



Uni-variate Prediction Time-Series Foundation Models in Finance

by

Filip Topic

September 2024

Dissertation Supervisor: Ramin Okhrati Dissertation submitted in part in
fulfilment of the Master of Science degree in Banking and Digital Finance
Institute of Finance and Technology University College London

¹⁰ **ABSTRACT**

Keywords:

ACKNOWLEDGEMENTS

CONTENTS

| | | |
|----|--|-----------|
| | Abstract | 1 |
| 15 | Acknowledgement | 2 |
| | List of Figures | 6 |
| | List of Tables | 7 |
| | 1 Introduction | 8 |
| | 1.1 Time-Series prediction | 8 |
| 20 | 1.2 Transformer-based Models | 8 |
| | 1.2.1 Time Series Foundation Models (TSFM) | 8 |
| | 2 Background | 9 |
| | 2.1 The Transformer | 9 |
| | 2.1.1 Encoder | 10 |
| 25 | 2.1.2 Decoder | 12 |
| | 2.1.3 Summary | 13 |
| | 2.2 Transformer-based models | 14 |
| | 2.2.1 NLP Transformer-based models | 14 |
| | 2.2.2 CV (computer vision) Transformer-based models | 14 |
| 30 | 2.2.3 Speech and Audio processing Transformer-based models | 14 |
| | 2.2.4 Time-series Transformer-based models | 14 |
| | 2.3 Time-series Foundation models (TSFM) | 15 |
| | 3 Literature review | 18 |
| | 3.1 TimeGPT-1 [20] | 18 |
| 35 | 3.1.1 Architecture | 18 |
| | 3.1.2 Training Data | 18 |
| | 3.2 Lag-Llama [43] | 19 |
| | 3.2.1 Tokenization | 19 |
| | 3.2.2 Architecture | 19 |
| 40 | 3.2.3 Training data | 20 |

| | | |
|----|---|-----------|
| | 3.3 Time Series Foundation Models in Finance | 20 |
| | 4 Methodology and Data | 22 |
| | 4.1 Time-series prediction evaluation | 22 |
| | 4.2 Data | 24 |
| 45 | 4.2.1 Data Pre-processing..... | 25 |
| | 4.2.2 Time-series data characteristics | 25 |
| | 4.2.3 Data exploration | 26 |
| | 4.3 Experiment Design | 27 |
| | 4.3.1 Time Series Cross Validation (TSCV) | 27 |
| 50 | 4.3.2 Time Series Foundation Models..... | 29 |
| | 4.3.3 Benchmark models | 30 |
| | 4.3.4 Experiment..... | 30 |
| | 4.3.5 Experiment configurations | 31 |
| | 4.3.6 Exploration phase..... | 32 |
| 55 | 4.3.7 Aggregation phase 1..... | 32 |
| | 4.3.8 Model-parameter optimization phase | 33 |
| | 4.3.9 Aggregation phase 2..... | 34 |
| | 5 Results and Discussion | 35 |
| | 5.1 Results | 35 |
| 60 | 5.2 Discussion..... | 35 |
| | 5.2.1 Information disparity between fine-tuned TSFMs and benchmark models + zero-shot TSFMs..... | 35 |
| | 5.2.2 Limitations | 35 |
| | 5.2.3 | 36 |
| 65 | A History of time-series forecasting | 44 |
| | A.1 Brief History of time-series prediction in Finance..... | 44 |
| | A.2 Brief History of modern time-series prediction methods..... | 45 |
| | A.2.1 Statistical models | 45 |
| | A.2.2 Machine Learning (ML) models..... | 46 |
| 70 | A.2.3 Deep Learning models..... | 47 |

| | | |
|----------|--|-----------|
| B | TimeGPT-1 | 48 |
| C | Lag-Llama | 49 |
| D | Methodology | 50 |
| | D.1 Time-series prediction | 50 |
| 75 | D.1.1 Prediction error | 50 |
| | D.1.2 MDA | 50 |
| | D.1.3 MES | 50 |
| | D.2 Data | 51 |
| | D.2.1 Data sourcing | 51 |
| 80 | D.2.2 Partial autocorrelation | 51 |
| | D.3 Experiment design | 51 |
| | D.3.1 ARIMA | 51 |
| | D.3.2 Naive Simple Autoregressor (NSA) | 52 |
| | D.3.3 Loss | 52 |
| 85 | D.3.4 Benchmark | 52 |

LIST OF FIGURES

| | | |
|----|---|----|
| | 2.1 Transformer architecture schema [58] | 10 |
| | 2.2 Scaled Dot product attention [58] | 12 |
| | 3.1 Lag-Llama architecture [43] | 20 |
| 90 | 4.1 Breakdown of all the time/series used by trend, presence of cyclical patterns and stationarity | 26 |
| | 4.2 Breakdown of all time-series by type of data and frequency..... | 26 |
| | 4.3 Breakdown of all time-series by time-series features, type of data and frequency. | 27 |
| 95 | 4.4 Schematic view of the TSCV strategy | 29 |
| | 4.5 Schematic view of a single fold from the TSCV technique..... | 31 |
| | 4.6 Schematic view of the Exploration stage..... | 33 |

LIST OF TABLES

| | | | |
|-----|-----|--|----|
| | 4.1 | Types of data used..... | 24 |
| 100 | 4.2 | Available frequencies of different types of data..... | 24 |
| | 4.3 | Time periods used for different frequencies of data..... | 24 |
| | 4.4 | Different experiment parameters and their values | 32 |

1. INTRODUCTION

1.1 Time-Series prediction

105 Since the beginning of financial markets, their participants have had the desire to predict the future values of instruments being traded there. And for this purposes, they have used many different techniques - majority of which have fallen short of randomly guessing the direction of the price movement. In the past few decades, many models have emerged which claim excellent capability of time-series prediction. First were the statistical models
110 (such as ARIMA), then the Machine Learning models (such as XGBoost), then Deep Learning models (such as DeepAR).

1.2 Transformer-based Models

The revolution in the field of Machine Learning came in 2017 with the invention of the Transformer architecture (more specifically: the attention mechanism). Transformer based
115 models have taken the whole field of ML by storm, however, their application has mostly been in the sub-field of Natural Language Processing. Recently they started seeing use in the field of time-series prediction (eg. Informer).

Time Series Foundation Models (TSFM)

Time-series foundation models are large transformer-based models which are designed for
120 the purpose of time-series prediction and which have been pre-trained on a large amount of time-series data. This research will explore their use on financial time-series data

2. BACKGROUND

2.1 The Transformer

In 2017, Vaswani et al. [58] introduced a ground-breaking model called the "Transformer".

Transformer is a model designed to solve the sequence-to-sequence mapping problem. In a sequence-to-sequence problem we have a sequence which consists of tokens¹ coming from a finite vocabulary² which are in a specific order, and we have an output sequence which is again a sequence of tokens in a specific order. The task of a sequence-to-sequence model is to learn the correct relationship between the input and output sequences so that when the model is presented with an unseen input sequence, it can predict what the correct output sequence is. In essence, sequence-to-sequence model's goal is to learn the "meaning" of the relationship between the input and output sequence. An example of sequence-to-sequence tasks are:

1. Language translation - where an input sequence could be a sentence in our native language and the output sequence would be the correct translation of that sentence in a foreign language.
2. Chatbots - where an input sequence could be a question and the output sequence the correct answer to that question.
3. Time-series forecasting - where an input sequence is a certain time-series and the output sequence is the future values of that time-series.

Before the Transformer, there have been many ML models which attempted to solve these tasks. Sutskever et al. (2014) [52] used multilayer LSTM³ (type of RNN⁴ model) for this task. LSTM [25] and GRU⁵ [10] used to be state-of-the art sequence-to-sequence

¹A token is a multi-dimensional numeric vector representation of an element in the sequence. Reason why sequence of tokens is used as input to a sequence-to-sequence model is because ML models can only "understand" numbers - therefore they can't work with some type of sequences (eg. human language sentences). Reason why tokens are vectors rather than simple numbers is because they are designed to represent the semantic meaning of a real-world element they represent and the semantic meaning of a real-world element could change depending on the context. Therefore different dimensions of a token account for different meaning of that element in different contexts.

²The vocabulary doesn't always need to be finite, as we will see later on.

³LSTM stands for Long Short Term Memory cell

⁴RNN stands for Recurrent Neural Network. It is a type of Artificial Neural Network architecture.

⁵GRU stands for Gated Recurrent Unit. GRU is also a type of RNN.

models, however they have some key issues. They require large memory⁶, their nature is
 145 sequential⁷, and they suffer from an information bottleneck due to the fixed size of the
 hidden state⁸. These problems were addressed by the Transformer.

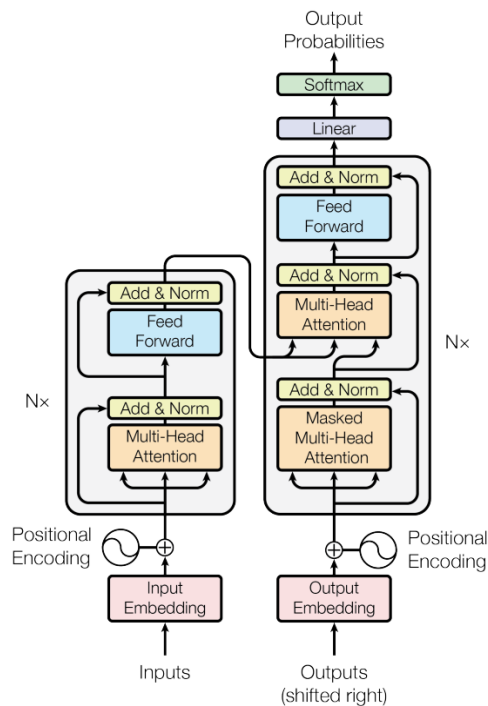


Figure 2.1: Transformer architecture schema [58]

Transformer consists of an Encoder and a Decoder⁹. Figure 2.1 is a schematic representation of the Transformer model. Left part is the Encoder and the right part is the Decoder.

Encoder

The encoder is a stack of multiple identical layers¹⁰, each layer consisting of two sub-layers:

1. Multi-head self-attention layer

⁶This means that it is difficult to parallelize the training process across training examples and smaller batch sizes had to be used.

⁷This means that it is impossible to parallelize training within the training examples.

