

COUB01  
F014621

# Predicting Inflation With Machine Learning

by

Ben M. Taylor

A Project Report

Submitted in partial fulfilment  
of the requirements for the award of

B.Sc.

in

Computer Science

of

Loughborough University

2nd December 2024

Copyright 2024 Ben M. Taylor

# Abstract

In October 2022, the UK hit an inflation rate of 11.1%, the country's highest in over 40 years. Now more than ever, the ability to accurately predict inflation and other financial indicators is a crucial skill required by the government and the individual to prepare themselves for the future financially. In a time where Artificial Intelligence and Machine Learning are ever flourishing, it is only natural to attempt to use these tools at our disposal to predict and combat the issues we face.

In this paper I will attempt to predict inflation through the use of machine learning eventually presenting my findings and evaluations.

**Keywords:** Inflation, Artificial Intelligence, Machine Learning

---

# Acknowledgements

I would like to thank my supervisor Mohamad Saada for his help and guidance throughout this project.

# Contents

<b>Abstract</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>4</b>
1.1 Motivation . . . . .	4
1.2 Aims and Objectives . . . . .	4
1.2.1 Aims . . . . .	4
1.2.2 Objectives . . . . .	5
<b>2 Literature Review</b>	<b>6</b>
2.1 Motivation . . . . .	6
2.2 Available Literature and Context . . . . .	6
2.2.1 Financial Indicator Prediction Papers . . . . .	7
2.2.2 Machine Learning Papers . . . . .	8
2.2.3 Inflation Papers . . . . .	9
2.3 Problem Domain . . . . .	9
2.3.1 UK Inflation . . . . .	10
2.3.2 Machine Learning Models . . . . .	10
2.3.3 Regression Metrics and Model Evaluation . . . . .	14
2.4 Summary and Conclusion . . . . .	15
<b>3 Main chapters</b>	<b>16</b>
<b>4 Data Selection and Pre-processing</b>	<b>17</b>
4.1 Data Selection . . . . .	17
4.2 Data Pre-processing . . . . .	17
4.2.1 Cleaning the Data . . . . .	18
4.2.2 Granger Causality . . . . .	18
<b>5 Conclusion</b>	<b>20</b>



# List of Figures

2.1	FTS Forecasting Methods . . . . .	7
2.2	ML arXiv articles per year . . . . .	9
2.3	Common Activation Functions . . . . .	12
4.1	The function written to Granger test the data. . . . .	19

---

# List of Tables

# References

- [1] Steven B Achelis. Technical analysis from a to z, 2001.
- [2] Oguz Akbilgic, Hamparsum Bozdogan, and Mehmet Erdal Balaban. A novel hybrid rbf neural networks model as a forecaster. *Statistics and Computing*, 24:365–375, 2014.
- [3] Mariette Awad and Rahul Khanna. *Machine Learning*, pages 1–18. Apress, Berkeley, CA, 2015.
- [4] World Bank. World development indicators, 2023.
- [5] David Beckett. Consumer price inflation, uk: October 2023, Nov 2023.
- [6] Leo Breiman. Bagging predictors. *Machine learning*, 24:123–140, 1996.
- [7] Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- [8] Lijuan Cao and Francis E.H. Tay. Financial forecasting using support vector machines. *Neural Computing and Applications*, 10(2):184 – 192, 2001. Cited by: 272.
- [9] G. Ciaburro, V.K. Ayyadevara, and A. Perrier. *Hands-On Machine Learning on Google Cloud Platform: Implementing smart and efficient analytics using Cloud ML Engine*. Packt Publishing, 2018.
- [10] John Harold Clapham. *The bank of England*. Number 7. CUP Archive, 1939.
- [11] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20:273–297, 1995.
- [12] Jeff Dean, David Patterson, and Cliff Young. A new golden age in computer architecture: Empowering the machine-learning revolution. *IEEE Micro*, 38(2):21–29, 2018.
- [13] Jeffrey Dean. The deep learning revolution and its implications for computer architecture and chip design. *ArXiv*, abs/1911.05289, 2019.



