds whether our modeling is enhanced by LSTM or instead is ARMA unassisted by LSTM. The Bayesian information-theoretic MML principle provides more reliable and highly accurate results in the model selection of the hybrid ARMA-LSTM model than other traditional methods (AIC, BIC, HQ). When doing ARMA without a hybrid with LSTM, MML also performs better than other traditional methods (AIC, BIC, HQ). The best performing method considered is the hybrid MML ARMA-LSTM model. These results hold on simulated data and on the real-world datasets considered.

Section 2 introduces the Box and Jenkins theory for the ARIMA model and discusses its limitations. Section 3 introduces the information-theoretic Minimum Message Length criterion in model selection, and Section 4 introduces the deep learning model LSTM. Section 5 provides the algorithm of the hybrid ARMA-LSTM model, and Section 6 provides the experimental results with a comparison.

2. ARIMA Modeling

This section reviews the theory of Autoregressive Integrated Moving Average (ARIMA) modeling from Box and Jenkins (1970) [13,15]. Let $\{Y_t\}$ be a homogeneous nonstationary time series and suppose that the d^{th} (d = 1, 2, ...) difference of the series is stationary and is given by $X_t = (1 - B)^d Y_t$, where *B* is the backshift operator. Then a stationary ARMA(p,q) model can be fitted for $\{X_t\}$, satisfying

$$X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon_t + \sum_{i=1}^q \theta_i \epsilon_{t-i},$$
(1)

where $\{\epsilon_t\} \sim WN(0, \sigma^2)$.

Let $\phi(B) = 1 - \phi_1 B - \ldots - \phi_p B^p$; $\theta(B) = 1 + \theta_1 B + \ldots + \theta_q B^q$, be two polynomials of degree *p* and *q*, respectively, such that the zeros of $\phi(B)$ and $\theta(B)$ are outside the unit circle. Then the ARMA(p,q) in Equation (1) can be written in a compact form as

Now the corresponding ARIMA(p,d,q) model for the original series $\{Y_t\}$ is given by

$$\phi(B)(1-B)^{d}Y_{t} = c + \theta(B)\epsilon_{t}.$$
(3)

It is known that ARIMA is a form of a linear regression model with the lag order of time series data and corresponding residuals. In an application where the ARIMA model fits well for the given data, then the corresponding residuals through the model should form a random scatter plot with a constant mean and a constant variance over the time, see, for example, [13]. If the ARIMA model is not well fitted for the data or an incorrect model has been fitted, then the residuals will not show a random scatter plot and instead indicate autocorrelations within the residuals. This reveals that the information hidden in the data has not been completely captured by the fitted ARIMA model, and we consider refitting an alternative ARIMA model [22].

The above family of ARIMA models are also capable of modeling a wide range of seasonal data using slight modifications. A seasonal extension of Model (3) can be written for a set of time series data with seasonality *m*. Incorporating both the seasonal and nonseasonal components together with additional polynomials, a new model is

$$\phi(B)\Phi(B^m)(1-B)^d(1-B^m)^D Y_t = c + \theta(B)\Theta(B^m)\epsilon_t,\tag{4}$$

where $\Phi(B^m) = 1 - \Phi_1 B^m - \ldots - \Phi_P B^{mP}$, $\Theta(B^m) = 1 + \Theta_1 B^m + \ldots + \Theta_Q B^{mQ}$, and *D* is the degree of seasonal differencing. For simplicity, this is written as

$$Y_t \sim SARIMA(p,d,q)(P,D,Q)_m \tag{5}$$

Model (4) is known as the Seasonal ARIMA or SARIMA model.

To estimate the parameters of Model (4), it is important to identify the changes of variance in the autocorrelation function (ACF) plot of data. This ACF provides an indication of linear dependencies among the observation of time series, which is related to the order of the model. In addition, the corresponding partial autocorrelation function (PACF) can be used to confirm the approximate order required in the model.

In this study, we use non-seasonal ARIMA modeling because the non-seasonal degree of differencing *d* can be predetermined in practice. We consider the stationary time series data. Assuming the data are generated from a mean zero stationary ARMA(*p*, *q*) process with Gaussian errors, we use the fact that the distribution of data is a multivariate Gaussian distribution with mean $\mu = 0$.

Suppose that we have a sample of *N* observations $y = (y_1, ..., y_N)$ generated through Model (2), with c = 0, and let $\beta = (\phi_1, ..., \phi_p, \theta_1, ..., \theta_q, \sigma^2)$ be the vector of all the parameters. Then the corresponding unconditional log-likelihood function, $L(y\beta)$, can be written as:

$$L(y\beta) = -\frac{N}{2}\log(2\pi\sigma^2) - \frac{1}{2}\log\Sigma - \frac{1}{2\sigma^2}y^T\Sigma^{-1}y,$$
 (6)

where Σ is the determinant of Σ and $\sigma^2 \Sigma$ is the $N \times N$ theoretical autocovariance matrix of *y*.

3. Minimum Message Length

The Bayesian information-theoretic Minimum Message Length (MML) principle [2,6,7,9,19,23] is based on coding theory and can be thought of in several equivalent ways. It can be thought of in terms of a transmitter encoding a two-part message and transmitting it to a receiver, where the first part of the message contains information encoding the model and the second part of the message encodes the data given the model. The length of the first part of the message can be thought of as the complexity of the model, and the length of the second part of the message (effectively, the statistical negative log-likelihood) is a measure of goodness of fit to the observed data. For example, with

 $X = \{A, B, C, D\}$, possible encodings would be, e.g., A = 00, B = 01, C = 10, and D = 11, or instead, e.g., A = 1, B = 01, C = 001, and D = 0001, with the length of code represented as I(), e.g., with A = 00, I(A) = 2. The code length is typically (close to) the negative logarithm of the probability.

MML thus gives a quantitative information-theoretic trade-off between model complexity (length of first part of message) and goodness of fit (length of second part of message) [24]. A smaller MML value (or, equivalently, a shorter message length) indicates the model is less complex and highly fitted to the data [6]. In practice, minimizing the message length can be expressed as:

$$\underset{\theta \in \Theta}{\arg\min} \{ \mathbf{I}(\theta) + \mathbf{I}(y^N \theta) \}, \tag{7}$$

where $I(\theta)$ is the length of encoding the assertion (or model), and $I(y^N\theta)$ is the length of encoding the detail (or data given the model). In MML, there is (Bayesian) prior knowledge (or a prior distribution), π , over the parameter space. Following Wallace and Freeman [9], MML has been shown to work well in time series models, such as autoregressive (AR) and moving average (MA) models [18,25,26]. We can thus estimate the parameters [7,9] by minimizing the message length:

$$MessLen(y,\beta) = -\log(\frac{h_3(\beta)f(y_1,...,y_N\beta)\epsilon^N}{\sqrt{F(\beta)}}) + \frac{k}{2}(1+\log\kappa_k) - \log h_1(p) - \log h_2(q), \quad (8)$$

where ϵ is measuring the accuracy of data, $h_3(\beta)$ is the Bayesian prior distribution over the parameter set β , we model the parameter set β using uniform prior [0, 1] in the stationarity region $h_3(\beta) = 1$, and $h_1(p) = 2^{-(1+p)}$ and $h_2(q) = 2^{-(1+q)}$ are the priors on the (non-negative integer) parameters p and q, k = p + q + 1 is the number of continuous-valued parameters, $f(y_1, ..., y_N\beta)$ is the standard statistical likelihood function, $L = -\log f$, $F(\beta)$ is the expected Fisher Information matrix (of expected second-order partial derivatives of L) and is a function of the parameter set β , $F(\beta)$ is the expected Fisher information, κ_k is the lattice constant (which accounts for the expected error in the log-likelihood function from ARMA model (Equation (6)) due to the quantization of the k-dimensional space, which is bounded above by $\frac{1}{12}$ and bounded below by $\frac{1}{2\pi e}$. For example, $\kappa_1 = \frac{1}{12}$, $\kappa_2 = \frac{5}{36\sqrt{3}}$, $\kappa_3 = \frac{19}{192*2^{1/3}}$, and $\kappa_k \to \frac{1}{2\pi e}$ as $k \to \infty$).

Ignoring the $-\log h_1(p)$, $-\log h_2(q)$, and $-N\log(\epsilon)$ terms, the message length for the ARMA model β can also be represented as:

$$I(y,\beta) = -\log h_3(\beta) + \frac{1}{2}\log F(\beta) + \frac{k}{2}\log \kappa_k + \frac{k}{2} - \log f(y\beta)$$
(9)

MML87 is model invariant and avoids explicitly constructing the quantized parameter space [7–9,23]. This is used for model selection and parameter estimation by choosing the model that minimizes the message length.