⁸The hidden state is an intermittent sequence within a sequence-to-sequence model. It is a product of applying the encoder on the input sequence. The output sequence is generated by applying the decoder to the hidden state. Having a hidden state of fixed size means that some information will inherently get lost when a very long input sequence is fed into the model.

⁹Similar to already mentioned RNNs.

¹⁰In the original paper there are 6 of these layers in the encoder stack

2. Feed-forward neural network

Multi-head Self-attention layer consists of N so-called "attention heads"¹¹. An attention head is essentially a series of arithmetic operations performed on the input sequence: Let's say we have a sequence of tokens $S = s_1, s_2, \dots, s_n$ the dimension of input tokens is $1 \times d_{\text{model}}$ ¹². An attention head consists of weight matrices¹³ W^Q , W^K and W^V ¹⁴ of the dimension $(d_{\text{model}} \times d_k)$ where $d_k = d_{\text{model}}/N$. Multiplying these weight matrices by token s_i from the input sequence creates Q_i , K_i and V_i vectors respectively of dimension $1 \times d_k$. For each token i , dot product of Q_i and K_j (for all j , $0 < j < n+1$) is calculated: $\text{dot}_{i,1}, \text{dot}_{i,2}, \dots, \text{dot}_{i,n}$. These dot-products are then scaled¹⁵, a softmax layer is applied¹⁶ to all these dot products. Each dot product $\text{dot}_{i,j}$ then multiplies the corresponding V_j to get the series of V vectors: $V_{i,1}, V_{i,2}, \dots, V_{i,n}$. All these V vectors are then element-wise summed up to Z_i (of dimension $1 \times d_k$) which represents the self-attention output for the token i . This is repeated for all n tokens and creates a series of vectors Z_1, Z_2, \dots, Z_n (one for each token) which are then vertically stacked into matrix Z of dimension $n \times d_k$. This matrix Z_h is the output of the single self-attention head h . If we vectorize this process¹⁷, we can express attention as

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

See Figure 2.2:

This process goes on in parallel in every head of the multi-head self-attention layer. Finally, the Z matrices of each self-attention head Z_1, Z_2, \dots, Z_N are horizontally concatenated to produce the output of the whole multi-head self-attention layer L : Z_L (of dimension $n \times N \times d_k$). Multi-head self-attention layer has a weight matrix W_O of dimension $N \times d_k \times d_{\text{model}}$ which gets multiplied by Z_L to produce the final output of the multi-head self-attention layer L .

¹¹In the Original paper, $N = 8$.

¹²In the original paper, $d(\text{model}) = 512$.

¹³These weight matrices are all learnable weights.

¹⁴ Q , K , and V stand for Query, Key and Value.

¹⁵Scaling factor is $1 / \text{square root of } d(k)$. According to the paper, this is done so that the weights are more stable during training.

¹⁶Applying Softmax on a series of numbers means scaling all of them so that they sum up to 1.

¹⁷When we need to do the same arithmetic operations multiple times with different vectors, we can simply concatenate these vectors into matrices and do the operations just once with these matrices.

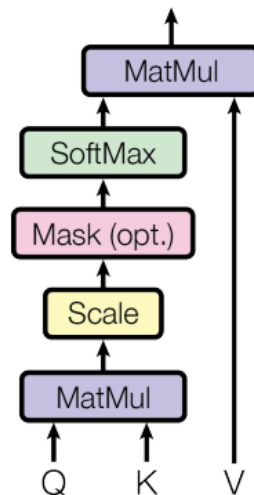


Figure 2.2: Scaled Dot product attention [58]

Before the output of the multi-head self-attention layer goes into a feed-forward neural network¹⁸, a residual connection¹⁹ [23] and layer normalization²⁰ [3] are applied. Residual connection and layer normalization are applied after every sub-layer of the Encoder and Decoder.

Decoder

The Decoder has a similar structure as the Encoder with a few notable differences.

Same as the Encoder, Decoder has N layers. Each layer has three sub-layers:

1. masked self-attention layer
2. encoder-decoder attention layer
3. feed forward layer

In principle, masked self-attention layer works similarly as the multi-head attention

¹⁸This network has two hidden layers with ReLU activation of the first layer and linear activation of the 2nd layer. In the original paper, dimensionality of hidden layers is 2048.

¹⁹Residual connection in this context means adding the input of the layer to the output of that layer before passing to the next layer. This is used to mitigate the vanishing gradient problem and to stabilize the training process.

²⁰Layer normalization is a scaling technique that normalizes the the data across features within each data-sample individually, as opposed to batch normalization which normalizes each feature across multiple data-samples.

layer in the Encoder, however, Q, K, and V come from the already generated output²¹, and the attention is calculated only for the current and previously generated outputs i.e. for an already generated token i , masked self-attention head only calculates Z_1, Z_2, \dots, Z_i .
190 Otherwise, the process is the same as in the Encoder multi-head self-attention sub-layer²². (See Figure 2.1.)

Encoder-decoder attention layer also works similarly as the multi-head self-attention layer in the Encoder, however in this layer Q comes from the previous masked self-attention layer but K and V come from the Encoder. This allows every position in the Encoder to
195 attend all positions in the input. (See Figure 2.1.)

Summary

Even though the Attention mechanism has been around for a while [4], the self-attention model, as implemented in the Transformer, was a revolutionary step - as it solved the issue of long range dependancies²³. This concept enables the model to learn the intrinsic
200 structure of the sequence it is presented with and therefore builds a certain level of actual "understanding" of the inputs rather than simple "fitting".

Another revolutionary aspect of the Transformer is the positional encoding. Positional encoding injects information about the absolute and relative position of each token in the sequence²⁴ thus allowing the model to learn about the importance of relative positions of
205 tokens²⁵.

Finally, and probably most importantly, the way Transformer processes inputs is inherently parallelizable - allowing significant scaling benefits to be exploited and very large models to be built. This is beneficial as larger models can learn more intricate

²¹This makes the Decoder particularly suitable for autoregressive tasks as it predicts tokens sequentially based on which tokens have been predicted before, as opposed to the Encoder which would predict a sequence of tokens which it thinks would have the highest probability of being true overall without regard for the order of the sequence.

²²i.e. Encoder also stacks the outputs of each head, applies the W_o matrix, the residual connection and layer normalization.

²³Earlier attention models had issues relating tokens which were too far apart due to attention being applied only on hidden states (hidden states are a feature of RNN models) rather than on inputs directly.

²⁴this is done by sumating each input token with a positional vector of equal dimension. Value of each number in the positional vector is a sin (or could be cos) function of the position of the token it is being added to as well as of the position inside the positional vector.

²⁵This is most important on the field of NLP (Natural Language Processing) tasks.

patterns in the data and they allow for richer vector representations of the input data²⁶ ²⁷

210 2.2 Transformer-based models

Many have built on the success of the Transformer and have come up with their own models which retain many concepts from the original Transformer.

NLP Transformer-based models

Among the most famous NLP transformer-based models are the BERT (2018) [15], T5
 215 (2019) [44] which themselves have been built upon many times (2019 alBERT [31], 2019
 roBERTa [34], 2019 distilBERT [46], 2020 mT5 [64], 2024 flan-T5 [11]), however, OpenAI's
 GPT²⁸ series and LLaMA series [57] take the crown as the most famous Transformer-based
 NLP models by far.

CV (computer vision) Transformer-based models

220 Among the most famous CV transformer-based models are 2020 ViT [18], 2020 DETR [7],
 and 2021 swin-Transformer [36].

Speech and Audio processing Transformer-based models

Time-series Transformer-based models

It was a matter of time Transformer-based architectures appeared in time-series
 225 prediction tasks. Li et al. (2019) [32] wrote a pioneering work on transformers
 (LogSparse Transformer) in time-series forecasting where they addressed the issue of the
 quadratic space complexity²⁹ of the Transformer³⁰ and proposed a more task-appropriate

²⁶Vector representations in the original Transformer had the dimension 512, whereas some modern transformer-based models support dimensions in the thousands.

²⁷Richer representations of the input data means that the model is able to better understand more nuanced differences between some similar input tokens.

²⁸GPT stands for Generative Pre-trained Transformer

²⁹Quadratic space complexity of a model means that the amount of working memory a model uses grows quadratically with the size of the input of the model. This is especially problematic for longer inputs for which the required memory could explode.

³⁰They used heuristic approach to reduce the storage complexity to $O(L \cdot \log(L))$.

version of self-attention³¹. Zhou et al. (2021) [69] proposed a model (Informer) which also reduces computational and space complexity to $O(L \cdot \log(L))$ using "ProbSparse" self-attention mechanism³² and self-attention distilling³³. Zhou et al. (2021) [70] incorporate a seasonal-trend decomposition approach [12] [59], along with Fourier analysis into the transformer-based model (FEDformer)³⁴. This approach saw great improvements in time-series prediction. Improvements were in terms of lower prediction errors as well as in terms of the distribution of predictions being closer to the distribution of ground-truth values according to the Kolmogorov-Smirnov distribution test³⁵ [38].

2.3 Time-series Foundation models (TSFM)

Foundation Models (FM) are a class of deep models which are pre-trained on a large amount of data - thus being able to generalize³⁶ well as they have been taught many diverse patterns. FMs saw success on the fields of CV and NLP so it is only natural that development of FMs start developing on the time-series front. And so was the case: Since 2022, over 50 models, which can be classified as Time-series Foundation Models, have been released. If we analyze the TSFM landscape, according to the taxonomy proposed by Liang et al. (2024) [33], we can see that the TSFM is a very diverse class of models. TSFMs can be classified according to:

- Type of time-series it is working with, which can be:

1. Standard time-series³⁷.
2. Spatial time-series³⁸.

³¹They proposed so-called "convolutional self-attention" which brings local context awareness to Q-K matching

³²ProbSparse self-attention uses Query Sparsity Measurement - a metric which helps determine which Qs are more "dominant" so that the attention can be calculated for only those Qs. This is an upgrade from LogSparse Transformer, which used a heuristic to determine Which QK pairs will be calculated.

³³Distilling means that the outputs of each the self-attention layers are smaller than their inputs - thus reducing the number of calculations in every subsequent layer relative to the previous one.

³⁴FED stands for Frequency Enhanced Decomposition.

³⁵Authors of the FEDformer claim even though Informer predictions were good, they don't have the same distribution as the ground-truth time-series.

³⁶Generalization is the ability of a model to perform well on data that it hasn't been trained on.