- [14] The Economist. Our big mac index shows how burger prices are changing.
- [15] Brigid Francis-Devine. Rising cost of living in the uk. <https://commonslibrary.parliament.uk/research-briefings/cbp-9428/>, 2023. [Accessed 17-12-2023].
- [16] C. W. J. Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37(3):424–438, 1969.
- [17] Nicholas Geoffrey Lemprière Hammond. *Alexander the Great*. Chatto & Windus London, 1981.
- [18] Bruce E. Hansen. Threshold effects in non-dynamic panels: Estimation, testing, and inference. *Journal of Econometrics*, 93(2):345–368, 1999.
- [19] Cynthia Harrington. Fundamental vs. technical analysis, 2003.
- [20] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- [21] Ryszard Stanislaw Michalski, Jaime Guillermo Carbonell, and Tom M Mitchell. *Machine learning: An artificial intelligence approach*. Springer Science & Business Media, 2013.
- [22] John J Murphy. *Technical analysis of the financial markets: A comprehensive guide to trading methods and applications*. Penguin, 1999.
- [23] Eir Nolsøe. Bank of England’s Covid money-printing spree ‘drove up inflation’ — telegraph.co.uk. <https://www.telegraph.co.uk/business/2023/04/18/bank-of-england-covid-money-printing-drove-up-inflation/#:~:text=During%20the%20pandemic%2C%20the%20Bank,to%20a%20record%20%C2%A3895bn.,2023>. [Accessed 17-12-2023].
- [24] Ceyda Oner. Inflation: Prices on the rise, Jul 2019.
- [25] Michael Parkin. *Inflation*, pages 1–10. Palgrave Macmillan UK, London, 2016.
- [26] Barbara Rockefeller. *Technical analysis for dummies*. John Wiley & Sons, 2019.
- [27] F. Rosenblatt. The perceptron - a perceiving and recognizing automaton. Technical Report 85-460-1, Cornell Aeronautical Laboratory, Ithaca, New York, January 1957.

- [28] Ehsan Sadrfaridpour, Korey Palmer, and Ilya Safro. Aml-svm: Adaptive multilevel learning with support vector machines, 2020.
- [29] Tom Seegmiller. 10 metrics to measure the financial efficiency of your organization, Sep 2023.
- [30] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [31] Yajiao Tang, Zhenyu Song, Yulin Zhu, Huaiyu Yuan, Maozhang Hou, Junkai Ji, Cheng Tang, and Jianqiang Li. A survey on machine learning models for financial time series forecasting. *Neurocomputing*, 512:363–380, 2022.
- [32] Michael C Thomsett. *Getting started in fundamental analysis*. John Wiley & Sons, 2006.
- [33] Ruey S Tsay. *Analysis of financial time series*. John wiley & sons, 2005.
- [34] Ahmed. S. Wafi, Hassan Hassan, and Adel Mabrouk. Fundamental analysis models in financial markets – review study. *Procedia Economics and Finance*, 30:939–947, 2015. IISES 3rd and 4th Economics and Finance Conference.
- [35] Zhongyuan Wei. *A SVM approach in forecasting the moving direction of Chinese stock indices*. Lehigh University, 2012.
- [36] Daniela Hristova \* Wojciech W. Charemza and Peter Burridge. Is inflation stationary? *Applied Economics*, 37(8):901–903, 2005.
- [37] Jeffrey Marc Wooldridge. *Introductory Econometrics: A Modern Approach*. ISE - International Student Edition. South-Western, 2009.
- [38] G. Udny Yule. Why do we sometimes get nonsense-correlations between time-series?—a study in sampling and the nature of time-series. *Journal of the Royal Statistical Society*, 89(1):1–63, 1926.

# Chapter 1

## Introduction

### 1.1 Motivation

Prices for goods and services are ever-changing, constantly affecting individuals, organizations, and governments. These clients would all benefit from the ability to predict the rise and fall of prices as it can impact their choices to spend, invest, or save. Ultimately allowing them to make the most out of the resources they have.

