

AI-Driven Exploration and Prediction of Company Registration Trends with Registrar of Companies (RoC)

TEAM LEADERS

510521104304: EZHUMALAI P

PHASE-2:DOCUMENT SUBMISSION



OBJECTIVE:

The problem is to perform an AI-driven exploration and predictive analysis on the master details of companies registered with the Registrar of Companies (RoC). The objective is to uncover hidden patterns, gain insights into the company landscape, and forecast future registration trends.

PHASE 2: Innovation

Consider exploring advanced AI algorithms like time series forecasting or ensemble methods for improved predictive accuracy.

Dataset Link: <https://tn.data.gov.in/resource/company-master-data-tamil-nadu-upto-28th-february-2019>

ABSTRACT

This paper describes a new type of ensembles that aims at improving the predictive performance of these approaches in time series forecasting. Ensembles are recognised as one of the most successful approaches to prediction tasks. Previous theoretical studies of ensembles have shown that one of the key reasons for this performance is diversity among ensemble members. Several methods exist to generate diversity. The key idea of the work we are presenting here is to propose a new form of diversity generation that explores some specific properties of time series prediction tasks. Our hypothesis is that the resulting ensemble members will be better at addressing different dynamic regimes of time series data. Our large set of experiments confirms that the methods we have explored for generating diversity are able to improve the performance of the equivalent ensembles with standard diversity generation procedures. Keywords: ensemble methods, time series forecasting, bagging, maximum embed

1. Introduction

Ensembles are known to be among the most competitive forms of solving predictive tasks. Several studies (e.g., Dietterich (2000); Brown and Kunchewa (2010)) have been carried out to understand and explain the reasons for their competitiveness in a wide range of application domains. Diversity among the individual components of ensembles is known to be a key element to generate a successful model. Bagging (Breiman, 1996) is a well known and simple type of ensembles that consists of a large set of standard tree-based models that are grown with the goal of generating a diverse set of models. Diversity in bagging is created through the use of different random bootstrap samples of the original training set to grow each tree. The main idea behind this paper is to propose a variant of bagging where the forms of generating diversity are biased towards specific characteristics of time series forecasting tasks. Namely, we aim at trying to have different forms of handling the diverse

2. Problem Description

The standard definition of time series forecasting assumes the existence of a set of time-ordered observations of a variable, y_1, y_2, \dots, y_t , where y_i is the value of Y measured at time i , and defines the predictive task as trying to forecast the future values of this variable for time stamps $s > t$. Many variants of this general task exist, including the use of other measured variables as potential predictors of the future values of the target series Y . Still, the general assumption is that there is an unknown function that "maps" the past observations into the future values of Y , i.e. $Y_{t+h} = f(\langle \text{DescriptorsOfThePast} \rangle)$, and the learning goal is to approximate this function using some prediction error criterion and a historical record of observed values. The predictors used for forecasting the future values of Y are usually the most recent observations of Y , as the basic assumption of time series forecasting is that of the existence of some form of correlation between successive observations of the series. This is the approach used on most approaches to time series forecasting, like for instance the well-known ARIMA models (e.g. Chatfield (2013)). This is also the idea of time delay coordinate embedding (Takens, 1981) that is a standard procedure for applying out of the box regression tools to time series forecasting tasks. This strategy assumes that future values of the series are only dependent on a limited number of previous values. In this context, delay coordinate embedding consists in using the k past values of a time series as descriptors of the state of the system at an instant t . If k is appropriate, it's possible to capture the dynamics of the time series from the embed vectors $r_t = \langle y_t, y_{t-1}, \dots, y_{t-k} \rangle$. Under this assumption, we can then use any regression tool to obtain a model of the form $Y_{t+h} = f(r_t)$ that specifies the relationship between a set of predictors (described by the embed vector) and the future values of the series.

3. Bagging for Time Series Forecasting

The approaches based on the use of the recent past values of a time series (the embed) as predictors require setting a critical parameter - how many past values to include, i.e. the size of the embed. Setting this parameter is not trivial and it may involve trying different alternative values for the embed size and use some reliable performance estimation process as a means for deciding the "optimal" value. The main drawback of these approaches is the fact that frequently there may not exist one single correct answer. In effect, non-stationary series and the occurrence of different regime shifts along time may lead to the best value being clearly time-dependent. This is one of the main motivations for our work. Another being previous work showing that having diversity in ensembles is a key ingredient to boost

Table 1: Summary of the main characteristics of the ensemble variants.

Embed size	Extra predictors
E	All models use k_{max} . None.
E+S	All models use k_{max} . All models use μY and $\sigma^2 Y$
	calculated with the respective embed.
	DE
	One third of the models use k_{max} , another third uses $k_{max}/2$, and the last third uses $k_{max}/4$.
	None.
	DE+S
	One third of the models use k_{max} , another third uses $k_{max}/2$, and the last third uses $k_{max}/4$.
	All models use μY and σ^2
	Y
	calculated with the respective embed.
	DE±S

One third of the models use k_{\max} ,
 another third uses $k_{\max}/2$,
 and the last third uses $k_{\max}/4$.
 Half of the models using a certain
 embed size use μ_Y and σ^2_Y
 calculated with the respec

4. Experimental Evaluation

The main goal of our experimental evaluation is to check the validity of the hypothesis that the variants of bagging we have described in Section 3 will outperform standard bagging on time series forecasting tasks. In this context, our baseline benchmark is a standard bagging implementation using the approach tagged as E in the list given in Section 3. All other four variants will use the same base data (the values of the past k_{\max} observations) as trainingset, but will use it in a different way, e.g. by using only part of it in some models or by using it to generate extra features. We also compare the ARIMA model, a more standard time series forecasting approach, to the same baseline approach to shed more light onto the overall competitiveness of our proposal. Since ARIMA models usually require a significant parameter tuning effort to

OLIVEIRA TORGO

obtain good results, we used the `auto.arima` function available in the R package `forecast`

(Hyndman et al., 2014) which automatically searches for an optimal model.