³⁷Standard time-series is a simple sequence of any number of data-points, each associated with a time-stamp, ordered chronologically.

³⁸Spatial time-series is a standard time-series but with a spatial dimension as well.

3. Trajectory³⁹.

4. Event sequence⁴⁰.

250

- Model architecture, which could be:

1. Transformer-based.

2. Non-Transformer-based⁴¹.

3. Diffusion-based [49] [24]⁴².

- The nature of the models' pre-training. which could be:

255

1. Pre-trained LM (Large Model)⁴³.

2. Self-supervised⁴⁴.

3. Fully supervised.

- Capabilities to adapt to a new time-series, which could be:

1. Fine-tuning.

260

2. Zero-shot learning.

3. Prompt engineering.

4. Tokenization.

In this dissertation, I will only consider a small subset of TSFMs; ones that belong to a class of standard time-series, Transformer-based, Self-supervised (Generative) models

³⁹Trajectory is a sequence of time-stamped locations which describe the movement of an object in space.

⁴⁰Event sequence is a chronologically ordered sequence of events within a specific context.

⁴¹These models are usually MLP, RNN or CNN -bases models.

⁴²Diffusion-based models are self-supervised models usually applied in image generation. They are trained by adding gaussian blur to the training-sample in a markov process manner and then applying (usually) CNN-based model to learn how to un-blur the same sample.

⁴³These models use a model (usually Large Language Model) which has already been trained on some other type of sequential data, for another task and they just adopt it for time-series. Adoption strategies usually involve changing the way the inputs are tokenized in order to work better with time-series data.

⁴⁴Self-supervised models can be further split into Generative, Contrastive and Hybrid models. Generative models are those which have been trained to reconstruct the input data. They are particularly useful for tasks that involve generating a "missing" part of the input data (such as time-series forecasting). Contrastive models are those that have been trained to distinguish between "similar" and "dissimilar" pairs of data. As such, they are more appropriate for classification tasks (In time-series analysis, they are usually used for anomaly detection). As the name suggests, Hybrid models use mix of these two strategies.

265 with fine-tuning capability. However, this class of TSFMs is itself very large [33] (in no particular order: PatchTST [41], MOIRAI [62], Lag-Llama [43], TimeSiam [17], Timer [35], TimesFM [14], UniTS [19], TimeGPT-1 [20], Chronos [1], MTSMAE [54]) so the efforts of this dissertation will be focused on two pioneering examples, trained on vast collection of time-series data spanning multiple domains [33]: Lag-Llama and TimeGPT-1.

3. LITERATURE REVIEW

3.1 TimeGPT-1 [20]

TimeGPT-1 is closed-source model - meaning that there is a lot of opacity to model's architecture, parameters and training data.

Architecture

TimeGPT-1 has encoder-decoder structure with multiple layers, each having residual connections and layer normalization. It utilizes self-attention mechanism based on the original Transformer [58].

Special capability of TimeGPT-1 is multivariate time-series forecasting - meaning that it is able to take into account "special events" and exogenous features⁴⁵ when making predictions on a target time-series. However, in order to take full advantage of this feature, we would have to know for certain the realizations of the exogenous variables in the future which is almost never the case in time-series forecasting⁴⁶. TimeGPT-1 solves this by separately forecasting the exogenous variables into the future and then basing the forecast of the target time-series on its own forecast of the exogenous features.

Training Data

Authors claim TimeGPT-1 has been trained on a data-set containing over 100B data-points - the largest collection of publically available time-series according to their knowledge. Data comes from the domains of finance, economics, demographics, healthcare, weather, IoT sensor data, energy, web traffic, sales, transport, and banking. Due to diversity of data domains, the training set contains time-series with many different characteristics: seasonalities, cycles, trends, noise, outliers. Having been trained on such a diverse dataset, TimeGPT-1 has very strong robustness and generalization capabilities.

⁴⁵Special events and exogenous features are time-series which are assumed to have some influence on the target time-series. Example of the former being bank holidays if the target time-series is number of flights per day and example of the latter being price of petrol if we are forecasting number of cars on the road.

⁴⁶because we can never predict anything in the future with certainty.

3.2 Lag-Llama [43]

Tokenization

295 Lag-Llama constructs "lagged features" from the past values of the time-series according to a set of frequency-dependant set of lag indices $L_{\text{indices}}=1, \dots, N$. The set of lagged features of a point y_t at timestamp t is then $\{y_{t-L_{\text{indices}}[i]}, \text{ for all integers } i \text{ such that } 0 < i < N+1\}$. Lag-Llama also constructs "date-time features" F which for a point y_t in time t contain information about second-of-the-minute, minute-of-the-hour, ..., month-of-the-year. Final
300 tokenization is done by simply concatenating the date-time features and lagged features into a single vector. One of Lag-Llama's hyperparameters, "Context Length" cl determines how many timepoints it will consider when making a prediction. When making a prediction for y_{t+1} , Lag-Llama will use $\{y_{t-cl+1}, y_{t-cl+2}, \dots, y_t\}$ in order to make a prediction. It is important to keep in mind that even though Lag-Llama only uses cl previous points for
305 making a prediction, due to how the tokenization process works - specifically the lagged features, tokens contain information from timestamps older than cl time-points in the past.

Architecture

Lag-Llama is an open-source, decoder-only TSFM based on LLaMA [57] architecture with
310 reduced number of parameters⁴⁷ (2.5m) . Similarly as LLaMA, Lag-Llama utilizes concepts such as RMSnorm [68]⁴⁸, RoPE [51]⁴⁹ and SwiGLU [47]⁵⁰ activation function[40]. On top of N decoder layers, Lag-Llama has a "parametric distribution head". Parametric distribution head is a module which predicts the distribution parameters of the next prediction, given

⁴⁷Original LLaMA series comes in sizes of 6.7B, 13.0B, 32.5B and 65.2B parameters.

⁴⁸RMSnorm is a layer normalization technique. Zhang and Sennrich (2019) demonstrate that centering component of the LayerNorm technique is dispensable (and computationally expensive) and propose their own layer normalization strategy called RMSnorm which delivers similar benefits (more stable training and better model convergence) but with lower computational cost.

⁴⁹RoPE stands for Rotary Positional Encoding. This is a method for positional encoding. RoPE enables valuable properties, including the flexibility of sequence length, decaying inter-token dependency with increasing relative distances, and the capability of equipping linear self-attention with relative position encoding.

⁵⁰SwiGLU activation function is a combination of swish and GLU activation functions. With beta hyperparameter equal to 1, it is of similar shape as ReLU but completely smooth (continuous) - therefore "behaving better" during model training.

the set distribution. Authors of the paper have chosen student-t distribution for the
 315 distribution head, meaning that at inference time, the model predicts the degrees of
 freedom, mean and scale [50] of the t-distribution from which it randomly draws n samples
 - idea being that this sample of n predictions gives a probabilistic insight into the model's
 prediction. See Figure 3.1:

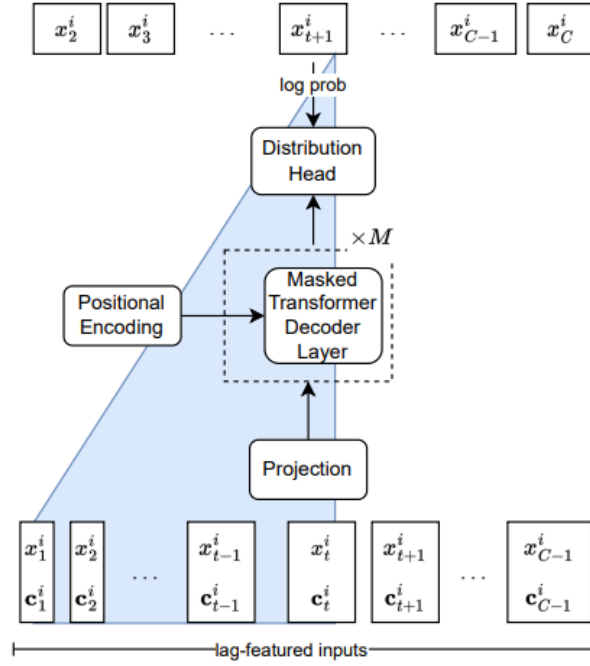


Figure 3.1: Lag-Llama architecture [43]

Training data

320 Lag-Llama is pre-trained on diverse set of time-series domains such as energy,
 transportation, economics, nature, air quality and cloud operations.

3.3 Time Series Foundation Models in Finance

There has been work done on time-series prediction using Foundation models in Finance.
 However, none of it was of uni-variate nature as all used some form of exogenous data
 325 and none of it was with the use of dedicated Time Series Foundation Models. Yu et
 al. (2023) [65] successfully used GPT-4 and LLaMA 13B LLMs with instruction-based
 fine-tuning, one-shot and few-shot inference on company profile, finance/economy news
 and index price data to predict NASDAQ-100 returns. However, they were predicting

bins and movement direction of future values instead of actual values. Chen et al. (2023)
330 [9] used ChatGPT-informed Graph Neural Network (GNN) with prompt engineering and
financial news headlines to predict direction of stock movements. Xie et al. (2023) [63]
used chatGPT with chain-of-thought - enhanced zero-shot prompting with twitter data
to predict stock price movement. Wimmer et al. (2023) [60] used CLIP (pre-trained CV
model) [42] in combination with LSTM to predict German Stock Market movements using
335 Open, High, Low and Close prices. As to my knowledge, this is the first paper strictly
examining the performance of uni-variate TSFMs on time-series prediction in Finance.

4. METHODOLOGY AND DATA

This Research aims to answer 3 main questions:

1. How good are TSFMs, namely Lag-Llama and TimeGPT-1, on financial time-series data?
2. In which cases do they perform better/worse?
3. How to improve TSFMs performance?

Methodology of this research is designed in a way to try best answer these three questions in a scientific, statistically significant manner.

4.1 Time-series prediction evaluation

Time-series predictions are evaluated using the traditional regression metrics due to the continuous nature of the target variables. All regression metrics are based on some variation of the prediction error (D.1.1). From the prediction error, following metrics are derived, and commonly used in time-series prediction tasks in Finance[26]: RMSE (Root Mean Square Error), MAE (Mean Absolute Error), MAPE (Mean Absolute Percentage Error) and R^2 .