MML has been used for a variety of problems, including clustering and mixture modeling [27,28] ([19] Section 6.8), clustering of protein dihedral angles [29], decision graphs (as an extension of decision trees, allowing for disjunctions, or "or") [30] (Section 7.2.4 [19]) and multi-way joins in decision graphs with dynamic attributes [31], causal Bayesian nets (or Bayesian networks, or causal nets) ([19] Section 7.4) and Bayesian nets with decision trees in their (leaf) nodes [32,33], inference of probabilistic finite state automata (or probabilistic finite state machines, PFSAs, PFSMs) ([19] Section 7.1) and hierarchical PFSAs [34], and (given sufficient data and time, and based to whatever degree on the above-mentioned inference of Bayesian nets) automation of database normalization [35], etc.

Part of the reason for the above list is the universality of the MML approach [7] ([19] Chapter 2) (seeking the single best theory) and that of the predictive approach (seeking a Bayesian weighted combination of theories) of Solomonoff [36,37] ([38] Section 3.1).

learning often runs relatively quickly. This motivates us to combine these approaches, as we do using the deep learning approach of long short-term memory (LSTM). This gives us something of a combination of the simplicity and accuracy of MML and the speed of deep learning.

We note in passing that an earlier effort at combining MML with neural nets is [39]. We further note that some approaches to deep learning use a (suitably weighted) combination of a squared error term and a Kullback–Leibler divergence term. Given that squared error comes (or can come) from a Gaussian log-likelihood, this version of deep learning regularization bears similarities to D. F. Schmidt's MML approximation [11] ([6] footnotes 64 and 65).(The MMLD version of MML ([6] Section 0.2.2, p. 528) [40] ([19] Sections 4.10, 4.12.2 and 8.8.2, p. 360) modified MML87 [9] to allow for cases when the Bayesian prior is not approximately constant over the relevant region. D. F. Schmidt's MML approximation, just discussed, is a further modification, and explicitly introduces Kullback–Leibler divergence into the expression.) We also ask, for future work, whether our approach might be combined with graph neural networks [41] or (higher-dimensional) hyper-graph neural networks.

4. Long Short-Term Memory (LSTM)

With the development of computational power in electronic equipment, powerful computers provide many learning algorithms and approaches in time series forecasting [42–44]. Deep learning is one of the popular approaches in recent years; it provides a complex model that has at least the potential to capture (and often does capture) more general information from the predictors than a traditional model, such as ARMA. Long short-term memory (LSTM) is a special kind of recurrent neural network introduced by Hochreiter and Schmidhuber in 1997 [45]. LSTM manages the two state vectors, the short-term state h_t and long term state c_t , and uses the gating mechanism by adding linear components from the previous layer in order to provide the long memory. LSTM has been widely used in time series forecasting because it is able to capture more information in the time series data, particularly for the financial econometrics area, where the price of financial assets depends on various different factors that are difficult to represent by a linear model [44,46]. Each LSTM layer, including the cells of the forget gate, input gate, and output gate, is shown in Figure 1.

- Forget gate: $f_t = \sigma(U^f x_t + W^f h_{t-1} + b^f);$
- Input gate: $i_t = \sigma(U^i x_t + W^i h_{t-1} + b^i);$
- Output gate: $o_t = \sigma(U^o x_t + W^o h_{t-1} + b^o)$.

The forget gate uses a sigmoid function $\sigma(x)$ from Equation (10). It has a value between 0 and 1, and it determines how much information should be forgotten. If the result from the sigmoid function is close to 0, then more information should be forgotten, and if the result from the sigmoid function is close to 1, then less information should be forgotten.

$$\tau(x) = \frac{1}{1 + e^{-x}} \tag{10}$$

The input gate also uses the sigmoid function, the input gate controls the value input from the input function of $g_t = tanh(Wh_{t-1} + Ux_t + b)$ using the tanh(x) function:

(

$$tanh(x) = \frac{sinh(x)}{cosh(x)} = \frac{e^{x} - e^{-x}}{e^{x} + e^{-x}}$$
 (11)

The input gate controls how much information should be remembered. The LSTM long-term state uses an element-wise operation with $c_t = f_t \odot c_{t-1} + g_t \odot i_t$, where \odot is element-wise multiplication (of two matrices of the same dimension), also known as the Hadamard product.



Figure 1. LSTM Structure.

The output gate o_t controls how much long-term information c_t should be carried forward to the next layer, and it also contributes to the short-term state of h_t . The result from the output gate function is also between 0 and 1, and the LSTM short-term state also uses element-wise multiplication, with $h_t = o_t \odot tanh(c_t)$. An LSTM with more than one layer is shown in Figure 2, and its structure enables the LSTM to capture long-term and short-term information in order to forecast. As usual, an LSTM is trained by back propagation as other neural network models are. An LSTM requires time series data to train the model, and its time series pattern will be modeled in every layer of the network.



Figure 2. LSTM Overlapping.

5. Hybrid ARMA-LSTM Model

In recent years, LSTM and its variants—along with some hybrid models—have been thought by many to largely dominate the financial time series forecasting domain [44]. The LSTM is able to capture the dependency of residuals across time, and the LSTM is trained by the time step [47]. In this paper, we are using the Moving Average lag order q from ARMA parameters selected by MML87, AIC, BIC, and HQ—if q = 0, then we only use ARMA to forecast the time series data without LSTM. Our LSTM model is composed of

a single input layer with an input shape of MA order and the sequence learning features. The following LSTM layer also contains the sequence learning features, and the third LSTM layer with the same unit is followed by the fourth dense layer with one unit.

We developed Algorithm 1 based on [17] using a different loss function and activation function in the regression task. The hybrid ARMA-LSTM model trains the LSTM model by the residuals from the ARMA model. (This is similar in spirit to the discussion in ([7] (Section 5.1))). The simple point here is that the LSTM has at least the potential to find dependencies that the ARMA model (on its own) can not express. In this paper, MML87, AIC, BIC, and HQ have been used to select the model parameter orders from the ARMA model; so, this paper not only compares the errors of the hybrid ARMA-LSTM model with those from the single ARMA model but also the hybrid model in terms of the selection(s) of MML87, AIC, BIC, and HQ. The forecast from the ARMA model is the fitted mean μ_{t+1} . Because information is hidden in the residuals from the ARMA model (in a similar vein to ([7] (Section 5.1))), the forecast of the hybrid model will be

$$\hat{Y}_{t+1} = \mu_{t+1} + E_{t+1} \tag{12}$$

where μ_{t+1} represents the linearity modeling of data from the ARMA model selected according to the information-theoretic MML87, AIC, BIC, and HQ. The term ϵ_t is the residual left by the ARMA model $Y_t - \hat{Y}_t$, and $E_{t+1} = f(\epsilon_t) = f(Y_t - \hat{Y}_t)$, which is forecasted by the LSTM based on the past residual values $\epsilon_t, \epsilon_{t-1}, ..., \epsilon_{t-q}$, where the parameter *q* is selected by MML87, AIC, BIC, and HQ. The hybrid ARMA-LSTM model combines both linear and non-linear tendencies in time series data [48].

Algorithm 1 Algorithm 1 with the LSTM Model [17].
Require: number of epochs $= 10$
while MA(q) order in order set selected by MML, AIC, BIC, and HQ do
model.add(LSTM(30, return_sequences=True, input_shape= $(q, 1)$))

model.add(LSTM(30, return_sequences=True)) model.add(LSTM(30)) model.add(Dense(1))

The algorithm of the hybrid model is shown below (Algorithm 2):

Algorithm 2 Algorithm 2 with the Hybrid ARMA-LSTM Model.
Require: number of data $n \ge 0$
while $N \leq$ number of different simulations do
while $n \leq$ number of dataset in simulation do
while $i \in MA$ orders selected from MML, AIC, BIC, and HQ do
if $i \neq 0$ then
Train LSTM model by the residuals of ARMA model
Rolling forecast the residual by LSTM
Calculate root mean squared error by Y_{t+1}
else if $i = 0$ then
Calculate root mean squared error by forecast from ARMA only

6. Experiments

The experiments have been designed to compare the results of the ARMA model itself with the hybrid ARMA-LSTM model and also to compare different versions of the hybrid model with the parameters variously selected by the MML87, AIC, BIC, and HQ. In order to analyze the accuracy of forecasting, we are using the root mean squared error, RMSE = $\sqrt{\frac{1}{T}} \sum_{t=1}^{L} (y_t - \hat{y}_i)^2$, to compare the different results, where *T* stands for the forecast

window size, and we are using rolling forecast in this experiment. To elaborate and clarify,

for the financial data in Section 6.2, we do integer differencing with d = 1 to obtain stationarity before using the ARMA model and, as such, use an ARIMA or autoregressive integrated moving average model. We compare the performance of ARMA, ARIMA-LSTM, and LSTM alone on simulated dataset(s) (Section 6.1) and also on real-world financial (Section 6.2) and air pollution (Section 6.3) datasets.