However, the value of inflation is an incredibly difficult indicator to predict with even establishments like the Bank of England, who have direct control over interest rates, failing to correctly forecast it. With the ability to accurately predict inflation having so many stakeholders and huge organizations working on the problem, I am under no illusion that an undergraduate like myself will be able to find a solution better than all those before me. Instead, I aim to understand which machine learning (ML) methods produce the best results when faced with this task and potentially find out why this is the case. This project's results will ultimately contribute to the ongoing research into inflation forecasting and help clarify which ML methods may be most suitable.

### 1.2 Aims and Objectives

#### 1.2.1 Aims

This project aims to create multiple machine learning models and compare their predictive abilities when it comes to UK inflation. The findings will then be presented in this report, outlining the advantages and disadvantages of each tested model. A conclusion will then be made about which models performed best and in which scenarios certain models should be used over others.

### 1.2.2 Objectives

This project can be broken down into a list of objectives that not only provide a strong path to follow to complete the project but also a way to evaluate the project's success post-completion. I have placed these objectives into phases corresponding to the project's work plan.

1. Phase 1: Research

- (a) Conduct a literature review to understand the current landscape of inflation forecasting.
- (b) Research prominent models in the literature.

2. Phase 2: Source Data

- (a) Source the data to be used in the models.
- (b) Evaluate the usefulness and appropriateness of the data.
- (c) Clean the data: dropping and retaining specific variables.

3. Phase 3: Creating Models

- (a) Choose at least 3 appropriate models to use.
- (b) Develop and tune the models.
- (c) Train the models on the dataset.

4. Phase 4: Evaluation and Report

- (a) Evaluate the models using statistical tests.
- (b) Conclude the findings.
- (c) Present the findings in this final report.

## Chapter 2

# Literature Review

### 2.1 Motivation

Embarking on a literature review before developing our project offers numerous benefits. Understanding existing knowledge in Machine Learning, specifically when used to predict financial indicators, helps to contextualise our research, positioning it within the existing field. Reviewing previous literature also provides the benefits of identifying gaps in current research and finding supporting arguments that can guide our work and help us to avoid, as much as possible, redundancy in our and others' works. Having completed the literature review, we should have a strong foundation to start and guide our project.

### 2.2 Available Literature and Context

There is certainly a strong monetary incentive to produce research on how best to predict financial indicators. The correct predictions can not only allow organizations and individuals to profit greatly but also to avoid loss. This results in a myriad of papers being written, experimenting with a variety of techniques to predict future values, most of which we can learn from to help structure our models.

This report's topic focuses on the prediction of inflation through the use of machine learning. To accomplish this we can view papers predominantly addressing two types of topics. The first type is papers that focus on the topic of predicting inflation or other economic indicators and time series. The second type of papers we can research are ones that deal with different machine learning techniques. Additionally, it is pertinent to survey the current literature on inflation: its causing factors, effects, and significance.

### 2.2.1 Financial Indicator Prediction Papers

According to "Analysis of Financial Time Series" by Ruey S. Tsay "Financial time series analysis is concerned with the theory and practice of asset valuation over time." [33] There are many financial time series (FTS for short) prediction methods both theoretical and practical that have attracted attention, the taxonomy of which is shown in figure 2.1. The predominant analysis strategies for predicting financial market behaviour are fundamental analysis and technical analysis[19].

