

Time Series: Defining a Search Engine

Philipp Beer

September 28, 2021

1 Thesis: Time Series Search Engine

1.1 Purpose

The purpose of this thesis is to explore the possibility of creating a time series search engine.

1.2 Introduction

Time series is often described as "anything that is observed sequentially over time" which usually are observed at regular intervals of time [1]. They can be described as collection of observations that are considered together in chronological order rather than as individual values or a multiset of values. Their representation can be described as ordered pairs: $S = (s_1, s_2, \dots, s_n)$ where $s_n = (t_n, v_n)$. t_n can be a date, timestamp or any other element that defines order. v_1 represents the observation at that position in the time series.

Time series are utilized to analyze and gain insight from historic events/patterns with respect to the observational variable(s) and their interactions. A second area of application is forecasting. Here time series are utilized to predict the observations that occur in future under the assumption that the historic information can provide insight into the behavior of the observed variables.

Fu in their work [2] categorized time series research into (1) representation, (2) indexing, (3) similarity measure, (4) segmentation, (5) visualization and (6) mining. Research in these different fields started taking off in the second half of the 20th century. For example in [3] the authors worked on questions of representation via sampling the sampling of time series in 1969. All these different research areas always have to deal with the challenges that inhibit time series data. Generally datasets in this domain are large. Through this time series data incorporates the similar obstacles as high dimensional data, namely the "curse of dimensionality" [4] and requiring large computational efforts in order to generate insights. And as will be discussed in 1.2.1 further.

1.2.1 Applications

Time series are encountered everywhere. Any metric that is captured over time can be utilized as time series. Granularity can be used as descriptor for the sampling rate of a series or more general how often measurements for a particular metric are taken. This granularity has a tendency to increase as well. As example consumer electronics that capture health and fitness data can be mentioned. Or sensors which are utilized in the automotive industry or heavy machinery where they are employed to capture information for predictive maintenance applications.

In the financial industry time series are a very fundamental component of decision making, like the development of stock prices over time or financial metrics of interest. The same is true for macro economic information or metrics concerning social structures in society, etc.

In the medical field time series are also ubiquitous. Whether they relate to patient data like blood pressure. The bio statistics field utilizes electrographical data like electrocardiography, electroencephalography and many others. In more aggregate medical analysis like toxicology analysis of drug treatments for vaccine approvals they are utilized and in many forms of risk management, for example, population level obesity levels.

In engineering fields the utilization is often times similar to the above but it also require that information that is captured in time series is transferred between locations in an efficient manner. For example voice calls are required to be transferred between the participants in a fast manner and with minimized levels of noise in the data. Another interesting industrial example in the biomedical technology field is Neuralink which aims to implement a brain-machine-interface (BMI) utilizing mobile hardware like smartphones as the basis for its computation. Here a large of amount of time series data is generated which requires quick processing to generate real-time information. Musk describes a recording system with 3072 electrodes generating time series data [5] that is used to capture the brain information and visualized in real-time [6].

Time series data is paramount to a wide variety of areas, relating to many different fields. Looking at the trajectory it seems likely that going forward more time series data on a higher granularity will be generated. This in turn increases the need to be able process, analyze, compare and respond to the data with methods that are faster than today’s standard options.

1.2.2 Focus

In this work we focus on allowing the comparison of time series generated from different processes with different underlying properties and introduce a new measure of similarity based on the Fourier transformation. Our method does not require assumptions on the utilized time series and requires significantly less computaional resources for the execution of the comparison.

1.2.3 Organization of this thesis

The rest of this thesis is organized as follows:

1.2.4 TODO to be integrated

- refer to previous work on measures of similarity and outcome
- measure of similarity required
- challenges with time series (domains, granularity, length, outliers)
- area of signal processing interesting methods

1.3 Exisiting work

For our task of enabling a fast comparison between time series mainly two areas of research are relevant: (1) measures of similarity, and (2) time series representation.

1.3.1 Measuring similarity

In this work measures of similarity are differentiated into (1) raw similarity metrics between series, (2) metrics that are derived via feature extraction, and (3) similarity that is measured via the comparison of the time series representation.