We argue elsewhere (([6] footnotes 75 and 76) ([23] Section 3) ([38] Section 4.1)) about various uniqueness and invariance properties of log-loss (or logarithm loss). Squared error is a popular method and is also a variant of log-loss.

6.1. Simulated Dataset(s)

In this section, we perform experiments using various previously described modeling methods on simulated data, and we begin (in terms of LNPPP space ([6] Section 0.2.7)) by describing the experiments. We use a uniform distribution on [-0.9, 0.9] (from minimum -0.9 to maximum 0.9) to randomize the parameters p and q of ARMA(p,q) for the data simulation by using the arima.sim function in R and then reject them if they are outside the stationarity region. There are $5 \times 2 = 10$ different parameter sets from $p_1, ..., p_5$ and q_1, q_2 . The values in the table are the average RMSE over 100 runs (with standard deviation in brackets) in the simulated dataset corresponding to the particular parameters. The dataset includes N = 50, 100, 200, 300, and 500 time series data points in one dataset and also includes forecast windows of window size(s) T = 3, 10, 30, and 50. Table 1 shows the average of RMSE trained by LSTM alone (with different numbers of LSTM time steps) with different forecast window sizes, T. The results suggest that the LSTM alone does not work well in ARMA simulated data. For convenience of reading, we have moved Tables A1-A8 to the Appendix; each value in Table A1 is the average RMSE of forecast errors over the datasets (with standard deviation in brackets). The bold texts indicate the smallest forecast errors from the different kinds of models. Tables A1-A4 provide a comparison of different forecast window sizes (or window size) with T = 3, 10, 30, and 50.

No. of LSTM Time Steps	T = 3	T = 10	T = 30	T = 50
1	1.2519	1.3677	1.4962	1.3911
2	1.1794	1.2442	1.3863	1.2718
3	1.3372	1.6324	1.2256	1.3018
4	1.2195	1.2301	1.3284	1.3951
5	1.1341	1.6294	1.4276	1.4494

Table 1. RMSE in LSTM for simulated data (p_1, q_1) with different time steps and N = 100.

Table A2 shows the results for the average RMSE in the datasets for different simulated ARMA parameter sets, with the forecast window of T = 10. Table A3 provides the comparison of root mean squared error results of those datasets in different criteria, also comparing different simulated datasets with the forecast window of T = 30.

A large forecast window usually decreases the accuracy for the time series model. A window size of T = 50 (Data provided by Table A4) is 50% of the size of the in-sample set, and the MML87 hybrid model still outperforms its rivals. This indicates that the MML information criterion is efficient in model selection, and the algorithm of the hybrid model is also efficient in time series analysis, with the result of T = 50, as shown in Table A4. Table 2 shows the average of the ten different parameters of the simulated dataset in the forecast window sizes of T = 3 (Data provided by Table A1), 10 (Data provided by Table A2), 30 (Data provided by Table A3) and 50 (Data provided by Table A4) with the in-sample size of N = 100.

	Average of RMSE											
		AR	MA			ARMA-LSTM						
	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87				
T = 3	1.086	1.076	1.090	1.072	1.121	1.123	1.134	1.115				
T = 10	1.149	1.136	1.144	1.121	1.159	1.136	1.143	1.134				
T = 30	1.338	1.331	1.340	1.325	1.234	1.221	1.224	1.220				
T = 50	1.308	1.296	1.297	1.295	1.225	1.195	1.219	1.221				

Table 2. Average of RMSE in forecast window size T = 3, 10, 30, and 50.

MML87 outperforms the rival methods in the in-sample size of N = 100 in all cases of T = 3, 10, 30, and 50. MML87 not only considers the goodness of fit of data but also considers the model complexity. Figure 3 shows that MML87 has a lower root mean squared error in most cases. The hybrid model selected by MML87 has the lowest error rate for T = 3, 10, and 30. These comparisons argue well for MML. The results of N = 100 with T = 50 seem to suggest that for a large size of the forecast window, the complex hybrid ARMA-LSTM model seems to perform better than the simple time series model. Given that the simulated data were generated from an ARMA model, it is not immediately apparent why adding LSTM to produce a hybrid model should be advantageous in the case of larger datasets (although we would typically expect this if not dealing with data that are purely from an ARMA model). Table 3 shows the average of the ten different parameters of the simulated dataset in the in-sample size of N = 50 (Data provided by Table A5), 100 (Data provided by Table A2), 200 (Data provided by Table A6), 300 (Data provided by Table A7), and 500 (Data provided by Table A8).



Figure 3. Comparison in forecast window sizes T = 3, 10, 30, and 50.

	Average of RMSE (and Standard Deviation)											
		AR	MA		ARMA-LSTM							
-	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87				
N = 50	1.301	1.291	1.299	1.280	1.224	1.202	1.216	1.244				
N = 100	1.149	1.136	1.144	1.121	1.159	1.136	1.143	1.134				
N = 200	1.177	1.187	1.189	1.183	1.159	1.161	1.164	1.154				
N = 300	1.163	1.152	1.155	1.147	1.131	1.125	1.124	1.123				
N = 500	1.194	1.197	1.1984	1.196	1.180	1.173	1.179	1.181				

Table 3. Average of RMSE for in sample size N = 50, 100, 200, 300, and 500 with forecast window size T = 10.

Tables A5–A8 compare six different models or model selection techniques in the RMSE of the dataset in N = 50, 200, 300, and 500, with the forecast window size T = 10. AIC tends to overfit for small datasets, such as N = 50 (Data provided by Table A5 in the Appendix A). Through an increase in the amount for the in-sample dataset, the RMSE decreases in the hybrid ARMA-LSTM model because the larger size of data helps the LSTM to train and fit an accurate model. Thus, the results show the RMSE for the MML87 model is lower than the other models in the range N = 100, 200, and 300. Because of the efficiency in controlling the model complexity in MML87, the model can avoid the overfitting problem for small datasets.

The hybrid model with LSTM overfits when the in-sample size is small, basically because there is a larger amount of parameters that need to be estimated compared to the pure ARMA model. On the other hand, the hybrid model tends to perform well for a large in-sample size because the deep learning model is often better off for a large in-sample size, such as N = 200 (Data provided by Table A6), 300 (Data provided by Table A7), and 500 (Data provided by Table A8).

For a small in-sample size, such as N = 50, the BIC performance is good on the hybrid ARMA-LSTM because BIC is able to select the model well without overfitting. The MML87-Hybrid has the smallest average RMSE for N = 100, 200, and 300 for the different randomized datasets. The hybrid models work efficiently when there is enough in-sample data; otherwise, it can also overfit small datasets. In the meantime, by comparing the RMSE from MML87-ARMA, AIC-ARMA, BIC-ARMA, and HQ-ARMA, the results favor MML87 rather than AIC, BIC, and HQ. MML87 has a good performance in time series model selection and is able to select the ARMA model with lower forecasting errors. However, as noted earlier in this section, given that the simulated data were generated from an ARMA model, it is not immediately apparent why adding LSTM to produce a hybrid model should be advantageous in the case of larger datasets (although we would typically expect this if not dealing with data that are purely from an ARMA model). Figure 4 shows the comparison of RMSE in the in-sample size N = 50, 200, 300, and 500.



Figure 4. Comparison in the in-sample size N = 50, 200, 300, and 500.

6.2. Financial Data-and Extension to ARIMA Models

Stock return prediction is one of the most popular research topics in economics and finance [49,50]. This section studies the performance of the hybrid model from MML87; the hybrid models from AIC, BIC, HQ; and the ARIMA models selected by MML87, AIC, BIC, and HQ. The stock prices were selected from the components of the Dow Jones Industrial Average, including Apple (APPL), Boeing (BA), Cisco System (CSCO), Goldman Sachs (GS), IBM, Intel (INTC), Johnson & Johnson (JNJ), JPMorgan Chase (JPM), Coca-Cola (KO), and 3 M (MMM). The data selected start at 23 September 2016 and finish at 22 September 2021, with a total of 1258 trading days. This experiment studies the different performances in forecast window sizes T = 3, 5, 10, 30, 50, 70, 100, 130, 150, and 200. Table 4 shows the characteristic of stock prices selected, including mean, standard deviation, and partial autocorrelation.

Table 4. Mean, standard deviation, PACF lag 1 to 3 for ten selected stocks.