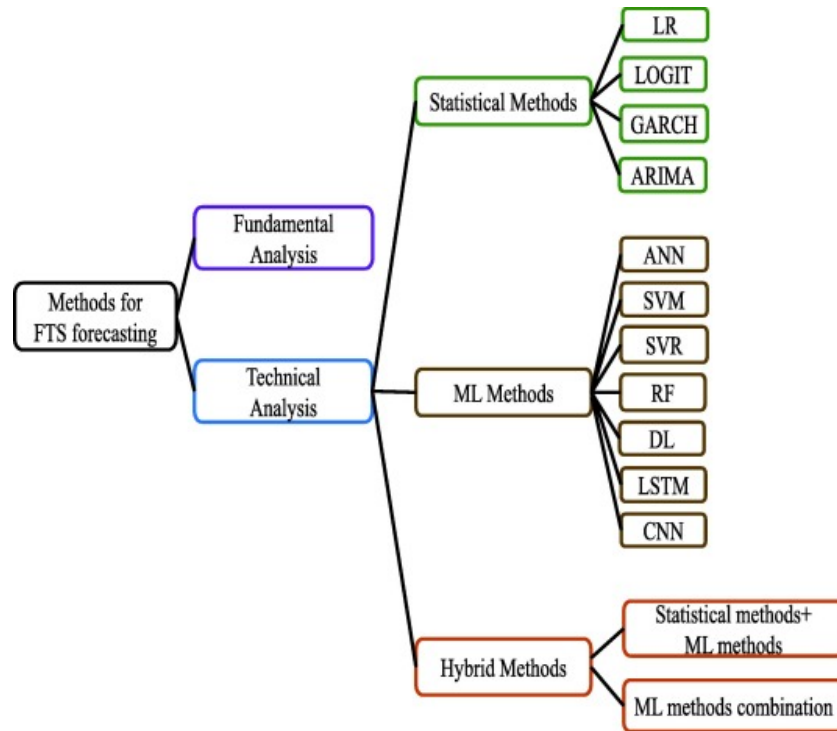


Figure 2.1: FTS Forecasting Methods.

Figure from page 3 of 'A survey on machine learning models for financial time series forecasting by Yajiao Tang et al.' [31]

#### Fundamental Analysis

Fundamental analysis[32] attempts to measure the intrinsic value of an asset by looking at current market and economic conditions. Additionally, fundamental analysis frequently makes use of techniques - such as sentiment analysis - that often deal with unstructured data. The success of fundamental analysis often relies on the financial efficiency of the target[34], which according to Tom Seegmiller is defined as "... how successful your organization is at turning expenses into revenue" [29].

### Technical Analysis

Technical analysis[1] attempts to identify opportunities and predict investments by viewing movements and trends in market data alongside using a variety of technical indicators. Unlike fundamental analysis, technical analysis does not take into account many of the same fundamentals that can help indicate an asset's current value such as quarterly revenue. This is partially because it is often argued that technical indicators such as inflation or a stock's value are already priced according to the fundamentals that cause or contribute to them[22]. From this, we can come to the understanding that while fundamental analysis is the idea of looking at the current factors affecting an asset and using them to evaluate to asset's true value; Technical analysis is built upon the idea that past performance can predict future performance. Traditionally, technical analysis has relied heavily on statistical models to forecast the future performance of assets[26]. Furthermore, the application of using past values to predict future values has been widely implemented for years, with one of the earliest uses of autoregressive models being used to predict time series created by U.G.Yule in the 1920s[38]. However, with the increase of big data and the internet, ever-larger amounts of financial predictive data are continually being produced. Nowadays, simple statistical models may struggle to produce accurate future predictions when faced with big data sets containing complex characteristics[2].

### 2.2.2 Machine Learning Papers

This brings us to machine learning algorithms[21]. According to Mariette Awad et al. machine learning "is a branch of artificial intelligence that systematically applies algorithms to synthesize the underlying relationships among data and information"[3]. Currently, there is a massive abundance of fresh machine learning papers constantly being produced in the field. [12]