In order to be able to describe the closeness of time series or multiple time series to each a measure for similarity is required. In the literature various general measures and corresponding computation methods can be found. Wang et al. reviewed time series measures and categorized the similarity measures into 4 categories: (1) lock-step measures, (2) elastic measures, (3) threshold-based measures, and (4) pattern-based measures.

Lock-step measures include the L_p -norms (Manhattan and Euclidean Distance) as well as Dissimilarity Measure (DISSIM). **Elastic measures** include metrics like Dynamic Time Warping (DTW) and edit distance based measures like Longest Common Subsequence (LCSS), Edit Sequence on Real Sequence (EDR), Swale and Edit Distance with Real Penalty. An example for **threshold-based measures** are threshold query based similarity search (TQuEST). And Spatial Assembling Distance (SpADe) is an example for pattern-based measures.

A more flexible category are the **elastic measures**. Their key advantage is the fact that comparison is applied on a one-to-many-basis allowing the comparison of regions from one series to regions of the other time series as well as comparing the all the points of one series to all the points of the other series.

- Euclidean Distance Euclidean Distance is the most widely used distance metric in the research of time series. (add list of papers here)
 - explain advantages
 - * linear complexity
 - mention shortcomings
 - * same length period
 - * handling of outliers and noise
 - * handling of stretching of series
 - * computational complexity
- Dynamic Time Warping
 - invented by [8] in 1994
 - warp series by computing the distance from one point to all other points in the other series and define a warped path that minimizes the distance

$$DTW(S_a, S_b) = \min \left\{ \sqrt{\sum_{k=1}^P \delta(\omega_k)} \right\} \quad (1)$$

- advantages: handles distortions, does not require same length ts
- disadvantages: outliers may create a false impression of similarity, computational complexity of $O(n^2)$ makes utilization for very long time series impractical and comparison with large sets of time series is also very time intensive
- pathological matchings
- Similarity through decomposition
 - introduce time series decomposition (reference in [1])
 - trend and seasonality (mention assumptions about period)

1.3.2 Time series representation

- Principal Component Analysis
- SAX
- Discrete Fourier Transform (DFT) and Discrete Wavelet Transform (DWT)
 - mention origin in signal processing and ubiquitous use in engineering (image and audio compression)

1.3.3 Challenges when building a time series

- length of series
- trend
- seasonality
- computational complexity -> issue because of data size
- granularity or sampling rates
- noise
- data quality
- similarity is task dependent (level)
- usual need for preprocessing the time series data (denoising, detrending, amplitude scaling)
 - > any pre-processing does modify the series

1.3.4 Data Analysis

what does M4 data look like

1.3.5 Challenges

- How many frequencies to compare?
- priorities of frequencies (power spectrum)
- different length of time series (leading to different frequencies) - ranges solved with logs

1.4 Methodology

1.4.1 Used Data

The research in time series has been numerous and focused on various properties of them as well as finding methods to accurately predict them. Aside of forecasting all researched areas are measures of similarity and retrieval of time series.

- **Forecasting** In the arena of forecasting the M-competition organized by Prof. Makridakis played a big role in the development of forecasting methods shortly after their inception in 1979.

One of the aspects that has been correct up until the 5th installment of the M-competition is that statistical methods in forecasting have outperformed more complex machine learning methods. So learning algorithms did not benefit sufficiently from learning from multiple series to generate more accurate point predictions and prediction intervals compared to the statistics-based alternatives.

One interesting question in this area is whether clustering of time series that have similar properties and training algorithms per cluster of "similar" series can help simplify the learning process for machine learning methods and in consequence improve their performance in future competitions.

However, expressing similarity for time series is a challenging questions with respect to which metrics to utilize, computational complexity as well as limiting assumptions that need to be made for time series.