	Mean	S.D	PACF1	PACF2	PACF3
AAPL	66.440217	37.060808	0.996875	0.044454	-0.004848
BA	258.704781	82.478194	0.995870	-0.031231	-0.061804
CSCO	40.585947	8.595774	0.994585	0.073202	-0.016488
GS	227.095242	56.820929	0.993579	0.039741	-0.043412
IBM	124.851224	10.369478	0.982339	0.070195	-0.040622
INTC	46.269478	9.305502	0.992194	0.178757	-0.053398
JNJ	130.715314	18.399352	0.993930	0.050988	-0.031304
JPM	104.046116	24.467471	0.993854	0.067756	-0.049235
КО	44.519034	6.089778	0.993828	0.031639	-0.039178
MMM	173.550240	20.467854	0.991641	0.004475	0.026664

The empirical results show that the hybrid ARIMA-LSTM model can substantially outperform the traditional ARIMA (Autoregressive Integrated Moving Average) time series

model, particularly in the forecast window sizes of T = 5, 30, 100, 130, 150, and 200. Many studies demonstrated that the stock return depends on various factors, such as dividend yield, the book to market ratio, and/or interest rate [49,51,52]. However, traditional linear time series models are not able to take into account the effect of all those factors, thus requiring a more complex model to capture the information hidden in residuals from the ARIMA model. The hybrid model with LSTM is able to model publicly available and other information, which we have no reason to believe will be restricted, coming from a purely ARMA or ARIMA model. In order to make the stock price stationary in time series analysis, the ARIMA models are using the parameter d = 1 (or, equivalently, first-order differencing). As the experimental results show, MML87 outperforms the other information-theoretic criteria AIC, BIC, and HQ in terms of lower root mean squared error for out-of-sample forecasting. Figure 5 demonstrates the log prices for stock prices selected in this experiment.



Figure 5. Log prices for ten selected stocks.

The hybrid model tends to outperform for a large forecast window size rather than the small forecast window size because a large lookahead in forecasting has higher uncertainty. For much—or perhaps even most—of the financial industry, there is high volatility in long forecasts. The notion of semi-strong market efficiency suggests that the stock price fully and fairly reflects publicly available information in the time horizontal in the forecast window and also reflects all past information (although by no means all authors agree with this in its entirety [53], partly due to principles of Solomonoff [37] and Wallace [7]). Thus, it is more likely that a complex model will at least be able to provide accurate results in predictions for a T greater than 100. Table 5 shows MML models have lower the RMSE in most cases for different forecast window sizes in financial data.

	Average of RMSE (& Standard Deviation)												
		AR	IMA		ARIMA-LSTM								
	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87					
T = 3	2.987	3.027	2.914	3.075	4.414	4.302	4.375	4.289					
	(3.446)	(3.555)	(3.567)	(3.572)	(4.75)	(4.608)	(4.757)	(4.616)					
T = 5	4.024	4.077	4.126	4.163	4.024	3.966	4.081	3.907					
	(5.091)	(5.228)	(5.218)	(5.086)	(5.45)	(5.42)	(5.739)	(5.449)					
T = 10	4.748	4.747	4.712	4.868	5.359	5.261	5.272	5.249					
	(4.707)	(4.858)	(4.815)	(4.347)	(5.429)	(5.268)	(5.443)	(5.262)					
T = 30	5.872	5.867	5.91	5.994	5.754	5.628	5.726	5.643					
	(6.797)	(6.6)	(6.576)	(5.662)	(4.822)	(4.687)	(4.776)	(4.677)					
T = 50	7.834	7.609	7.726	6.659	7.328	7.411	7.405	7.384					
	(7.511)	(7.298)	(7.269)	(6.966)	(6.787)	(6.879)	(6.789)	(6.898)					
T = 70	9.991	9.909	10.024	9.645	10.393	10.221	10.42	10.085					
	(9.491)	(9.316)	(9.173)	(7.99)	(8.048)	(7.789)	(8.061)	(7.612)					
T = 100	14.465	13.991	14.197	9.866	9.304	9.087	9.235	9.253					
	(17.187)	(15.428)	(13.637)	(10.854)	(9.256)	(9.35)	(9.396)	(9.486)					
T = 130	14.482	14.301	17.672	13.551	13.768	13.811	13.9	14.581					
	(9.714)	(10.571)	(13.139)	(10.238)	(10.598)	(11.124)	(11.516)	(10.972)					
T = 150	22.985	22.985	23.021	18.045	17.778	17.526	17.98	17.461					
	(28.173)	(28.077)	(28.071)	(17.856)	(16.771)	(16.582)	(16.734)	(15.931)					
T = 200	31.144	30.502	30.712	30.286	26.831	26.424	26.662	26.507					
	(37.567)	(38.314)	(38.322)	(32.564)	(31.63)	(31.547)	(31.645)	(31.59)					

Table 5. RMSE for forecast window sizes T = 3, 5, 10, 30, 50, 70, 100, 130, 150, and 200.

Table 6 provides the average of RMSE for the selected stocks in different sizes, *T*, of the forecast window (shown in different columns) and numbers of LSTM time steps (shown in different rows). The LSTM models are trained by scalers in the range of 0 to 1, and the LSTM model performs worse in the case without scaling, which indicates that the neural network LSTM is scale insensitive and that combining the traditional ARMA time series model makes the neural network more scale-sensitive [54]. The results from Table 6 suggest that the LSTM model alone (unenhanced by ARMA and ARIMA) is not particularly able to capture the time series pattern for the stock price. The figures of the average RMSE are significantly higher than traditional ARMA and ARIMA-LSTM models. Figure 6 shows the comparison between the ARIMA model and the hybrid ARIMA-LSTM model in this experiment.

No. Steps	T = 3	T = 5	T = 30	T = 10	T = 50	T = 70	T = 100	T = 130	T = 150	T = 200
1	8.5789	10.1965	56.7817	104.3681	123.4805	119.2673	151.1338	107.2951	114.8106	73.2335
3	5.7604	3.5166	3.6097	13.325	10.6368	31.9361	33.4419	26.0112	31.5578	26.6354
5	4.0695	3.0575	8.5064	11.9009	15.5075	17.3077	19.0942	48.0012	30.0622	36.6099
7	3.9708	6.4145	10.6368	6.8547	13.2163	16.5474	19.0724	32.7076	20.5954	44.1875
10	5.3985	6.4576	5.9597	13.8295	16.0972	20.6271	12.8859	28.2251	28.2803	25.5409

Table 6. LSTM with different time steps for financial data in varying forecast windows.



Figure 6. Comparison in different forecast windows.

6.3. PM2.5 Pollution Data

In this section, we use environmental data of PM2.5 pollution levels in the city of Beijing, China, with ten sensors located in different areas. The data are hourly PM2.5 levels in 53 days in 2013.

We are using the same data length and information-theoretic methods from Section 6.2 in order to demonstrate the performance of MML compared to rival methods. Table 7 shows the comparison between MML, AIC, BIC, and HQ. The hourly PM2.5 data have a seasonality; the level of PM2.5 reaches its highest near midday and decreases to its lowest near midnight. The results suggest that MML is a good model selection technique in this case.

Table 7. RMSE for forecast window sizes T = 3, 5, 10, 30, 50, 70, 100, 130, 150, and 200.

	Average of RMSE & Standard Deviation												
		AR	IMA		ARIMA-LSTM								
	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87					
T = 3	26.805	26.689	26.569	23.104	25.768	23.066	24.791	22.965					
	(7.496)	(7.532)	(7.545)	(7.843)	(7.833)	(8.693)	(7.883)	(6.711)					
T = 5	28.036	27.538	27.479	23.478	24.636	22.309	24.113	21.666					
	(6.986)	(6.805)	(7.186)	(8.426)	(8.518)	(7.596)	(8.584)	(7.424)					
T = 10	30.633	31.502	31.585	30.074	26.970	27.566	27.924	25.102					
	(12.679)	(12.283)	(12.518)	(14.917)	(10.502)	(13.487)	(10.107)	(9.855)					
T = 30	40.730	40.788	40.157	37.989	31.022	29.382	31.689	28.572					
	(14.001)	(13.372)	(14.195)	(19.180)	(11.409)	(12.196)	(14.124)	(12.229)					
T = 50	39.097	38.662	39.007	42.986	35.639	33.335	36.036	40.568					
	(4.238)	(4.660)	(4.232)	(6.062)	(5.599)	(5.184)	(6.339)	(9.24)					
T = 70	34.004	33.551	34.773	32.030	48.942	45.305	49.723	42.987					
	(4.223)	(4.105)	(3.514)	(9.404)	(12.377)	(10.567)	(12.068)	(8.705)					
T = 100	32.002	32.444	31.170	37.925	56.024	51.543	59.705	49.513					
	(2.434)	(2.425)	(2.865)	(4.444)	(13.199)	(12.714)	(13.435)	(11.45)					
T = 130	44.023	44.162	43.635	43.802	36.183	33.716	39.401	46.496					
	(2.583)	(2.576)	(2.836)	(1.853)	(7.184)	(4.591)	(8.488)	(7.168)					
T = 150	44.463	44.736	43.928	41.150	32.574	31.225	30.923	33.584					
	(1.612)	(1.862)	(1.773)	(5.221)	(6.211)	(4.301)	(6.598)	(7.679)					
T = 200	42.150	42.372	41.863	43.75	46.711	43.721	53.393	46.363					
	(2.620)	(2.522)	(2.787)	(4.07)	(13.86)	(11.349)	(12.458)	(12.472)					

Table 8 shows the LSTM model alone in the PM2.5 data, and the results suggest that the LSTM model (on its own, unenhanced by ARMA and ARIMA) outperforms in the smaller-sized forecast windows, such as T = 3, 5, and 10. The RMSEs in larger window sizes ($T \ge 50$) are much larger for the LSTM than for the ARMA model and hybrid ARMA-LSTM.