Error-based evaluation metrics are scale-dependant as the errors themselves are inherently scale-dependant. Since the experiments were ran on different types of data of different frequencies, each experiment is on a different scale, therefore, the raw evaluation metrics for different experiments cannot be used for model comparison across different time-series. It could be argued that MAPE accounts for this as it calculates absolute errors as a percentage of ground-truth actual values, however, due to different time-series having different time-series features such as noise, trend and seasonality, some time-series are inherently more difficult to predict than the others hence the MAPE would be on a different scale on different time-series that reason. Similar case could be made for R^2 as it scales the sum of square errors (SSE) by the total sum of squares (TSS), however, it also doesn't account for the fact that some time-series are inherently more predictable than the others. A perfect example of this is CHAPS data and Stock index data: most models score

quite high R^2 on CHAPS data (some even > 0.7), whereas pretty much all score negative
 365 R^2 on Stock index data (as we will see later on). Due to these reasons, in every experiment,
 models are given a rank from 1 to 7⁵¹ for each evaluation metric, which denotes how
 well each model did (relative to other models) on that specific experiment. This is a
 common practice in time-series results evaluation and is widely employed - most famously
 in the M-series competitions (prestigious time-series forecasting competition) [37]. Other
 370 scale-invariant metrics such as RMSSE⁵² and MASE⁵³ [28] were considered, however,
 due to the way they are calculated, could facilitate unfair comparison if time-series are
 non-stationary or have heteroskedastic errors. Example of the ranking system is if a model
 has MSE Rank 1 on experiment E, it means that the model was the best performing
 model in experiment E, according to the MSE metric. Since the scale of the ranking
 375 system is constant 1-7, it can be used to compare relative model performance on different
 time-series. However, a flaw of the ranking method is that it doesn't account for exactly
 how much models are better/worse than each other.

An additional metric that will be used in this research is MDA (Mean Directional
 Accuracy) (D.1.2). See [13] for Directional accuracy calculation. Finally, I suggest a new
 380 metric called MES (Mean Equal Sign) to be used in scenario where the time-series being
 predicted is financial returns. For a set of n time-series predictions $\{y_{\text{predicted}}^1, y_{\text{predicted}}^2, \dots, y_{\text{predicted}}^n\}$
 and ground-truth values $\{y_{\text{actual}}^1, y_{\text{actual}}^2, \dots, y_{\text{actual}}^n\}$, I define MES as
 $\text{mean}\{\text{sign}(y_{\text{predicted}}^1 \times y_{\text{actual}}^1), \text{sign}(y_{\text{predicted}}^2 \times y_{\text{actual}}^2), \dots, \text{sign}(y_{\text{predicted}}^n \times y_{\text{actual}}^n)\}$.
 This metric is introduced because MDA alone doesn't fully capture the model's ability to
 385 predict direction of underlying asset value movement (see D.1.3). MDA and MES should
 be especially important in the field of Finance as we are not necessarily only interested in
 predicting the actual values of an asset, but rather we are interested in predicting the
 direction in which the price will go next time step.

⁵¹There are 7 models in each experiment: 2 zero-shot TSFMs, 2 fine-tuned TSFMs and 3 benchmark models - as we will see later on.

⁵²Root Mean Square Scaled Error

⁵³Mean Absolute Scaled Error

4.2 Data

390 This project uses 4 types of financial time-series (See table 4.1). Index price data was sourced from Yahoo Finance using yfinance Python package. Commodity and Exchange rate data was sourced using Alpha Vantage Python API. CHAPS⁵⁴ data was sourced from UK Office for National Statistics⁵⁵. Depending on the frequency and the type of financial data, there is different availability of it. See table 4.2 for information on which frequencies were available in sufficient ammount depending on the type of data. See table 4.3 for information on which time-periods were used depending on the frequency of choice. See D.2.1 for exceptions. All in all, 56 distinct time-series were used for this research.

| Stock Index | Commodity | Exchange Rate | CHAPS |
|-------------|-----------|---------------|--------------|
| S&P 500 | WTI | USD/GBP | Aggregate |
| FTSE 100 | | | Delayable |
| DOWJ | | | Social |
| NASDAQ | | | Staple |
| | | | Work-related |

Table 4.1: Types of data used

| Frequency | Type of Data | | | |
|-----------|--------------|-----------|---------------|-------|
| | Stock Index | Commodity | Exchange Rate | CHAPS |
| Daily | YES | YES | YES | YES |
| Weekly | YES | YES | YES | YES |
| Monthly | YES | YES | NO | NO |

Table 4.2: Available frequencies of different types of data

| Daily | Weekly | Monthly |
|--------------------------|--------------------------|--------------------------|
| 2018-01-01 to 2020-01-01 | 2015-01-01 to 2020-01-01 | 1987-01-01 to 2024-01-01 |
| 2020-01-01 to 2022-01-01 | 2017-01-01 to 2022-01-01 | |
| 2022-01-01 to 2024-01-01 | 2019-01-01 to 2024-01-01 | |

Table 4.3: Time periods used for different frequencies of data

⁵⁴CHAPS is UK debit and credit card spending index

⁵⁵<https://www.ons.gov.uk/economy/economicoutputandproductivity/output/datasets/ukspendingoncreditanddebitcards>

Data Pre-processing

In the field of Finance, there is a concept of "return". Return of an asset A at a certain point in time T: A_T is calculated as

$$(A_T - A_{T-1})/A_{T-1}$$

i.e. the percentage change in the value of that asset between the points in time T and T-1. If we think about each one of the fore-mentioned time-series as prices of an underlying asset, we can represent them as time-series of returns rather than their actual values which is a common practice in financial time-series prediction [26]. This method is not applied to CHAP data as it wouldn't make sense given the context of that data.

Time-series data characteristics

Time-series data exhibits several key characteristics that are essential to understand for effective analysis and forecasting. **Trend** is a long-term increase or decrease in the data. It could take many shapes: none, linear, exponential, polynomial. This research uses Kendall's Tau coefficient of correlation [30]. Kendall's Tau measures the strength of the relationship between two ordinal variables. If a time-series of length L has statistically significant Kendall Tau correlation (p-value < 0.05) with a monotonically increasing sequence {1, 2, ..., L}, this is a sign that there is a linear trend present in the data. If the square roots of every value in a time-series have statistically significant Kendall Tau correlation with a monotonically increasing sequence {1, 2, ..., L}, this is a sign that there is a quadratic trend present in the data. **Seasonality** is the presence of recurring patterns at regular intervals in the data. Chosen method for determining the presence of seasonality is by looking at partial autocorrelation values (see D.2.2 for explanation) in the data and declaring a seasonal pattern present if there is a significant positive or negative partial autocorrelation present. There is no way to determine the significance in this case. Rather, I chose values >0.3 or <-0.3 to be significant. **Stationarity**: A time-series is said to be stationary if the statistical properties such as mean and variance remain constant over time. Stationarity is determined using the Dickey-Fuller test [16]. Other characteristics of time-series are **Noise** and **Volatility** however, they are relatively

425 difficult to classify and are scale-dependant.

Data exploration

56 distinct time-series were used for this research. 38% of all time-series contain quadratic trend, 62% contain no trend (Linear trend was also tested for but none had it). 46% of time-series contain cyclical patterns and 54% contain no discernable cyclical pattern. 88% of time-series are stationary whereas only 12% are non-stationary (see Figure 4.1).

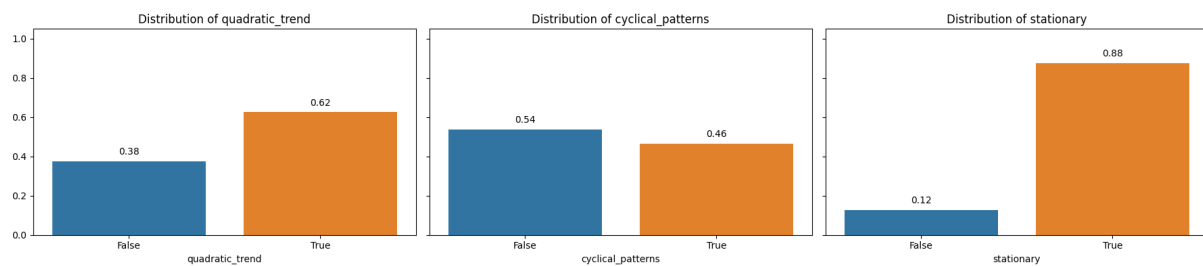


Figure 4.1: Breakdown of all the time/series used by trend, presence of cyclical patterns and stationarity

Half of all time-series used are Stock indices, 27% are CHAPS, 12% Commodity and 11% Exchange rate. 50% of time-series used had daily frequency, 41% had weekly frequency and 0.09% were monthly See Figure 4.2).

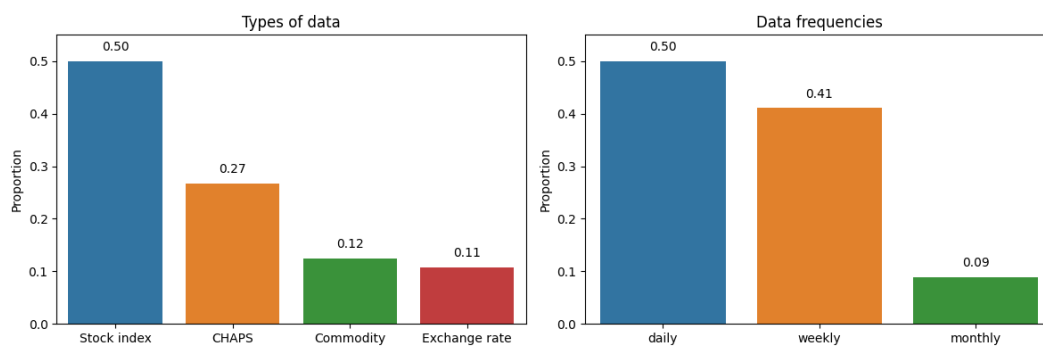


Figure 4.2: Breakdown of all time-series by type of data and frequency.

See Figure 4.3 for breakdown of time-series features by Type of data and Frequency of data. Expected observations are that Commodity, Exchange rate and Stock index data are all stationary and CHAPS data is the only type which contains some non-stationary time-series. This is due to the pre-processing method applied on the former.