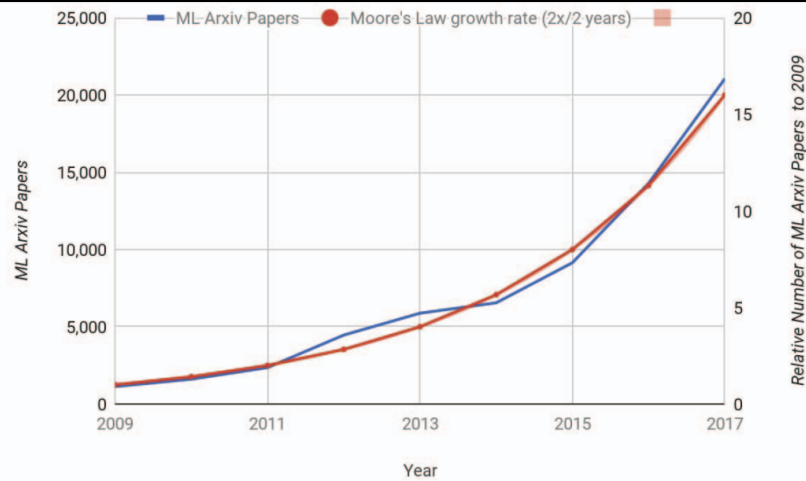


Figure 2.2: ML arXiv articles per year.

Figure from page 4 of 'A New Golden Age in Computer Architecture: Empowering the Machine Learning Revolution' by Jeff Dean, David Patterson, and Cliff Young[12]

As you can see in figure 2.2 articles on machine learning posted to arXiv (an archive for scholarly articles) have more than doubled every two years. Additionally, in 2018 the number of articles released reached 100 per day, summing to more than 33,000 by the end of the year. The number of articles released has steadily continued to increase in the years since[13]. Naturally, to read this many articles is impossible, however, the sheer quantity bodes well for this project as it means there will be plenty of guidance on how best to select and develop our predictive models.

### 2.2.3 Inflation Papers

According to Ceyda Oner at the International Monetary Fund "Inflation is the rate of increase in prices over a given period of time." [24] Often inflation is used to broadly indicate a countries of the global state of price fluctuation, however, inflation can also be used for certain services, goods, or food. Inflation affects everyone leading to an abundance of different papers linking to the topic. The papers focusing on inflation cover a variety of topics such as inflation's various effects, forecasting inflation, the causes and trends of inflation, and more. [25]

## 2.3 Problem Domain

This section of the literature review will go more in-depth into the area that this report will tackle as well as which methods could potentially be used.



### 2.3.1 UK Inflation

Inflation has been around since money has been used with one of the earliest recordings of inflation being caused by the death of Alexander the Great in around 300 BC[17]. Because of inflation's pervasiveness and its intrinsic link to both macro and microeconomics, there is a surplus of literature surrounding the subject. This report will focus on inflation in the UK as opposed to global inflation. Several indicators can be used to measure inflation, the most common of which are Consumer Price Index (CPI), Consumer Price Index with Housing (CPIH), and Retail Price Index (RPI). There are also more novel measurements of inflation as well such as viewing the Big Mac index[14] and recording its change over time.

In the UK, inflation is the responsibility of the Bank of England who set monetary policy eight times a year with the goal of controlling and stabilising inflation. The Bank of England (BoE) calculates inflation using their own in-house models[10]. The BoE are not transparent with the specifics of the models or exactly how they calculate inflation, however, they state their models do factor in market expectations. In October 2022 inflation rates reached a peak of 11.1% the highest in over 40 years. Yet by October 2023, the annual rate was the lowest since October 2021 at 4.7%[5]. Granted this is still higher than the targeted 2% imposed on the BoE by the government. According to a research briefing published by the House of Commons Library, the main causes for the extremely high current inflation were: Covid-19 lockdowns causing supply chain disruptions and Russia's invasion of Ukraine causing increased energy prices due to the UK's prior reliance on Russian fuel[15]. Additionally, several pundits have suggested the BoE's slow reaction to combating high inflation rates and their increased money-printing during the pandemic[23]. The main way that the BoE combats inflation is through the manipulation of interest rates. The general premise is that as inflation rates grow, the growth of inflation should slow and inflation should eventually decrease. This is because increased interest rates mean the overall spending in the economy lessens as money is instead being spent on the increased interest.

Through this project we hope to be able to forecast future inflation through the use of machine learning models.