No. Steps	T = 3	T = 5	T = 10	T = 30	T = 50	T = 70	T = 100	T = 130	T = 150	T = 200
1	3.0976	5.7806	16.5048	47.8431	53.7436	67.2412	81.7044	92.6897	73.4192	71.5536
3	4.4983	8.1565	17.0462	34.0492	36.3896	47.2558	64.68533	90.7986	78.9972	78.5648
5	4.7719	9.2208	18.7955	33.6065	50.4786	56.7465	59.3666	75.4321	102.0098	88.1695
7	5.9126	9.8355	15.3696	25.8551	38.6874	53.06845	50.962	87.998	92.101	101.2337
10	8.4289	11.4749	11.4479	38.3303	44.5138	65.6299	70.6415	74.6879	90.1211	84.0196

Table 8. LSTM for PM2.5 Beijing data in different time steps and forecast windows.

7. Conclusions

We have investigated time series modeling in the Minimum Message Length framework using Wallace and Freeman's (1987) approximation [9]. The hybrid ARMA-LSTM model has been compared with the traditional ARMA (Autoregressive Moving Average) time series model based on the information-theoretic approaches: AIC, BIC, HQ and MML87. We performed experiments on simulated data and also on two real-world datasets (financial and environmental data). We conducted the experiments based on hybrid ARMA-LSTM (with LSTM) and ARMA without LSTM (long short-term memory). This could be broadly thought of as constituting two experiments each on three datasets or with six experiments. For each of the six experiments, the results show that MML87 outperforms the other information-theoretic criteria. The hybrid ARMA-LSTM model performs better than the traditional ARMA model, and the MML hybrid ARMA-LSTM model performed best out of everything considered. It is worth noting that the LSTM model alone with unscaled data performed worse than everything else considered. In summary, MML87 is able to select the lower forecasting errors better than the AIC, BIC, and HQ, as the experimental results show.

Author Contributions: Conceptualization: Z.F.; methodology: Z.F. and D.L.D.; computation: Z.F. and D.R.; validation: Z.F., D.L.D. and S.P.; investigation: Z.F., S.P. and D.L.D.; writing and preparation: Z.F. and D.L.D.; writing and review: Z.F., D.L.D. and S.P.; supervision, D.L.D., S.P. and D.R. All authors have read and agreed to the published version of the manuscript and have endeavored to make the work as error-free as possible.

Funding: This research received no external funding.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: Not applicable.

Acknowledgments: The authors thank three anonymous referees for their valuable comments and useful suggestions to improve the quality of this version of the paper. The (first two) authors would further like to thank the Department of Data Science and Artificial Intelligence, Faculty of IT, Monash University, for their support.

Conflicts of Interest: The authors declare no conflict of interest.

16 of 21

Appendix A

	Average of RMSE (and Standard Deviation)											
Order of Stat-		AF	MA			ARMA	-LSTM					
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87				
<i>p</i> ₁ , <i>q</i> ₁	0.982	1.108	1.112	1.033	1.204	1.217	1.215	1.234				
	(0.646)	(0.529)	(0.499)	(0.469)	(0.308)	(0.502)	(0.511)	(0.7)				
<i>p</i> ₁ , <i>q</i> ₂	1.133	1.053	1.172	1.166	1.301	1.289	1.366	1.384				
	(0.592)	(0.669)	(0.601)	(0.635)	(0.922)	(0.95)	(0.862)	(0.838)				
<i>p</i> ₂ , <i>q</i> ₁	1.027	1.024	1.029	1.023	1.025	1.012	1.021	1.005				
	(0.423)	(0.421)	(0.445)	(0.418)	(0.408)	(0.376)	(0.411)	(0.48)				
<i>p</i> ₂ , <i>q</i> ₂	1.333	1.278	1.286	1.271	1.241	1.182	1.211	1.194				
	(0.793)	(0.841)	(0.879)	(0.848)	(0.745)	(0.711)	(0.735)	(0.674)				
<i>p</i> ₃ , <i>q</i> ₁	0.955	0.956	0.951	0.944	0.965	0.975	0.971	0.986				
	(0.377)	(0.377)	(0.375)	(0.37)	(0.341)	(0.35)	(0.351)	(0.426)				
<i>p</i> ₃ , <i>q</i> ₂	1.293	1.241	1.245	1.238	1.114	1.211	1.172	1.105				
	(0.331)	(0.296)	(0.307)	(0.296)	(0.284)	(0.266)	(0.269)	(0.259)				
<i>p</i> ₄ , <i>q</i> ₁	0.901	0.916	0.913	0.871	0.948	0.944	0.945	0.932				
	(0.483)	(0.448)	(0.451)	(0.398)	(0.397)	(0.41)	(0.413)	(0.442)				
<i>p</i> ₄ , <i>q</i> ₂	1.207	1.226	1.224	1.206	1.252	1.261	1.257	1.251				
	(0.539)	(0.515)	(0.531)	(0.513)	(0.777)	(0.778)	(0.792)	(0.772)				
<i>p</i> ₅ , <i>q</i> ₁	1.006	0.907	0.911	0.903	1.122	1.117	1.112	1.018				
	(0.54)	(0.626)	(0.642)	(0.578)	(0.538)	(0.553)	(0.564)	(0.467)				
<i>p</i> ₅ , <i>q</i> ₂	1.026	1.052	1.054	1.061	1.042	1.021	1.027	1.046				
	(0.583)	(0.553)	(0.587)	(0.559)	(0.559)	(0.592)	(0.566)	(0.53)				

Table A1. Simulated data with N = 100 and T = 3 from Section 6.1.

Table A2. Simulated data with N = 100 and T = 10 from Section 6.1.

	Average of RMSE (and Standard Deviation)										
Order of Stat-		AR	MA		ARMA-LSTM						
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87			
<i>p</i> ₁ , <i>q</i> ₁	1.234	1.208	1.206	1.221	1.201	1.132	1.135	1.138			
	(0.178)	(0.165)	(0.167)	(0.293)	(0.404)	(0.184)	(0.188)	(0.317)			
<i>p</i> ₁ , <i>q</i> ₂	1.571	1.553	1.556	1.398	1.555	1.549	1.551	1.494			
	(0.375)	(0.386)	(0.391)	(0.304)	(1.109)	(1.125)	(1.123)	(0.834)			
<i>p</i> ₂ , <i>q</i> ₁	1.025	1.041	1.044	1.043	1.013	1.02	1.025	1.037			
	(0.194)	(0.203)	(0.216)	(0.193)	(0.174)	(0.182)	(0.174)	(0.265)			
<i>p</i> ₂ , <i>q</i> ₂	1.353	1.327	1.322	1.325	1.274	1.257	1.268	1.255			
	(0.438)	(0.373)	(0.391)	(0.368)	(0.213)	(0.206)	(0.211)	(0.205)			
<i>p</i> ₃ , <i>q</i> ₁	0.947	0.895	0.918	0.901	1.018	0.946	0.966	0.989			
	(0.194)	(0.129)	(0.157)	(0.134)	(0.135)	(0.116)	(0.151)	(0.154)			
<i>p</i> ₃ , <i>q</i> ₂	0.978	1.06	1.065	1.048	1.149	1.137	1.141	1.135			
	(0.266)	(0.239)	(0.225)	(0.226)	(0.328)	(0.338)	(0.352)	(0.328)			
<i>p</i> ₄ , <i>q</i> ₁	1.083	1.059	1.063	1.075	1.081	1.029	1.035	1.061			
	(0.206)	(0.2)	(0.218)	(0.179)	(0.261)	(0.179)	(0.219)	(0.128)			
<i>p</i> ₄ , <i>q</i> ₂	1.121	1.112	1.124	1.104	1.093	1.088	1.095	1.096			
	(0.192)	(0.17)	(0.186)	(0.174)	(0.212)	(0.191)	(0.252)	(0.181)			
<i>p</i> 5, <i>q</i> 1	1.279	1.244	1.264	1.242	1.169	1.167	1.172	1.166			
	(0.322)	(0.296)	(0.251)	(0.29)	(0.335)	(0.327)	(0.335)	(0.306)			
<i>p</i> ₅ , <i>q</i> ₂	0.903	0.867	0.882	0.877	1.053	1.033	1.046	0.972			
	(0.078)	(0.067)	(0.059)	(0.074)	(0.231)	(0.192)	(0.188)	(0.126)			