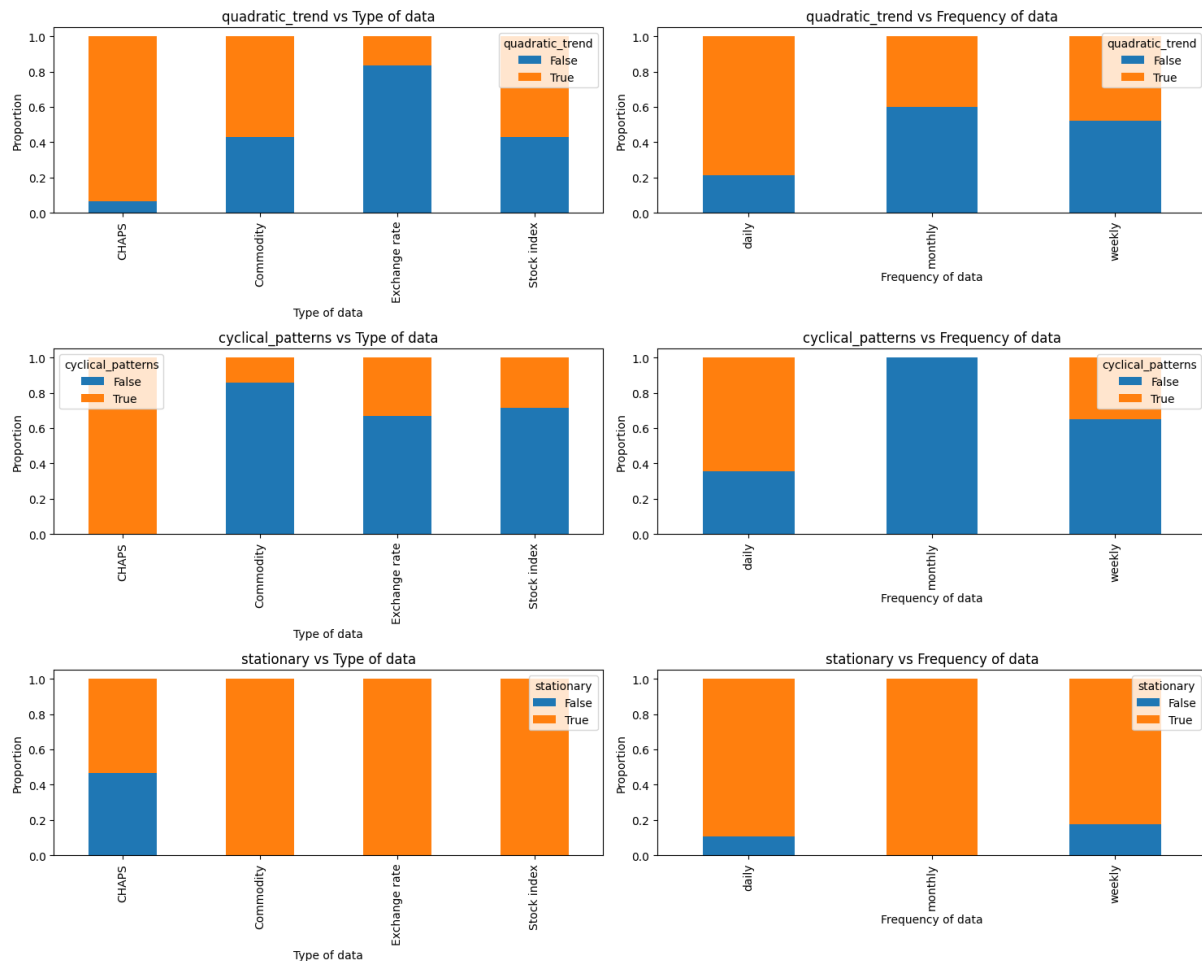


Figure 4.3: Breakdown of all time-series by time-series features, type of data and frequency.

4.3 Experiment Design

Time Series Cross Validation (TSCV)

K-fold cross-validation (CV) is a statistical technique used to evaluate the performance of a model by dividing the dataset into K equally sized subsets, or "folds." The model is trained on K-1 "train" folds and tested on the remaining "validation" fold, and this process repeats K times, with each fold serving as the test set once. The results are then averaged to provide a more reliable estimate of model performance. It is important because it helps

prevent overfitting by providing multiple evaluations of the model's ability to generalize to unseen data. However, in time-series prediction, where data points are temporally dependent, standard K-fold cross-validation cannot be applied directly since randomly partitioning the data breaks the time order and dependencies. To address this, K-fold cross-validation for time-series is modified by partitioning the data in a way that respects temporal structure. This allows for a valid evaluation of time-series models while still leveraging the benefits of cross-validation.

In this research, TSCV works the following way: say we have a time-series S of length l : $S = \{y_1, y_2, \dots, y_l\}$, we specified the length of the train sequence to be n , and the length of prediction horizon to be h . K-fold TSCV method will sample K time-series $\{F_1, F_2, \dots, F_K\}$ from S where each fold F_i ($1 \leq i \leq K$ for all integer i) consists of the "train" and "validation" parts. Fold F_i^j which starts from timepoint j is defined as $F_i^j = \{y_j, y_{j+1}, \dots, y_{j+n-1}, y_{j+n}, \dots, y_{j+n+h-1}\}$. The model is trained on $\{y_j, y_{j+1}, \dots, y_{j+n-1}\}$ and validated on $\{y_{j+n}, \dots, y_{j+n+h-1}\}$ (Size constrains being $K \leq (l-n)/h$). This research only next-point prediction i.e. $h=1$ due to the nature of financial markets where we are primarily interested in the next price. This research uses the version of TSCV where l is fixed which is called "rolling window" TSCV. We employ the following fold sampling strategy: say we have two additional cross-validation parameters: f and r . Folds we sample from S are split into r groups: $\{G_1, G_2, \dots, G_r\}$ such that all folds within a group are sequential i.e. the start-point of the following fold is the time-point right after the start-point of the preceding fold, each group containing f folds. The first group of folds G_1 is always $F_1^1, F_2^2, \dots, F_f^f$ (i.e. first f folds are folds with start-points 1, 2, ..., f respectively). This sampling procedure is repeated r times in a way that the time-distance between the first fold of G_k and the first fold of G_{k+1} is equal to the time-distance between G_{k+1} and G_{k+2} for all $1 \leq k \leq r-2$, and so that the distance between two consecutive groups is maximised subject to constraint $rf \leq (l-n)/h$. r and f parameter values used in the experiments are $r=6$ ⁵⁶ and $f=5$. In simple words: we are sampling five consecutive folds (a group), 6 times in a way that the groups are uniformly spread out over the whole time-series. See Figure 4.4:

⁵⁶An exception to this is with the experiments ran on the monthly data where $r=8$.

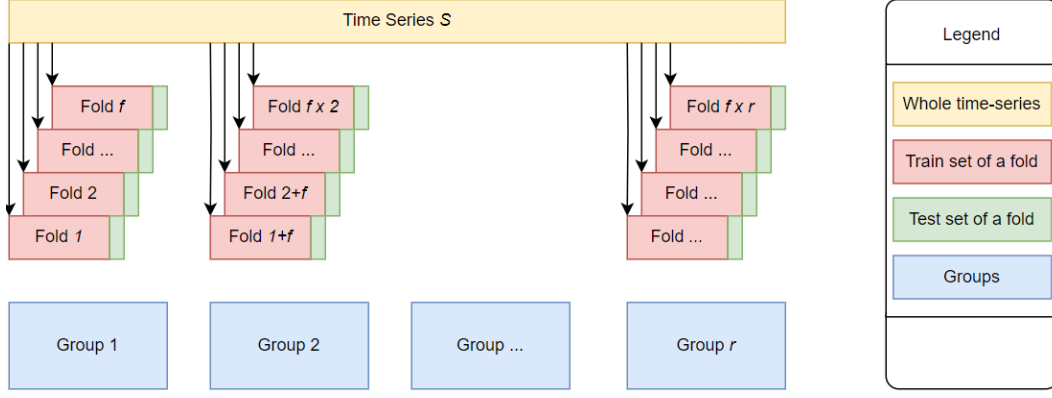


Figure 4.4: Schematic view of the TSCV strategy

Time Series Foundation Models

TSFMs used in this research are Lag-Llama [43] and TimeGPT-1 [20]. Both models have fine-tuning capabilities with accompanying fine-tuning hyper-parameters. Lag-Llama enables the user to choose Batch Size BS , Max Epochs ME , and Learning Rate LR . Batch Size refers to the number of training examples processed in one iteration before the model updates its' parameters. A smaller batch size allows for more frequent updates, while a larger batch size processes more data per update. An Epoch is a single pass through the entire fine-tuning dataset during model fine-tuning. Before an epoch i : E_i ($1 \leq i \leq ME$), model makes a record of all its' parameters. After an epoch i : E_i ($1 \leq i \leq ME$), model's parameters have changed and the loss (see D.3.3) for that epoch L_i is recorded. If $L_i < L_{i-1}$, then the model will keep the updated parameters from E_i , but if $L_i > L_{i-1}$, then the model will revert to the weights it had before E_i . On E_1 , model parameters are updated no matter what as there is no previous loss as a reference, but then this procedure repeats for all $i \leq ME$. A higher me means the model will have more opportunities to learn patterns from the data, but it can also risk overfitting. Learning rate is a hyperparameter that controls how much the model adjusts its weights with respect to the loss gradient; a small learning rate can lead to slow learning, while a large learning rate can cause unstable training (overshooting). These parameters together affect how efficiently and accurately a model learns from data. TimeGPT-1 enables users only to choose number of "Fine-tune Steps" which corresponds to the number of epochs. TimeGPT uses the Adam optimizer with unknown learning hyperparameters, and batch size. Because we don't have sufficient information on TimeGPT-1 fine-tuning process and

fine-tuning hyperparameters, we cannot use it with same Max Epochs (Fine-tune Steps) as Lag-Llama. According to its' authors [20], TimeGPT-1 performance strictly improves with the number of Fine-tune Steps and reaches the plateau at around 100 Fine-tune Steps, whereas Lag-Llama reaches its' full potential with only 4 Max Epochs - as we will see later in this research. This research uses both zero-shot and fine-tuned versions of both Lag-Llama and TimeGPT-1 in order to gain insight into how the models' performance changes with fine-tuning.

Benchmark models

Calculating the evaluation metrics for time-series prediction is not sufficient to answer whether a time-series prediction model is good or not. The results of the evaluation need to be put into context by comparing them against a benchmark in order to make a judgment on how good they are. Chosen benchmark models are: autoARIMA [27] (D.3.1), Naive Simple Autoregressor (D.3.2) and Meta's Prophet model [55].