### 2.3.2 Machine Learning Models

By utilising machine learning techniques in financial forecasting we can endeavor to improve upon the performance of traditional statistical models. Generally, the goal of FTS forecasting can be placed into two main categories: 1. Price prediction 2. Price movement prediction (this includes volatility predictions) These

two goals also reflect two types of machine-learning problems: 1. Regression 2. Classification. This paper will focus mainly on the price prediction/regression categories regarding inflation. This means that we will aim to predict future values of inflation as opposed to whether inflation will increase or decrease. Naturally, the areas we study will be with this regression problem in mind. We shall now cover some of the potential ML forecasting models that can aid us in predicting inflation.

### **Artificial Neural Networks**

An artificial neural network is a machine learning model that is made up of an interconnected group of nodes (also known as neurons) organised into layers. The model's inspiration stems from how neurons in the human brain interact with one another. In 1957 Frank Rosenblatt invented the perceptron, one of if not the first implementations of an artificial neural network [27].

ANNs are composed of 3 types of layers: the input, hidden, and output layers. The input layer receives input data and passes it through to the first hidden layer. Hidden layers receive weighted inputs, perform an activation function on said input, and then pass the new data to the next layer. The output layer receives data from the final hidden layer and then produces the resulting prediction. The neurons within an ANN can either be excited or inhibited. Neurons are connected between layers and the strength of these connections (the weight) is decided by how excited or inhibited a neuron is. Each neuron in the hidden and output layers contains biases and activation functions (with the activation function in the output layer typically being different from the one in the hidden layer). Activation values are passed from node to node through the connections in the network. When a neuron receives the activation value, it sums and modifies it based on the neuron's activation function and bias. The predicted results can then be compared to the true values and the weights and biases of the network are updated accordingly. Artificial neural networks can be modified with a variety of techniques to alter their accuracy. One of these alterations is changing the activation function of the neurons. Figure 2.3 shows some common activation functions. Other commonly used activation functions include sigmoid, Gaussian, and leaky ReLU.




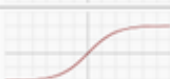
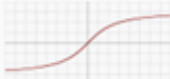




Name	Plot	Equation
Identity		$f(x) = x$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$
ArcTan		$f(x) = \tan^{-1}(x)$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) <sup>[2]</sup>		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) <sup>[3]</sup>		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$

Figure 2.3: Common Activation Functions [9]

ANNs have many advantages:

- Thanks to the many interconnected neurons, ANNs have strong learning capabilities[18].
- ANNs do not have a fixed structural equation making them very adaptable.
- ANNs can be finely tuned with many changes to the network in order to find the best-fitting model. For example, changing the number of hidden layers, the number of nodes in a layer, the activation functions, and so on.

ANNs also come with some disadvantages:

- The higher complexity of ANNs means that they require more resources than traditional statistical models.
- Due to the nature of the hidden layers, ANNs can be hard to interpret.
- ANNs can be overfitted[30].

### Support Vector Regression

Created by Vladimir Vapnik and Alexey Chervonenkis in the 1960s and later built upon further by Vapnik et al. with the addition of the kernel trick and soft margin [11], Support Vector Regression and Support Vector Machines are supervised learning methods used for regression and classification respectively. Both models use non-linear mapping to transform the dimension of the input data. Then utilise a hyperplane in order to either best fit or categorize the data. The hyperplane for these models is found by using an  $\varepsilon$ -insensitive tube, meaning that any errors within the range of the tube are ignored. This is unlike a standard line of best fit that takes into account the  $\varepsilon$  (distance) of all points to the line, instead, only errors outside of the tube are considered pertinent. The hyperplane is then placed in a way in which the sum of all points outside of the tube is minimised. The points outside of the tube are called support vectors hence the name support vector regressions.

Advantages of Support Vector Regression:

- SVRs are simple and easy to implement as well as producing easily interpretable results.
- SVRs require less computational resources than other models.
- SVRs can maintain stability despite noisy input data thanks to the  $\varepsilon$ -incentive tube[35].

Disadvantages of Support Vector Regression:

- Deciding a suitable kernel function can cause difficulty [8].
- SVRs may struggle with big data[28].