Table 2: Data sets used

ID Time series Data source Data characteristics

1 Temperature

Bike Sharing

(Fanaee-T and

Gama, 2013)

2 Humidity Daily values from Jan. 1, 2011

3 Windspeed to Dec. 31, 2012 (731 values)

4 Count of total bike rentals

5 Temperature

6 Humidity Hourly values from Jan. 1, 2011

7 Windspeed to Dec. 31, 2012 (7379 values)

8 Count of total bike rentals

9 Flow of Vatnsdalsa river Icelandic river

(Tong et al., 1985)

Daily values from Jan. 1, 1972

to Dec. 31, 1974 (1095 values)

10 Minimum temperature

Porto weather1 Daily values from Jan. 1, 2010

to Dec. 28, 2013 (1457 values)

11 Maximum temperature

12 Maximum steady wind

13 Maximum wind gust

14 Total precipitation

All five alternative forms of bagging and the ARIMA model were tested on fourteen real world time series obtained from three different data sources as described in **Table 2**. Each one of these series of data was treated separately from the others in their respective data source (e.g. information on the weather was not used to predict the total number of bike rentals). Moreover, please note that each time series of the Bike Sharing data source is available in two formats (daily and hourly), which were treated as different time series forecasting tasks. All fourteen time series were pre-processed to overcome some well-known issues with this type of data. Specifically, we have created all data sets used in our experiments with the series of the differences between successive values, and not from the original absolute values, in order to avoid trend effects. We have not, however, assumed any shape of these effects, if they exist. The target variable for all tasks was set to the next value of the series of differences.

Table 3:

Paired comparisons results in format Nr.Wins (Statistically Significant Wins)/ Nr.Losses (Statistically Significant Losses) M kmax Variant Wins/Losses

1020

20

E+S 13 (11) / 1 (1)

DE 7 (7) / 7 (3)

DE+S 13 (10) / 1 (0)

DE±S 14 (12) / 0 (0)

ARIMA 7 (3) / 7 (4)

30

E+S 11 (9) / 3 (2)

DE 10 (6) / 4 (3)

DE+S 10 (5) / 4 (2)

DE±S 10 (9) / 4 (2)

ARIMA 6 (3) / 8 (4)

1500

20

E+S 13 (10) / 1 (1)

DE 8 (6) / 6 (3)

DE+S 13 (10) / 1 (0)

DE±S 14 (12) / 0 (0)

ARIMA 7 (3) / 7 (4)

30

E+S 11 (9) / 3 (2)

DE 9 (7) / 5 (3)

DE+S 10 (7) / 4 (2)

Finally, Figure 1 presents more detailed values of the percent difference of MSE with relation to the baseline for each time series (summarized on **Table 5**). For visualization

purposes, a signed common logarithm was applied to the results. The represented metric

is, therefore,

$$\text{sgn}(M_{\text{SEX}} - M_{\text{SEE}}) \cdot \log$$

$$(|$$

$$|$$

$$|$$

$$|$$

$$100 \cdot (M_{\text{SEX}} - M_{\text{SEE}})$$

$$M_{\text{SEE}}$$

$$|$$

$$|$$

$$|$$

$$| + 1$$

$$)$$

Once more, our proposals, in general, seem to do well in comparison to standard bagging. The more apparent exceptions to this are the results obtained for $k_{\text{max}} = 30$ on time series 5, 6 and 8, with which the ARIMA model seems to struggle as well. The high variance of performance of the ARIMA approach is well illustrated in this figure. Although it achieves a higher decrease in MSE on datasets 9-14, our proposed approaches are also able to perform well on those while almost always beating the results of the ARIMA model on datasets 1-8 with a significant margin.

5. Related Work

Using different predictors on each member of an ensemble is not a novel idea. Random forests (Breiman, 2001) for instance, grow each tree in such a way that for each node a random subset of the features is used to select the best split. Random subspaces (Ho, 1998) is another example of an ensemble method that uses diverse sets of features among the models. This approach consists in randomly selecting subsets (usually of the same size) of the feature space to build each base learner. Contrary to our approach, none of these previous works address time series tasks. Moreover, our subsets of predictors are not

6. Conclusions

This paper describes an initial attempt at proposing ensembles for time series forecasting tasks. The main motivation of this ongoing work is the observation that handling time series tasks requires several decisions in terms of how we describe the recent dynamics of the observed values of the series. Settling on a single answer to these decisions may be dangerous in real world time series where one frequently observes changes in the dynamic properties of the variable being measured. Ensembles are a well-known answer to this type of problems by taking advantage of diversity among models to reduce both the bias and variance components of the prediction error. Motivated by these observations we have proposed an initial set of forms of injecting diversity into ensembles that takes into account some specific challenges posed by time series data. Namely, we have considered alternative ways of representing the recent observations of the target series among the members of the ensemble. These include the use of different sizes of the embed and also the addition of variables summarising the recent observed values.