Average of RMSE (and Standard Deviation)									
Order of Stat-		AR	MA		ARMA-LSTM				
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87	
p_1, q_1	1.263	1.252	1.253	1.256	1.217	1.118	1.125	1.192	
	(0.167)	(0.156)	(0.173)	(0.159)	(0.295)	(0.119)	(0.133)	(0.247)	
p_1, q_2	2.641	2.554	2.631	2.694	1.771	1.848	1.822	1.803	
	(0.905)	(0.838)	(0.972)	(0.961)	(1.135)	(0.739)	(0.959)	(1.373)	
p_2, q_1	1.221	1.186	1.199	1.184	1.102	1.088	1.094	1.124	
	(0.139)	(0.096)	(0.121)	(0.102)	(0.084)	(0.083)	(0.089)	(0.101)	
<i>p</i> ₂ , <i>q</i> ₂	1.044	1.145	1.093	1.041	1.138	1.153	1.148	1.136	
	(0.091)	(0.108)	(0.117)	(0.088)	(0.255)	(0.211)	(0.262)	(0.256)	
<i>p</i> ₃ , <i>q</i> ₁	1.086	1.066	1.073	1.061	1.038	1.036	1.036	1.035	
	(0.181)	(0.19)	(0.195)	(0.182)	(0.172)	(0.171)	(0.188)	(0.145)	
<i>p</i> ₃ , <i>q</i> ₂	1.112	1.096	1.139	1.101	1.202	1.153	1.166	1.099	
	(0.295)	(0.309)	(0.331)	(0.306)	(0.38)	(0.328)	(0.369)	(0.264)	
p_4, q_1	1.053	1.038	1.044	1.035	1.058	1.051	1.055	1.063	
	(0.22)	(0.189)	(0.167)	(0.185)	(0.14)	(0.124)	(0.139)	(0.152)	
p_4, q_2	1.263	1.247	1.251	1.238	1.204	1.191	1.211	1.183	
	(0.2)	(0.194)	(0.229)	(0.21)	(0.133)	(0.114)	(0.138)	(0.152)	
p_5, q_1	1.613	1.679	1.669	1.599	1.541	1.531	1.539	1.521	
	(0.27)	(0.301)	(0.343)	(0.342)	(0.884)	(0.609)	(0.915)	(0.848)	
p_5, q_2	1.092	1.047	1.052	1.047	1.074	1.041	1.045	1.041	
	(0.132)	(0.234)	(0.337)	(0.114)	(0.144)	(0.117)	(0.196)	(0.115)	

Table A3. Simulated data with N = 100 and T = 30 from Section 6.1.

Table A4. Simulated data with N = 100 and T = 50 from Section 6.1.

	Average of RMSE (and Standard Deviation)								
Order of Stat-	ARMA				ARMA-LSTM				
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87	
<i>p</i> ₁ , <i>q</i> ₁	1.189	1.191	1.193	1.182	1.164	1.091	1.155	1.173	
	(0.217)	(0.228)	(0.221)	(0.222)	(0.304)	(0.212)	(0.292)	(0.241)	
<i>p</i> ₁ , <i>q</i> ₂	2.307	2.308	2.305	2.298	1.862	1.868	1.889	1.852	
	(0.458)	(0.457)	(0.464)	(0.466)	(1.169)	(1.16)	(1.171)	(1.073)	
<i>p</i> ₂ , <i>q</i> ₁	1.113	1.092	1.095	1.094	1.058	1.045	1.172	1.059	
	(0.087)	(0.103)	(0.107)	(0.104)	(0.073)	(0.096)	(0.225)	(0.139)	
<i>p</i> ₂ , <i>q</i> ₂	1.191	1.189	1.192	1.191	1.176	1.178	1.183	1.201	
	(0.096)	(0.103)	(0.107)	(0.1)	(0.24)	(0.259)	(0.285)	(0.289)	
<i>p</i> ₃ , <i>q</i> ₁	1.094	1.093	1.095	1.097	1.101	1.061	1.065	1.093	
	(0.159)	(0.157)	(0.156)	(0.155)	(0.192)	(0.144)	(0.177)	(0.115)	
p3,q2	1.127	1.123	1.129	1.125	1.121	1.129	1.126	1.132	
	(0.06)	(0.055)	(0.057)	(0.058)	(0.134)	(0.155)	(0.143)	(0.153)	
p_4, q_1	1.188	1.189	1.185	1.192	1.136	1.095	1.099	1.139	
	(0.182)	(0.188)	(0.187)	(0.186)	(0.137)	(0.181)	(0.173)	(0.113)	
<i>p</i> ₄ , <i>q</i> ₂	1.232	1.221	1.222	1.212	1.268	1.19	1.197	1.203	
	(0.165)	(0.133)	(0.137)	(0.134)	(0.457)	(0.269)	(0.234)	(0.319)	
p_5, q_1	1.593	1.521	1.533	1.528	1.331	1.275	1.277	1.338	
	(0.304)	(0.199)	(0.216)	(0.209)	(0.428)	(0.234)	(0.214)	(0.383)	
<i>p</i> ₅ , <i>q</i> ₂	1.051	1.033	1.029	1.032	1.035	1.021	1.029	1.023	
	(0.083)	(0.064)	(0.096)	(0.063)	(0.055)	(0.067)	(0.077)	(0.069)	

Average of RMSE (and Standard Deviation)										
Order of Stat-		AR	MA		ARMA-LSTM					
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87		
p_1, q_1	1.068	1.071	1.073	1.067	1.198	1.122	1.137	1.175		
	(0.147)	(0.118)	(0.125)	(0.115)	(0.222)	(0.305)	(0.336)	(0.259)		
p_1, q_2	1.994	1.994	2.056	2.04	1.93	1.932	1.926	1.921		
	(0.655)	(0.655)	(0.692)	(0.705)	(1.553)	(1.563)	(1.572)	(1.566)		
p_2, q_1	1.242	1.242	1.274	1.235	1.106	1.116	1.119	1.154		
	(0.213)	(0.213)	(0.252)	(0.17)	(0.193)	(0.196)	(0.189)	(0.278)		
<i>p</i> ₂ , <i>q</i> ₂	1.185	1.183	1.196	1.232	1.163	1.194	1.172	1.254		
	(0.355)	(0.359)	(0.361)	(0.476)	(0.386)	(0.499)	(0.534)	(0.601)		
p_3, q_1	1.348	1.254	1.269	1.304	1.257	1.139	1.212	1.256		
	(0.557)	(0.604)	(0.661)	(0.575)	(0.499)	(0.605)	(0.657)	(0.449)		
p3,q2	1.283	1.283	1.281	1.291	1.198	1.198	1.211	1.215		
	(0.234)	(0.234)	(0.265)	(0.233)	(0.27)	(0.27)	(0.298)	(0.285)		
p_4, q_1	1.263	1.251	1.288	1.044	1.091	1.079	1.096	1.129		
	(0.461)	(0.469)	(0.477)	(0.172)	(0.264)	(0.27)	(0.288)	(0.243)		
p_4, q_2	0.987	0.987	0.989	0.989	1.007	1.017	1.022	0.999		
	(0.132)	(0.132)	(0.132)	(0.137)	(0.137)	(0.126)	(0.139)	(0.138)		
p_5, q_1	1.533 (0.457)	1.426 (0.535)	1.454 (0.561)	1.464 (0.509)	1.227 (0.442)	1.178 (0.445)	1.192	1.254 (0.434)		

1.137

(0.185)

1.061

(0.168)

1.068

(0.175)

1.072

(0.183)

1.08

(0.117)

Table A5. Simulated data with N = 50 and T = 10 from Section 6.1.

Table A6. Simulated data with N = 200 and T = 10 from Section 6.1.

1.111

(0.186)

1.098

(0.151)

1.101

(0.153)

 p_{5}, q_{2}

	Average of RMSE (and Standard Deviation)									
Order of Stat-		AR	MA	ARMA-LSTM						
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87		
p_1, q_1	1.244	1.277	1.286	1.248	1.153	1.13	1.146	1.151		
	(0.365)	(0.42)	(0.417)	(0.404)	(0.381)	(0.376)	(0.392)	(0.353)		
p_1, q_2	1.359	1.359	1.366	1.359	1.474	1.491	1.477	1.474		
	(0.445)	(0.445)	(0.462)	(0.445)	(0.813)	(0.882)	(0.893)	(0.813)		
p_2, q_1	0.927	0.915	0.916	0.92	0.939	0.955	0.969	0.933		
	(0.183)	(0.172)	(0.185)	(0.182)	(0.126)	(0.15)	(0.163)	(0.128)		
<i>p</i> ₂ , <i>q</i> ₂	1.184	1.191	1.193	1.189	1.134	1.114	1.116	1.106		
	(0.41)	(0.398)	(0.366)	(0.402)	(0.368)	(0.393)	(0.407)	(0.37)		
p_3, q_1	1.137	1.136	1.129	1.117	1.082	1.082	1.088	1.085		
	(0.347)	(0.347)	(0.351)	(0.355)	(0.314)	(0.316)	(0.361)	(0.325)		
p3,q2	0.915	1.038	1.011	0.991	1.088	1.083	1.075	1.054		
	(0.198)	(0.08)	(0.081)	(0.093)	(0.184)	(0.172)	(0.199)	(0.161)		
p_4, q_1	1.199	1.166	1.174	1.19	1.086	1.109	1.115	1.107		
	(0.558)	(0.557)	(0.531)	(0.562)	(0.591)	(0.507)	(0.691)	(0.732)		
p_4, q_2	1.108	1.101	1.132	1.129	1.184	1.186	1.192	1.184		
	(0.196)	(0.191)	(0.615)	(0.24)	(0.358)	(0.359)	(0.379)	(0.36)		
p_5, q_1	1.581	1.584	1.584	1.586	1.383	1.391	1.396	1.382		
	(0.481)	(0.475)	(0.422)	(0.48)	(0.83)	(0.802)	(0.811)	(0.832)		
p5,q2	1.123	1.101	1.107	1.101	1.063	1.069	1.065	1.063		
	(0.263)	(0.174)	(0.155)	(0.174)	(0.234)	(0.133)	(0.129)	(0.128)		