### Random Forest

The first random forest (RF) algorithm was created by Tin Kam Ho in 1995[20] which was later developed upon by Leo Breiman[7]. The random forest model makes predictions by consulting multiple decision trees. Each tree is trained on a random subset of data taken from the training set, this is called bagging or bootstrap aggregation. The final prediction is an average taken from all of the trees' predictions.

Advantages of random forest:

- Random forest can prevent overfitting by combining the results of several weak learners instead of using one powerful learner[6].

- Rf models have good prediction accuracy as the result is an average making it unlikely to be an outlier.
- Rfs are stable as changes to the data set may affect one tree but are unlikely to affect many trees.

Disadvantages of random forest:

- Rfs suffer from increased training time. This is due to the fact that to make a prediction you need predictions from all of the trees to get an average.

### 2.3.3 Regression Metrics and Model Evaluation

Evaluating a model is extremely important not only to know the quality of your model's predictions but also in order to make improvements to the model to achieve a more desirable result. Evaluation metrics can be used to evaluate the performance of the model. Some useful metrics that can be applied to measure a model's performance are: Mean absolute error (MAE), Mean squared error (MSE), Mean absolute percentage error (MAPE), Root mean absolute error (RMAE), Normalised mean square error (NMSE), Root mean squared error (RMSE), Relative root mean squared error (RRMSE), Correlation coefficient of prediction (R) The most commonly used metrics are MSE, MAE, MAPE, and R.

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - X_i|$$

$n$  is the number of data points.  
 $Y_i$  is the  $i$ th true value.  
 $X_i$  is the  $i$ th predicted value.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i)^2$$

$n$  is the number of data points.  
 $Y_i$  is the  $i$ th true value.  
 $X_i$  is the  $i$ th predicted value.

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - X_i}{Y_i} \right|$$

$n$  is the number of data points  
 $Y_i$  is the true value.  
 $X_i$  is the predicted value.

$$R = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

$X_i$  and  $Y_i$  are the data points.  
 $\bar{X}$  is the mean of the x-value and  $\bar{Y}$  the mean of the y values.

FOR MSE, MAE, and MAPE a lower result indicates a more accurate model, and when a model has no error the value will be zero.  $R$  is always between -1 and 1, when  $R=0$  it indicates that there is no linear relationship between the values. If  $R$  is -1 then there is a perfect negative linear relationship and if  $R$  is 1 then there is a perfect positive linear relationship. These formulas can be used to understand the predictive skills of a model.

## 2.4 Summary and Conclusion

Through this literature review, we have covered a sample of the literature on machine learning methods, financial time series forecasting, and inflation. There is a multitude of literature available in these areas due to the potential monetary gain as well as the new and emerging technologies being explored in the fields. Therefore, it would be nearly impossible to extensively cover all of the relevant articles to this project. Instead, the main focus was on understanding the most common techniques used to predict FTS and which machine learning models are often applied. We also covered the UK's relationship with inflation, who is tasked with controlling it, and how they attempt to do so as well as some of the factors contributing to it.

Having completed the literature review we are now equipped with general knowledge of the problem of inflation, its causes, and how it is controlled. This will help us to select appropriate data to train and build our models on. Furthermore, due to reviewing current literature in the field of FTS forecasting and machine learning, we have the knowledge of which available models may be best suited to forecasting inflation as well as an understanding of how to evaluate the models we choose to use.

---

## Chapter 3

### Main chapters

# Chapter 4

## Data Selection and Pre-processing

### 4.1 Data Selection

The saying "garbage in garbage out" succinctly states the importance of selecting data for machine learning models. It does not matter how powerful a model is if the data selected is poor or inappropriate. Selecting data is a key step to building effective machine-learning models. Thankfully, there are plenty of large open-source data sets available online, yet picking an appropriate set may still present a challenge. The dataset that will be used for this project is the United Kingdom subset of "World Development Indicators"[4]. My reason for selecting this set is due to its large number of metrics, reliable provider, and the fact that it is open-source. The dataset is classified as public under the access to information classification policy. The World Bank is known for reliable and accurate data which should avoid issues of bias, insufficient data, or poor quality. The original UK dataset contains over 1400 metrics with data running from 1960-2022 including 2 metrics for inflation (CPI and GDP deflator (the rate of price change in the economy as a whole)). However, due to the all-encompassing nature of the dataset, it will need to be cleaned and many columns will need to be removed to improve the set's relevance.