Experiment

"Running an experiment" on a certain time-series S involves the following: We employ the TSCV technique, in a manner described before, on S . Each fold F_i ($1 \leq i \leq rf$) represents a slice of S so with the train-set length n and prediction horizon $h=1$, $F_i = \{y_1, y_2, \dots, y_n, y_{n+1}\}$. Train set T_i of the Fold F_i is then $T_i = \{y_1, y_2, \dots, y_n\}$. Firstly, the whole train set T_i is used for TSFM fine-tuning. Then, only the last cl (See Section 3.2.1) points of T_i : $\{y_{n-cl+1}, y_{n-cl+2}, \dots, y_n\}$ are used as the context window for fine-tuned TSFMs, zero-shot TSFMs and as the actual train-set for benchmark models (see Section 5.2.1) (see Figure 4.5). After all the models have made their predictions for the fold F_i , their predictions are recorded along with the ground-truth actual value (Test set) from F_i . This is done for every fold of the TSCV process. Nuance worth mentioning is that fine-tuned Lag-Llama is getting fine-tuned only on the first fold of every group of folds (i.e. every 5 folds) whereas fine-tuned TimeGPT-1 is fine-tuned on every single fold (see Section 5.2.2). After the models' predictions have been made and recorded for all folds, evaluation metrics for each model's predictions are calculated (see Section 4.1) and also

recorded.

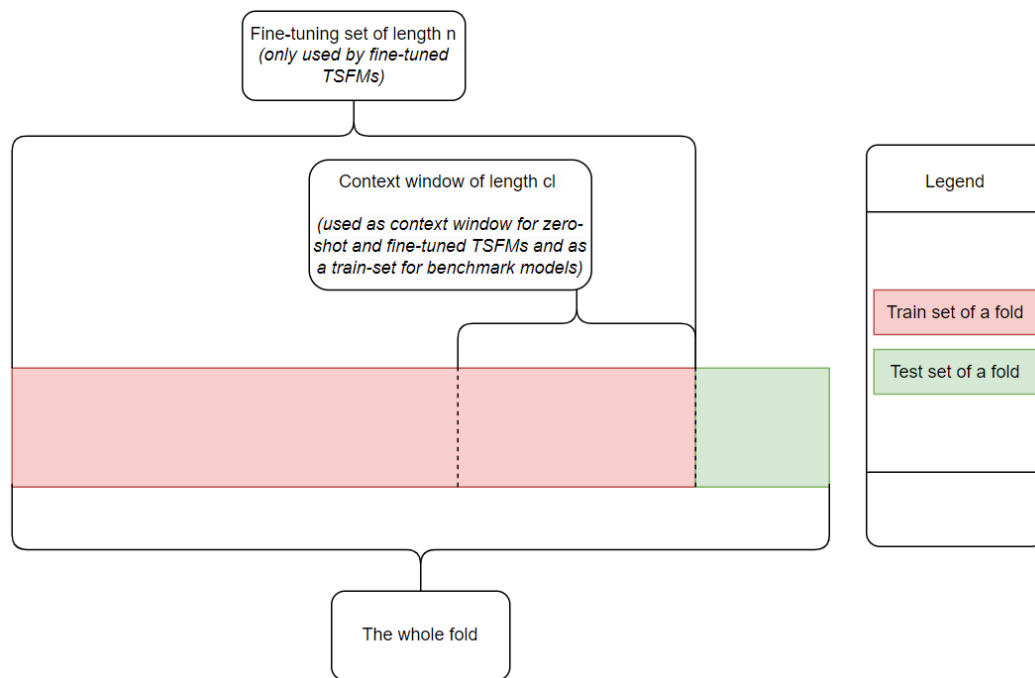


Figure 4.5: Schematic view of a single fold from the TSCV technique

Experiment configurations

Experiment-configuration refers to a unique combination of experiment-parameters. Experiment-parameter is a feature of an experiment which has a certain impact on how one or more models perform. Experiment-parameters can be divided into two groups: data-parameters, and fine-tuning hyper-parameters. Data-parameters are *Type of Data*, *Frequency of Data* and *Context Length*. Fine-tuning hyper-parameters are the following: *Fine-tuning length* denotes the length of the fine-tuning used for TSFM fine-tuning and is applicable to both Lag-Llama and TimeGPT-1. *Max Epochs*, *Batch Size* and *Learning Rate* are only applicable to Lag-Llama and *Fine-tune Steps* is only applicable to TimeGPT-1. See Table 4.4 for details on experiment-parameter values used. Running the experiments across many different experiment configurations allows us to analyze the results of all the experiments according to different experiment-parameters. Breaking down the results by data-parameters would answer the 2nd main question of this research: "In which cases do TSFMs perform better/worse?". Breaking down the results by fine-tuning hyper-parameters would answer the 3rd main question of this research: "How to improve TSFMs' performance?". Breaking down the results by two (or more)

experiment-parameters would give insight into the joint effect that those two (or more) parameters have on model performance which is very relevant because it is unreasonable to assume each experiment-parameter has an isolated effect on model performance. In order to get full granular insight into all joint effects of experiment-parameters, experiments need to be run across every single experiment configuration. Since we have one experiment-parameter with 4 different values and seven experiment-parameters with 3 different values, this gives the total of $4 \times 3^7 = 8748$ different experiment-configurations⁵⁷ - which is computationally unfeasible given that we run all the models, apart from TimeGPT-1, locally.

| Parameters | Parameter values | | | |
|--------------------|------------------|-----------|---------------|-------|
| Type of data | Stock index | Commodity | Exchange rate | CHAPS |
| Frequency of data | Daily | Weekly | Monthly | |
| Context length | 32 | 64 | 128 | |
| Fine-tuning length | 200 | 128 | 64 | |
| Batch size | 5 | 10 | 20 | |
| Max epochs | 4 | 8 | 16 | |
| Learning rate | 5e-4 | 5e-3 | 5e-5 | |
| Fine-tune steps | 100 | 50 | 10 | |

Table 4.4: Different experiment parameters and their values

Exploration phase

Since we are computationally restricted, we divide the whole experimentation process into phases, the first phase being the Experimentation phase. In this phase, we employ a heuristic approach where we keep all the fine-tuning hyper-parameters (last 5 rows of Table 4.4) fixed (respective fine-tuning hyper-parameter values being the ones in the first column of the last four rows of Table 4.4) and run experiments with different combinations of data-parameters. See Figure 4.6:

Aggregation phase 1

After having recorded and aggregated the results (model evaluations) and experiment-configurations in the Exploration phase we group the aggregated results by any single, or

⁵⁷This is not accounting for the fact that Stock index is represented by 4 different time-series and that CHAPS is represented by 5 different time-series.

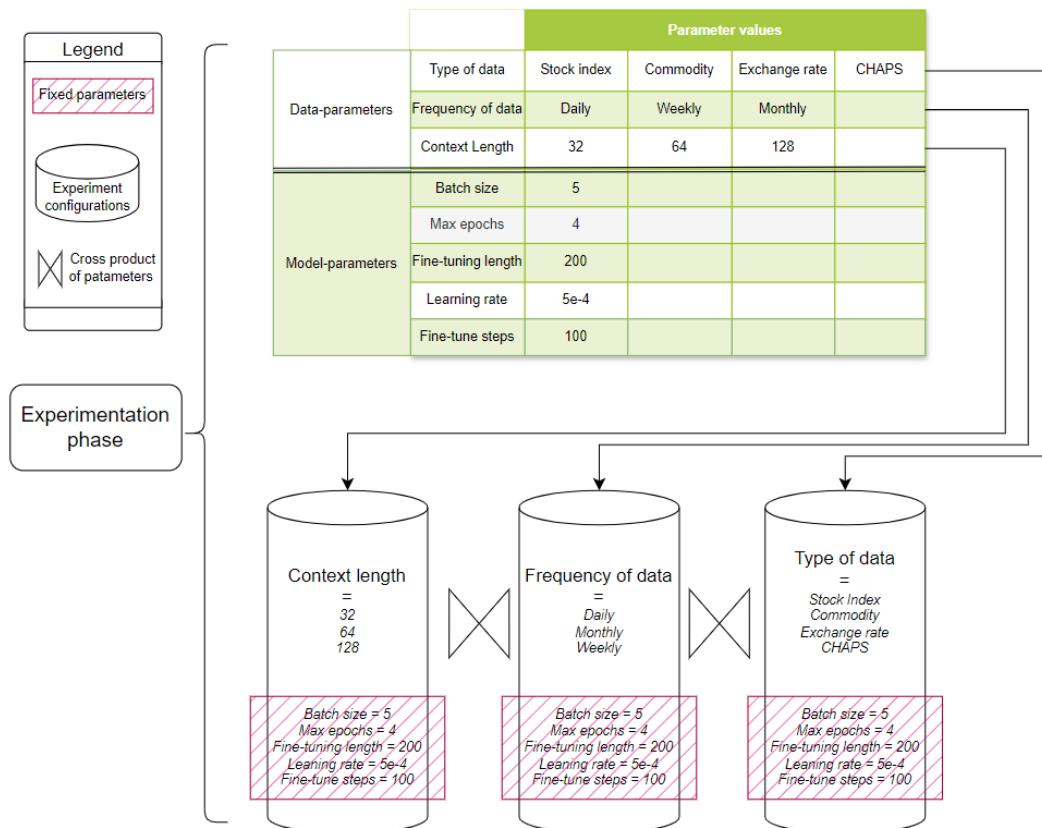


Figure 4.6: Schematic view of the Exploration stage

560 combination of data-parameters and take the mean of the evaluation metric ranks (see Section 4.1) in order to understand how models perform with different data-parameters.

Model-parameter optimization phase

Last phase in the heuristic experimentation approach is to pick (fix) a single data-configuration and run different experiment configurations with fixed data-configuration and variable model-configuration. The strategy for choosing the data-configuration for each TSFM's fine-tuning hyper-parameter optimization is by taking the data-configuration on which the TSFM performed best on - not in terms of the lowest mean rank but in terms of how much lower TSFM's rank was compared to the next best model i.e. we are choosing the data-configuration with the biggest gap between the TSFM and the next best model. After we have chosen the best data-configuration for Lag-Llama and the best data-configuration for TimeGPT-1, we keep it fixed for each TSFM and run experiments across different parameter configuration. This process is identical to the process in the Exploration phase, but this time the data-parameters are fixed, and we

570

vary the fine-tuning hyperparameters in each experiment.