Average of RMSE (and Standard Deviation)									
Order of Stat-	ARMA				ARMA-LSTM				
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87	
p_1, q_1	1.024	1.028	1.029	1.031	1.033	1.02	1.021	1.024	
	(0.312)	(0.332)	(0.321)	(0.322)	(0.316)	(0.27)	(0.291)	(0.32)	
p_1, q_2	2.008	1.995	1.996	1.988	1.709	1.72	1.725	1.72	
	(1.123)	(1.024)	(1.031)	(1.028)	(0.918)	(0.896)	(0.812)	(0.854)	
p_2, q_1	1.022	1.025	1.027	1.016	1.011	1.012	1.017	1.014	
	(0.144)	(0.138)	(0.132)	(0.133)	(0.125)	(0.121)	(0.126)	(0.297)	
p ₂ ,q ₂	1.172	1.168	1.171	1.166	1.164	1.177	1.172	1.17	
	(0.398)	(0.383)	(0.326)	(0.384)	(0.422)	(0.443)	(0.461)	(0.413)	
p_3, q_1	0.886	0.868	0.882	0.865	0.964	0.932	0.952	0.914	
	(0.198)	(0.205)	(0.217)	(0.215)	(0.261)	(0.183)	(0.191)	(0.188)	
p3,q2	1.07	1.068	1.065	1.059	1.096	1.095	1.097	1.092	
	(0.408)	(0.412)	(0.407)	(0.401)	(0.284)	(0.289)	(0.277)	(0.284)	
p_4, q_1	1.215	1.191	1.194	1.184	1.22	1.091	1.124	1.166	
	(0.445)	(0.468)	(0.462)	(0.464)	(0.621)	(0.42)	(0.468)	(0.453)	
p_4, q_2	1.191	1.167	1.172	1.162	1.182	1.188	1.184	1.184	
	(0.338)	(0.308)	(0.311)	(0.278)	(0.427)	(0.473)	(0.113)	(0.433)	
p_5, q_1	1.169	1.159	1.161	1.152	0.997	1.071	1.011	1.01	
	(0.225)	(0.216)	(0.232)	(0.216)	(0.131)	(0.213)	(0.159)	(0.146)	
<i>p</i> ₅ , <i>q</i> ₂	0.874	0.846	0.852	0.844	0.936	0.939	0.935	0.938	
	(0.25)	(0.249)	(0.297)	(0.247)	(0.213)	(0.197)	(0.199)	(0.196)	

Table A7. Simulated data with N = 300 and T = 10 from Section 6.1.

Table A8. Simulated data with N = 500 & T = 10 from Section 6.1.

	Average of RMSE (and Standard Deviation)									
Order of Stat-		AR	MA		ARMA-LSTM					
ionary ARMA	AIC	BIC	HQ	MML87	AIC	BIC	HQ	MML87		
p_1, q_1	0.988	0.966	0.967	0.968	1.016	1.012	1.017	1.014		
	(0.229)	(0.233)	(0.252)	(0.232)	(0.178)	(0.182)	(0.169)	(0.179)		
p_1, q_2	1.546	1.549	1.552	1.562	1.841	1.838	1.838	1.838		
	(0.728)	(0.713)	(0.736)	(0.703)	(0.915)	(0.875)	(0.876)	(0.877)		
p_2, q_1	1.002	1.017	1.017	1.016	1.008	1.008	1.008	1.05		
	(0.37)	(0.349)	(0.355)	(0.351)	(0.329)	(0.325)	(0.334)	(0.349)		
<i>p</i> ₂ , <i>q</i> ₂	1.156	1.165	1.163	1.165	1.167	1.156	1.159	1.156		
	(0.188)	(0.176)	(0.182)	(0.176)	(0.337)	(0.355)	(0.363)	(0.355)		
<i>p</i> ₃ , <i>q</i> ₁	1.091	1.093	1.099	1.09	1.064	1.06	1.066	1.058		
	(0.175)	(0.18)	(0.175)	(0.176)	(0.225)	(0.157)	(0.173)	(0.22)		
p3,q2	1.23	1.235	1.235	1.235	1.197	1.209	1.21	1.209		
	(0.372)	(0.364)	(0.364)	(0.364)	(0.365)	(0.393)	(0.378)	(0.393)		
p_4, q_1	1.041	1.07	1.069	1.063	1.135	1.053	1.096	1.139		
	(0.272)	(0.25)	(0.261)	(0.257)	(0.342)	(0.239)	(0.298)	(0.343)		
<i>p</i> ₄ , <i>q</i> ₂	1.253	1.256	1.261	1.255	1.134	1.136	1.131	1.126		
	(0.265)	(0.265)	(0.273)	(0.266)	(0.218)	(0.214)	(0.243)	(0.235)		
p_5, q_1	1.559	1.551	1.55	1.541	1.159	1.199	1.196	1.161		
	(0.363)	(0.385)	(0.374)	(0.365)	(0.331)	(0.421)	(0.411)	(0.292)		
<i>p</i> ₅ , <i>q</i> ₂	1.073	1.068	1.071	1.068	1.083	1.062	1.067	1.062		
	(0.179)	(0.188)	(0.192)	(0.188)	(0.136)	(0.167)	(0.166)	(0.167)		