### 4.2 Data Pre-processing

Often, datasets have several outstanding issues or properties that make them imperfect for use in a machine-learning model. Data pre-processing is the process of making changes to and cleaning a dataset before using it in a model.



### 4.2.1 Cleaning the Data

There were numerous outstanding issues with the original dataset taken from the World Bank. Several steps were taken to clean the initial dataset taken from the World Bank. An example of this is that almost all indicators in the set had null values from years where data was not collected. As part of the cleaning process, any indicators that met the following criteria were dropped from the dataset:

- Indicators with 20 or more years of missing data.
- Duplicate indicators.
- Indicators that were constant throughout the years (e.g. "secondary education duration")
- Indicators with 5 or fewer years of unique data.
- Indicators with no data in the last 5 years.

This reduced the number of indicators from 1458 to 396. Next all remaining null values were replaced with zeros to improve readability. Finally, a Granger causality test was carried out on the data.

### 4.2.2 Granger Causality

Granger causality is a statistical concept in economics used to show if time series A is useful at forecasting time series B. Clive Granger originally proposed the test in 1969 in the article "Investigating Causal Relations by Econometric Models and Cross-spectral Methods"[16]. The test only shows predictive causality and not true causality. Additionally, the test only provides information about forecasting ability and not the actual causal relationship. Another potential issue with the use of the Granger test in the context of inflation is that the test works best on stationary data, yet whether inflation is best treated as stationary or non-stationary data is currently inconclusive[36]. But for our purposes, these limitations are fine as we only want a loose idea on how useful our data will be for forecasting.

#### Granger Testing The Data

With two indicators X and Y, X causes Y if a series of tests on lagged values of X produce a p-value of less than 0.05. The closer the p-value is to zero the more likely it is for X to granger cause Y. Lagged values are values from a time series shifted forwards or backward in time. In the case of the Granger test, Jeffery Woolridge proposes that fewer lags should be used for annual data compared to quarterly

or monthly data in order to not lose degrees of freedom[37]. The Granger test was run on the remaining indicators comparing them each to CPI. The function ran 5 lags to see if data has a potential causal relationship with inflation within 5 years. All indicators that did not produce a p-value of less than 0.05 within the 5 lags were dropped from the dataset. We know that the data dropped does not Granger cause CPI so is unlikely to be useful in forecasting CPI. However, this does not mean that the retained data has a causal relationship with inflation. It simply means that at some point during the 5 lags a p-value $\leq$ 0.05 was produced and thus the data has the potential to be useful in forecasting CPI.

```
for label in Tdf.columns[3:]:#inflation and year metrics are in column 0-3
    print("\n#####")
    print(label)
    inflat1 = "Inflation, consumer prices (annual %)"
    ts_df = pd.DataFrame(columns=['Tdf[inflat1]', 'Tdf[label]'], data=zip(Tdf[inflat1], Tdf[label]))
    gc_res = grangercausalitytests(ts_df, 5)
    p_values = [round(gc_res[i+1][0]['ssr_chi2test'][1], 4) for i in range(5)]
    min_p_value = np.min(p_values)
    print("p:", min_p_value, "label:", label)
    if min_p_value > 0.05:
        Tdf = Tdf.drop([label], axis=1)
        print("dropped")
    else:
        print("granger caused")
        count += 1
    print("#####")
```

Figure 4.1: The function written to Granger test the data.

With this test complete the dataset has now dropped from 1458 indicators to 182.

---

**Chapter 5**

**Conclusion**

---

## Chapter 6

## References