575

Aggregation phase 2

Similarly as in Aggregation Phase 1, we can group the aggregated results by any single, or combination of fine-tuning hyper-parameters and take the mean of the evaluation metric ranks in order to understand how models perform with different fine-tuning hyper-parameters on the data-configuration that we chose. Doing this type of analysis enables

580

us to answer the question "How to improve TSFMs' performance?".

5. RESULTS AND DISCUSSION

5.1 Results

5.2 Discussion

Information disparity between fine-tuned TSFMs and benchmark models + zero-shot TSFMs

Theoretically, fine-tuned TSFMs are systematically advantaged over benchmark models and zero-shot TSFMs as they have access to more information as the set that they are fine-tuned on is larger than the size of the context window which zero-shot TSFMs and benchmark models use. However this is unavoidable as the technical implementation of fine-tuning in TSFMs is such that the length of the fine-tuning set has to be larger than its' context length. If we simply let benchmark models train on the whole fine-tuning set, then zero-shot TSFMs would be systematically disadvantaged in comparison to benchmark models as they would have access to less information than them. The final option is to let zero-shot TSFMs have the context length of same length as the length of the fine-tuning set, however, then the comparison between zero-shot and fine-tuned TSFMs wouldn't be appropriate as fine-tuned TSFMs would have shorter context lengths and the models would therefore be structurally different.

finetuning process of lag Llama and timeGPT cannot be directly compared.

Limitations

Theoretically, fine-tuned versions of TSFMs should be fine-tuned on every single fold as there being distance between the fine-tuning data and the timepoint of prediction increases the likelihood that the model is learning outdated patterns. However, as Lag-Llama is run locally on my machine, fine-tuning process is taking a very long time and it was temporally unfeasible to fine-tune it before every single prediction. So a compromise had to be made and fine-tune it only every 5 predictions.

REFERENCES

- [1] A. F. Ansari, L. Stella, C. Turkmen, X. Zhang, P. Mercado, H. Shen, O. Shchur, S. S. Rangapuram, S. P. Arango, S. Kapoor, et al. Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
- [2] V. Assimakopoulos and K. Nikolopoulos. The theta model: a decomposition approach to forecasting. *International journal of forecasting*, 16(4):521–530, 2000.
- [3] J. Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [4] D. Bahdanau. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [5] G. Box and G. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-Day series in time series analysis and digital processing. Holden-Day, 1970.
- [6] R. G. Brown. *Exponential smoothing for predicting demand*. Little, 1956.
- [7] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [8] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [9] Z. Chen, L. N. Zheng, C. Lu, J. Yuan, and D. Zhu. Chatgpt informed graph neural network for stock movement prediction. *arXiv preprint arXiv:2306.03763*, 2023.
- [10] K. Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [11] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- [12] R. B. Cleveland, W. S. Cleveland, J. E. McRae, I. Terpenning, et al. Stl: A seasonal-trend decomposition. *J. off. Stat*, 6(1):3–73, 1990.

- [13] M. Costantini, J. C. Cuaresma, and J. Hlouskova. Forecasting errors, directional
635 accuracy and profitability of currency trading: The case of eur/usd exchange rate.
Journal of Forecasting, 35(7):652–668, 2016.
- [14] A. Das, W. Kong, R. Sen, and Y. Zhou. A decoder-only foundation model for
time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023.
- [15] J. Devlin. Bert: Pre-training of deep bidirectional transformers for language
640 understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [16] D. A. Dickey and W. A. Fuller. Distribution of the estimators for autoregressive time
series with a unit root. *Journal of the American statistical association*, 74(366a):427–
431, 1979.
- [17] J. Dong, H. Wu, Y. Wang, Y. Qiu, L. Zhang, J. Wang, and M. Long. Timesiam: A pre-
645 training framework for siamese time-series modeling. *arXiv preprint arXiv:2402.02475*,
2024.
- [18] A. Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition
at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [19] S. Gao, T. Koker, O. Queen, T. Hartvigsen, T. Tsiligkaridis, and M. Zitnik. Units:
650 Building a unified time series model. *arXiv preprint arXiv:2403.00131*, 2024.
- [20] A. Garza and M. Mergenthaler-Canseco. Timegpt-1. *arXiv preprint arXiv:2310.03589*,
2023.
- [21] R. Gençay and M. Qi. Pricing and hedging derivative securities with neural networks:
Bayesian regularization, early stopping, and bagging. *IEEE transactions on neural*
655 *networks*, 12(4):726–734, 2001.
- [22] J. Hasanhodzic and A. W. Lo. Can hedge-fund returns be replicated?: The linear
case. *The Linear Case (August 16, 2006)*, 2006.
- [23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition.
In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
660 pages 770–778, 2016.

- [24] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [25] S. Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.
- [26] Z. Hu, Y. Zhao, and M. Khushi. A survey of forex and stock price prediction using deep learning. *Applied System Innovation*, 4(1):9, 2021.
- [27] R. J. Hyndman and Y. Khandakar. Automatic time series forecasting: the forecast package for r. *Journal of statistical software*, 27:1–22, 2008.
- [28] R. J. Hyndman and A. B. Koehler. Another look at measures of forecast accuracy. *International journal of forecasting*, 22(4):679–688, 2006.
- [29] R. J. Hyndman, A. B. Koehler, R. D. Snyder, and S. Grose. A state space framework for automatic forecasting using exponential smoothing methods. *International Journal of forecasting*, 18(3):439–454, 2002.
- [30] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93, 1938.
- [31] Z. Lan. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [32] S. Li, X. Jin, Y. Xuan, X. Zhou, W. Chen, Y.-X. Wang, and X. Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in neural information processing systems*, 32, 2019.
- [33] Y. Liang, H. Wen, Y. Nie, Y. Jiang, M. Jin, D. Song, S. Pan, and Q. Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6555–6565, 2024.
- [34] Y. Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [35] Y. Liu, H. Zhang, C. Li, X. Huang, J. Wang, and M. Long. Timer: Transformers for time series analysis at scale. *arXiv preprint arXiv:2402.02368*, 2024.

- [36] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- [37] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. M5 accuracy competition: Results, findings, and conclusions. *International Journal of Forecasting*, 38(4):1346–1364, 2022.
- [38] F. J. Massey Jr. The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78, 1951.
- [39] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- [40] J. A. Miller, M. Aldosari, F. Saeed, N. H. Barna, S. Rana, I. B. Arpinar, and N. Liu. A survey of deep learning and foundation models for time series forecasting. *arXiv preprint arXiv:2401.13912*, 2024.
- [41] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- [42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [43] K. Rasul, A. Ashok, A. R. Williams, A. Khorasani, G. Adamopoulos, R. Bhagwatkar, M. Biloš, H. Ghonia, N. V. Hassen, A. Schneider, et al. Lag-llama: Towards foundation models for time series forecasting. *arXiv preprint arXiv:2310.08278*, 2023.
- [44] A. Roberts, C. Raffel, K. Lee, M. Matena, N. Shazeer, P. J. Liu, S. Narang, W. Li, and Y. Zhou. Exploring the limits of transfer learning with a unified text-to-text transformer. *Google, Tech. Rep.*, 2019.
- [45] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.

- [46] V. Sanh. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [47] N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 720 [48] E. SLUTSKY. on a case where the law of large numbers applies to mutually dependent quantities. the extension of a theorem by magoichirô watanabe. *Tohoku Mathematical Journal, First Series*, 28:26–32, 1927.
- [49] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- 725 [50] Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- [51] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [52] I. Sutskever. Sequence to sequence learning with neural networks. *arXiv preprint arXiv:1409.3215*, 2014.
- 730 [53] I. Svetunkov, N. Kourentzes, and J. K. Ord. Complex exponential smoothing. *Naval Research Logistics (NRL)*, 69(8):1108–1123, 2022.
- [54] P. Tang and X. Zhang. Mtsmae: Masked autoencoders for multivariate time-series forecasting. In *2022 IEEE 34th International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 982–989. IEEE, 2022.
- 735 [55] S. J. Taylor and B. Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- [56] P. Temin. Price behavior in ancient babylon. *Explorations in Economic History*, 39(1):46–60, 2002.
- 740 [57] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

- [58] A. Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 745 [59] Q. Wen, J. Gao, X. Song, L. Sun, H. Xu, and S. Zhu. Robuststl: A robust seasonal-trend decomposition algorithm for long time series. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 5409–5416, 2019.
- [60] C. Wimmer and N. Rekabsaz. Leveraging vision-language models for granular market change prediction. *arXiv preprint arXiv:2301.10166*, 2023.
- 750 [61] H. Wold. *A study in the analysis of stationary time series*. PhD thesis, Almqvist & Wiksell, 1938.
- [62] G. Woo, C. Liu, A. Kumar, C. Xiong, S. Savarese, and D. Sahoo. Unified training of universal time series forecasting transformers. *arXiv preprint arXiv:2402.02592*, 2024.
- 755 [63] Q. Xie, W. Han, Y. Lai, M. Peng, and J. Huang. The wall street neophyte: A zero-shot analysis of chatgpt over multimodal stock movement prediction challenges. *arXiv preprint arXiv:2304.05351*, 2023.
- [64] L. Xue. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*, 2020.
- 760 [65] X. Yu, Z. Chen, Y. Ling, S. Dong, Z. Liu, and Y. Lu. Temporal data meets llm—explainable financial time series forecasting. *arXiv preprint arXiv:2306.11025*, 2023.
- [66] G. U. Yule. On the time-correlation problem, with especial reference to the variate-difference correlation method. *Journal of the Royal Statistical Society*, 84(4):497–537, 1921.
- 765 [67] G. U. Yule. Vii. on a method of investigating periodicities disturbed series, with special reference to wolfer’s sunspot numbers. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 226(636-646):267–298, 1927.
- 770 [68] B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

- [69] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115, 2021.
- [70] T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International conference on machine learning*, pages 27268–27286. PMLR, 2022.