1.4.2 Main contribution of the thesis

- transformation into Fourier-space
- transfer frequencies into frequency range band with increasing range width (using log scale)
- computation of frequency energy levels (sort and keep top 5) -> ask Prof. how to name this parameter
- conversion of ordered frequencies into frequency range band
- for each series to compare -> compare whether the frequency matches on the ordered positions -> provide exponential value per position -> match on more powerful frequencies is valued higher

1.4.3 additional computations

- utilization of FFT utilizes only frequency space (future work should consider comparison of energy levels per frequency)
- additional simple statistics computed (mean, std, quantiles)
- ts decomposition for trend estimation (requires parameter for period) -> then best line fit for slope of the time series
- computation of deltas for each series to search with statistics and slope of all other time series (review computational complexlity)
- ranking of matching series based highest frequency range match and ONE statistic

1.4.4 Preprocessing

- M4 data wide format vs. long format

1.4.5 Parallelization

- computation times
- scalability
- Samples for results only (stratification vs. non-stratification)
- Threads vs. Processes

1.4.6 Technology (check with Prof. if required)

R vs. Python vs. Mathematica, Matlab

1.4.7

- load
- transform to FFT vector space
- compare most important frequencies
- compare candidates
- select winner (which criteria)

1.5 Exploratory Data Study

- what do results look like

1.6 Formal Evaluation

- (maybe) improvement in forecasting approach
- find dataset with ground truth and compare DTW to this approach
- Distance metrics
- computational complexity

1.7 Conclusion & future work

1.7.1 Successes

1.7.2 Failures

1.7.3 Flaws

- final computation

1.7.4 What is missing

- denoising of time series
- adjustment of number of frequencies used
-

1.8 Results & Discussion

1.9 References

- [1] R.J. Hyndman and G. Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2014. ISBN: 9780987507105. URL: <https://books.google.de/books?id=gDuRBAAAQBAJ>.
- [2] Tak-chung Fu. “A review on time series data mining”. In: *Engineering Applications of Artificial Intelligence* 24.1 (Feb. 2011), pp. 164–181. ISSN: 0952-1976. DOI: 10.1016/j.engappai.2010.09.007. URL: <http://dx.doi.org/10.1016/j.engappai.2010.09.007>.
- [3] K.J. Åström. “On the choice of sampling rates in parametric identification of time series”. In: *Information Sciences* 1.3 (July 1969), pp. 273–278. ISSN: 0020-0255. DOI: 10.1016/S0020-0255(69)80013-7. URL: [http://dx.doi.org/10.1016/S0020-0255\(69\)80013-7](http://dx.doi.org/10.1016/S0020-0255(69)80013-7).
- [4] Yongqiang Tang, Yuan Xie, Xuebing Yang, et al. “Tensor Multi-Elastic Kernel Self-Paced Learning for Time Series Clustering”. In: *IEEE Transactions on Knowledge and Data Engineering* (2019), pp. 1–1. ISSN: 2326-3865. DOI: 10.1109/tkde.2019.2937027. URL: <http://dx.doi.org/10.1109/tkde.2019.2937027>.
- [5] Elon Musk. “An Integrated Brain-Machine Interface Platform With Thousands of Channels”. In: *Journal of Medical Internet Research* 21.10 (Oct. 2019), e16194. ISSN: 1438-8871. DOI: 10.2196/16194. URL: <http://dx.doi.org/10.2196/16194>.
- [6] Joshua H Siegle, Aarón Cuevas López, Yogi A Patel, et al. “Open Ephys: an open-source, plugin-based platform for multichannel electrophysiology”. In: *Journal of Neural Engineering* 14.4 (June 2017), p. 045003. ISSN: 1741-2552. DOI: 10.1088/1741-2552/aa5eea. URL: <http://dx.doi.org/10.1088/1741-2552/aa5eea>.
- [7] Xiaoyue Wang, Abdullah Mueen, Hui Ding, et al. “Experimental comparison of representation methods and distance measures for time series data”. In: *Data Mining and Knowledge Discovery* 26.2 (Feb. 2012), pp. 275–309. ISSN: 1573-756X. DOI: 10.1007/s10618-012-0250-5. URL: <http://dx.doi.org/10.1007/s10618-012-0250-5>.
- [8] Donald J. Berndt and James Clifford. “Using Dynamic Time Warping to Find Patterns in Time Series”. In: KDD Workshop. 1994, pp. 359–370.