References

- Siami-Namini, S.; Tavakoli, N.; Namin, A.S. A comparison of ARIMA and LSTM in forecasting time series. In Proceedings of the 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), Orlando, FL, USA, 17–20 December 2018; pp. 1394–1401.
- 2. Wallace, C.S.; Boulton, D.M. An information measure for classification. Comput. J. 1968, 11, 185–194.
- 3. Akaike, H. A new look at the statistical model identification. IEEE Trans. Autom. Control 1974, 19, 716–723.
- 4. Schwarz, G. Estimating the dimension of a model. Ann. Stat. 1978, 6, 461-464.
- 5. Hannan, E.J.; Quinn, B.G. The determination of the order of an autoregression. J. R. Stat. Soc. Ser. B Methodol. 1979, 41, 190–195.
- 6. Dowe, D.L. Foreword re C. S. Wallace. Comput. J. 2008, 51, 523–560.
- 7. Wallace, C.S.; Dowe, D.L. Minimum message length and Kolmogorov complexity. Comput. J. 1999, 42, 270–283.
- 8. Wong, C.K.; Makalic, E.; Schmidt, D.F. Minimum message length inference of the Poisson and geometric models using heavy-tailed prior distributions. *J. Math. Psychol.* **2018**, *83*, 1–11.
- 9. Wallace, C.S.; Freeman, P.R. Estimation and inference by compact coding. J. R. Stat. Soc. Ser. B Methodol. 1987, 49, 240–252.
- 10. Fang, Z.; Dowe, D.L.; Peiris, S.; Rosadi, D. Minimum Message Length Autoregressive Moving Average Model Order Selection. *arXiv* 2021, arXiv:2110.03250.
- 11. Schmidt, D.F. Minimum Message Length Inference of Autoregressive Moving Average Models. Ph.D. Thesis, Faculty of IT, Monash University, Melbourne, Australia, 2008.
- 12. Fathi, O. Time Series Forecasting Using a Hybrid ARIMA and LSTM Model; Velvet Consulting: Paris, France, 2019.
- 13. Box, G.E.; Jenkins, G.M.; Reinsel, G.C. *Time Series Analysis Prediction and Control*; John Wiley and Sons: Hoboken, NJ, USA, 1976.
- 14. De Gooijer, J.G.; Hyndman, R.J. 25 years of time series forecasting. Int. J. Forecast. 2006, 22, 443-473.
- 15. Box, G.E.; Jenkins, G.M.; Reinsel, G.C.; Ljung, G.M. *Time Series Analysis: Forecasting and Control;* John Wiley & Sons: New York, NY, USA, 2015.
- 16. Wang, J.Q.; Du, Y.; Wang, J. LSTM based long-term energy consumption prediction with periodicity. Energy 2020, 197, 117197.
- Chen, K.; Zhou, Y.; Dai, F. A LSTM-based method for stock returns prediction: A case study of China stock market. In Proceedings of the 2015 IEEE International Conference on Big Data (Big Data), Santa Clara, CA, USA, 29 October–1 November, 2015; pp. 2823–2824.
 Sale M.; Davie D.L.; Bay S. Minimum measure length maying sugress time series data mining. In Proceedings of the In 2005 ICEC.
- Sak, M.; Dowe, D.L.; Ray, S. Minimum message length moving average time series data mining. In Proceedings of the In 2005 ICSC Congress on Computational Intelligence Methods and Applications, Istanbul, Turkey, 15–17 December, 2005; 6p.
- 19. Wallace, C.S. Statistical and Inductive Inference by Minimum Message Length; Springer: New York, NY, USA, 2005; pp. 93–100.
- 20. Aho, K.; Derryberry, D.; Peterson, T. Model selection for ecologists: The worldviews of AIC and BIC. Ecology 2014, 95, 631-636.
- 21. Grasa, A.A. Econometric Model Selection: A New Approach; Springer Science & Business Media: New York, NY, USA, 2013; Volume 16.
- 22. Hernandez-Matamoros, A.; Fujita, H.; Hayashi, T.; Perez-Meana, H. Forecasting of COVID19 per regions using ARIMA models and polynomial functions. *Appl. Soft Comput.* **2020**, *96*, 106610.
- Dowe, D.L. MML, hybrid Bayesian network graphical models, statistical consistency, invariance and uniqueness. In *Handbook of the Philosophy of Science*; Volume 7: Philosophy of Statistics; Elsevier: New York, NY, USA, 2011; pp. 901–982.
- 24. Baxter, R.A.; Dowe, D.L. Model selection in linear regression using the MML criterion. In Proceedings of the Data Compression Conference, IEEE, Institute of Electrical and Electronics Engineers, Snowbird, UT, USA, 29–31 March 1994.
- 25. Fitzgibbon, L.J.; Dowe, D.L.; Vahid, F. Minimum message length autoregressive model order selection. In Proceedings of the International Conference on Intelligent Sensing and Information Processing, Chennai, India, 4–7 January 2004; pp. 439–444.
- Schmidt, D.F. Minimum message length order selection and parameter estimation of moving average models. In *Algorithmic Probability and Friends*; Bayesian Prediction and Artificial Intelligence; Springer: Berlin/Heidelberg, Germany, 2013; pp. 327–338.
- Wallace, C.S.; Dowe, D.L. Intrinsic classification by MML-the Snob program. In Proceedings of the 7th Australian Joint Conference on Artificial Intelligence World Scientific, Armidale, Australia, 1 January 1994; pp. 37–44.
- Wallace, C.S.; Dowe, D.L. MML clustering of multi-state, Poisson, von Mises circular and Gaussian distributions. *Stat. Comput.* 2000, 10, 73–83.
- 29. Dowe, D.L.; Allison, L.; Dix, T.I.; Hunter, L.; Wallace, C.S.; Edgoose, T. Circular clustering of protein dihedral angles by minimum message length. In *Pacific Symposium on Biocomputing*; World Scientific: Singapore, 1996; pp. 242–255.
- Oliver, J.J.; Dowe, D.L.; Wallace, C.S. Inferring decision graphs using the minimum message length principle. In Proceedings of the 5th Australian Joint Conference on Artificial Intelligence, Hobart, NSW, Australia, 16–18 November 1992; pp. 361–367.
- 31. Tan, P.J.; Dowe, D.L. MML inference of decision graphs with multi-way joins and dynamic attributes. In *Australasian Joint Conference* on *Artificial Intelligence*; Springer: Berlin/Heidelberg, Germany, 2003; pp. 269–281.
- 32. Comley, J.W.; Dowe, D.L. General Bayesian networks and asymmetric languages. In Proceedings of the 2nd Hawaii International Conference on Statistics and Related Fields, Honolulu, HI, USA, 5–8 June 2003; p. 18.
- Comley, J.W.; Dowe, D.L. Minimum Message Length and Generalized Bayesian Nets with Asymmetric Languages. *Minimum* 2005, 265, 265–294.
- Saikrishna, V.; Dowe, D.L.; Ray, S. MML learning and inference of hierarchical Probabilistic Finite State Machines. In *Applied Data Analytics: Principles and Applications*; River Publishers: Aalborg, Denmark, 2020; pp. 291–325.
- 35. Dowe, D.L.; Zaidi, N.A. Database normalization as a by-product of minimum message length inference. In *Australasian Joint Conference on Artificial Intelligence*; Springer: Berlin/Heidelberg, Germany, 2010; pp. 82–91.
- 36. Li, M.; Vitányi, P. An Introduction to Kolmogorov Complexity and Its Applications; Springer: New York, NY, USA, 2008; Volume 3.

- 37. Solomonoff, R.J. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Inf. Theory* **1978**, *24*, 422–432.
- Dowe, D.L. Introduction to Ray Solomonoff 85th memorial conference. In Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence; LNAI 7070; Springer: Berlin/Heidelberg, Germany, 2013; pp. 1–36.
- Makalic, E.; Allison, L.; Dowe, D.L. MML inference of single-layer neural networks. In Proceedings of the 3rd IASTED International Conferences Artificial Intelligence and Applications, Benalmadena, Spain, 8–10 September 2003; pp. 636–642.
- Fitzgibbon, L.J.; Dowe, D.L.; Allison, L. Univariate polynomial inference by Monte Carlo message length approximation. In International Conference Machine Learning; ICML: Sydney, Australia, 2002; pp. 147–154.
- 41. Wu, Z.; Pan, S.; Chen, F.; Long, G.; Zhang, C.; Philip, S.Y. A comprehensive survey on graph neural networks. *IEEE Trans. Neural Netw. Learn. Syst.* 2020, 32, 4–24.
- 42. Chong, E.; Han, C.; Park, F.C. Deep learning networks for stock market analysis and prediction: Methodology, data representations, and case studies. *Expert Syst. Appl.* 2017, *83*, 187–205.
- Qiu, X.; Zhang, L.; Ren, Y.; Suganthan, P.N.; Amaratunga, G. Ensemble deep learning for regression and time series forecasting. In Proceedings of the 2014 IEEE Symposium on Computational Intelligence in Ensemble Learning (CIEL), Orlando, FL, USA, 9–12 December 2014; pp. 1–6.
- 44. Sezer, O.B.; Gudelek, M.U.; Ozbayoglu, A.M. Financial time series forecasting with deep learning: A systematic literature review: 2005–2019. *Appl. Soft Comput.* 2020, *90*, 106181.
- 45. Hochreiter, S.; Schmidhuber, J. Long short-term memory. Neural Comput. 1997, 9, 1735–1780.
- Li, J.; Bu, H.; Wu, J. Sentiment-aware stock market prediction: A deep learning method. In Proceedings of the 2017 International Conference on Service Systems and Service Management, Dalian, China, 16–18 June 2017; pp. 1–6.
- 47. Zhang, X.; Tan, Y. Deep stock ranker: A LSTM neural network model for stock selection. In *International Conference on Data Mining and Big Data*; Springer: Cham, Switzerland, 2018; pp. 614–623.
- Bukhari, A.H.; Raja, M.A.Z.; Sulaiman, M.; Islam, S.; Shoaib, M.; & Kumam, P. Fractional neuro-sequential ARFIMA-LSTM for financial market forecasting. *IEEE Access* 2020, 8, 71326–71338.
- 49. Cheng, T.; Gao, J.; Linton, O. Nonparametric Predictive Regressions for Stock Return Prediction; Working Paper; University of Cambridge: Cambridge, UK, 2019.
- 50. Gao, J. Modelling long-range-dependent Gaussian processes with application in continuous-time financial models. *J. Appl. Probab.* **2004**, *41*, 467–482.
- 51. Fama, E.F.; French, K.R. Dividend Yields and Expected Stock Returns; University of Chicago Press: Chicago, IL, USA, 2021; pp. 568–595.
- 52. Keim, D.B.; Stambaugh, R.F. Predicting returns in the stock and bond markets. J. Financ. Econom. 1986, 17, 357–390.
- Dowe, D.L.; Korb, K.B. Conceptual difficulties with the efficient market hypothesis: Towards a naturalized economics. In Proceedings of the Information, Statistics and Induction in Science, World Scientific, Melbourne, Australia, 20–23 August 1996; pp. 212–223.
- Lai, G.; Chang, W.C.; Yang, Y.; Liu, H. Modeling long-and short-term temporal patterns with deep neural networks. In Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, New York, NY, USA, 8–12 June 2018; pp. 95–104.