775

A. HISTORY OF TIME-SERIES FORECASTING

A.1 Brief History of time-series prediction in Finance

When the agricultural revolution took place (circa 10,000 BC), for the first time in history, humans have started producing more food and items than they required just to survive. This gave birth to the concept of "storing" things and eventually trading them for other things which were perceived to be of greater value. When sufficient amount of people started transacting goods in this manner, the concept of "market price" took hold. It did not take long before humans realized that they didn't have to work in the field the whole day to amass wealth, but that they could do it through trade alone. They could buy items from one person at a low price and sell to another at a high price (arbitrage) or buy it now at a low price and hopefully sell later at a high price (asset appreciation). Opportunities to exploit the first option became more rare and less profitable as more and more people started engaging in trade and eliminating these arbitrage opportunities. However, the 2nd option has captured the imagination of traders ever since 10,000BC as no one seemed to be able to tell with certainty whether the price of a good would indeed go up in the future or not. As humans learned to write and keep records, they started noting down the prices of goods that were being traded along with the time when those prices were prevailing. We have records of ancient Babylonians keeping historic records of end-of-month barley, dates, cuscutea, sesame, cardamom and wool prices on clay tablets spanning some 400 years [22]. Investigating these prices, we can see that they follow a random walk pattern with exogenous shocks (such as the death of Alexander the Great) [56] and they seem almost impossible to predict with a naked eye, however, it is not impossible that someone might have tried and unknowingly became the first person in history to work on time-series prediction in the field of Finance.

Over the centuries, methods to solve the problem of time-series prediction in finance seem to have developed universally across many different cultures. Christopher Kurz, a 16th century merchant from Antwerp, thought about the idea of "trend" and "seasonality" in agricultural prices and claimed he could predict prices 20 days in advance. In 18th century Japan, Munehisa Honma - a man dubbed "God of the markets" by his contemporaries, developed a principle stating that when prices become extremely high, they must come

down again which we know today to be the principle of "mean reversion". In the late 19th century, Charles Dow, co-founder of Wall Street Journal and Dow Jones Company⁵⁸ popularized and coined the term "technical analysis" (a collection of methods that employ math and statistics to predict future prices) which later evolved into quantitative analysis - which is used today in institutions engaged in future price prediction.

A.2 Brief History of modern time-series prediction methods

Statistical models

Although the field of future price prediction has moved far beyond simply considering just the historical prices⁵⁹, there has been notable work on the general field of univariate⁶⁰ time-series prediction. It is hard to pinpoint the exact start of modern time-series prediction, but I think a good candidate would be Udny Yule's⁶¹ 1927 [67] paper which is considered first to have used simple autoregressor model (AR)⁶² to predict Wolfer's Sunspot numbers⁶³. Even though Slutsky (1927) [48] and Udny (1921) [66] demonstrated the use of Moving Average as a way to show trend and cyclicity in time-series data, Herman Wold's work (1938)⁶⁴ [61] indirectly invented the Auto Regressive Moving Average (ARMA) model⁶⁵. ARMA model is theoretically sound if the time-series being predicted is stationary⁶⁶, which is not the case with all time-series. In 1970, George Box and Gwilym Jenkins introduced probably the most famous time-series prediction model of all times: ARIMA [5]. ARIMA(p, i, q) model is an ARMA(p, q) model that employs the method of

⁵⁸Dow Jones Company is the publisher of the prominent DOWJ stock market index

⁵⁹Quantitative analysis now-adays considers wide palette of data and information

⁶⁰Univariate time-series prediction is a practice of using only the past values of a certain time-series in order to predict its future value(s).

⁶¹Udny studied at UCL under a well-known (and controversial) professor Karl Pearson - father of mathematical statistics and generally one of the most well-known personas in statistics.

⁶²An AR(p) time-series prediction model predicts future values based on a linear combination of the previous p observations in the series, assuming that the current value depends on its own past values and a stochastic error term.

⁶³Certain time-series data from the field of Astronomy.

⁶⁴Wold's Decomposition Theorem states that any stationary time series can be decomposed into a deterministic part (AR component) and a stochastic part (MA component)

⁶⁵The MA(q) part of the ARMA(p, q) model forecasts future values by modeling the current value as a linear combination of the past q error terms (shocks) and a stochastic error term.

⁶⁶A time series is said to be stationary if its statistical properties, such as mean, variance, and autocorrelation, remain constant over time

differencing⁶⁷ on the date before training and prediction.

Another direction of innovation in the field of univariate time-series prediction was towards "Exponential Smoothing" (ES). Concept of ES originates from the work of Brown 1956 [6] and has been thoroughly built upon since. Most notable invention was by Hyndman et al. in 2002 [29] - implementing exponential smoothing within a state-space framework, giving us the ETS⁶⁸ model. Most recent culmination in the ES family of models was by Ivan Svetunkov in 2022 [53] with the CES model⁶⁹. Final honorable mention among the time-series prediction models of "statistical" nature is the "Theta" model. Vassilis Assimakopoulos, and Konstantinos Nikolopoulos (2000) [2] Introduced the Theta model which decomposes a time-series into its' cyclical and trend components which are then forecasted separately and finally combined⁷⁰.

Machine Learning (ML) models

As the hardware technology advanced in the 21st century, so did the time-series models as the scientists have finally gotten machines with enough computational power to run more complex algorithms. The most famous of these algorithms being the Artificial Neural Network (ANN). Although the idea (in its primitive form⁷¹) dates back all the way to McCulloch and Pitts (1943) [39], ANN first real ancestor was the "Perceptron"⁷², invented by a psychologist - Frank Rosenblatt in 1958 [45]. The perceptron had ability to take continuous inputs and proper learning algorithm - known even to this day as the "Rosenblatt algorithm". Through the rest of the 20th century, there have been general

⁶⁷Differencing is a technique in time-series analysis used to transform a non-stationary series into a stationary one by deducting the preceding observation from each observation. The order "i" of the differencing refers to the ammount of times this procedure is applied to a time-series.

⁶⁸ETS stands for Error Trend Seasonality. ETS is a family of 18 models. ETS model can account for errors being additive or multiplicative, and trend or seasonality being none, additive or multiplicative. ETS models with no multiplicative errors or seasonality have their equivalent withing the ARIMA family of models

⁶⁹CES stands for Complex Exponential Smoothing model. CES method models the time-series as a series of complex numbers where the real part is the prediction and imaginary parts are errors. CES models can also be expressed in state-space form to adress seasonality. CES advantages are that it is flexible and has been empirically shown to perform well.

⁷⁰These components are called "theta lines". Theta model can use exponential smoothing or an AR model to predict these two lines. Predictions are finally combined using weighted average where weights could be equal or adjusted to favor either the trend line or the cycle line.

⁷¹The original vision of an ANN was based on simple binary inputs with fixed thresholds activation - a far cry from modern ANNs

⁷²Perceptron can be thought of as a single-neuron ANN. ANNs are comprised of one or more neuron layers, each layer containing one or more neurons.

periods of great interest and disinterest in ANNs for many reasons, but ANNs finally came on their own in the 21st century when computers become powerful enough to be able to train larger versions of these models⁷³ in reasonable time⁷⁴. Naturally it was not long before people realized ANN-type models can also be used for time-series prediction⁷⁵ which led to many variations of the ANN architecture to be invented. Bayesian regularization ANN has been used in Finance [21]. Another successful variation of ANN is the Radial Basis Function ANN (RBFANN).

Another

Chen and Guestrin (2016) [8] introduced the Xgboost model which has

Deep Learning models

⁷³Size of an ANN matters when dealing with very complex tasks (tasks with many input features).

⁷⁴Another reason for explosion in popularity was the beginning of the internet era and sudden availability of large ammounts of data which are needed for ANN training

⁷⁵It is important to note that univariate time-series prediction ANNs are simply regression ANNs where lags of the data we're trying to predict are explanatory features (inputs) to the model

B. TIMEGPT-1

C. LAG-LLAMA

D. METHODOLOGY

D.1 Time-series prediction

Prediction error

Prediction error is the difference between the value a model predicts and the actual (ground-truth) value.

MDA

MDA is time-series specific. Since time-series data is inherently ordered (as opposed to regular regression data), we can measure the direction in which a time-series is moving at each time-step. Hence we can compare the direction of movement of our predicted time-series against the actual time-series and we can calculate in what percentage of cases did the two time-series move in the same direction (up or down).

MES

Reason why MDA doesn't fully capture the model's ability to predict direction of underlying asset value movement is because the models are working with returns time-series, therefore MDA measures whether the model accurately predicts the directional change of returns and not the actual underlying values. Imagine scenario where $y_{\text{actual}}^T = 1\%$, $y_{\text{actual}}^{T+1} = 0.5\%$ and $y_{\text{predicted}}^T = 0.7\%$, $y_{\text{predicted}}^{T+1} = -0.5\%$. In this case MDA would account this as a correct directional guess as the model predicted y would go down in period $T+1$, and y actually did go down at $T+1$. However, if we inspect this case in greater detail, at the moment $T+1$, the value of the underlying asset which y represents, went up at $T+1$, as the return in that period is still positive, however the model actually implicitly predicted the price to go down. This is a nuance which MDA doesn't capture, but MES does.

D.2 Data

Data sourcing

Exceptions:

- Daily CHAPS data not available 2018-01-01 to 2020-01-01.
- Weekly CHAPS data not available 2017-01-01 to 2022-01-01 and 2015-01-01 to 2020-01-01.
- Weekly CHAPS data used was available only 2020-01-01 to 2024-01-01.
- Monthly CHAPS and Exchange rate data was not available in sufficient quantity to conduct experiment.

Partial autocorrelation

Partial autocorrelation measures the correlation between a time series and its lagged values, while controlling for the influence of intermediate lags. In contrast, autocorrelation simply measures the correlation between the time series and its lagged versions without accounting for the effects of other lags. Partial autocorrelation provides a clearer picture of the direct relationship between observations and their lagged counterparts by isolating the specific influence of each lag. This makes it particularly useful for detecting seasonality in time-series data because it helps identify the direct impact of observations at seasonal intervals. By highlighting the most relevant lags for seasonality, partial autocorrelation allows for a more precise analysis of recurring patterns, which is harder to achieve using autocorrelation alone.

D.3 Experiment design

ARIMA

Originally, this package is for R. I used the equivalent Python implementation: <https://github.com/alkaline-ml/pmdarima>.

Naive Simple Autoregressor (NSA)

905 Given a training time-series sequence of length n : $T = \{y_1, y_2, \dots, y_n\}$, and a prediction horizon of length h , the prediction of NSA will be: $P = \{y_n \text{ repeated } h \text{ times}\}$, i.e. NSA predicts the future value(s) the same as the last value.

Loss

Loss in the context of fine-tuning is a measure of model's performance on the current
910 epoch it has trained on. Lower the loss, the better the model's performance in that epoch.

Benchmark