

Peeking inside FFORMS: Feature-based FOfRecast Model Selection

Abstract

Features of time series are considered to be an important factor in identifying suitable forecast models. We proposed a classification framework, called FFORMS (Feature-based FOfRecast Model Selection), which selects forecast models based on features calculated from the time series. FFORMS framework builds a mapping that relates the features of time series to the “best” forecast-model using random forest algorithm. In this paper we make an effort to explore what is happening under the hood of FFORMS framework and thereby gain an understanding of what features led to the different choices of forecast-models and how different features influence the predicted outcome. This is accomplished by using model-agnostic machine learning interpretability approaches. Partial-dependency plots are used to visualize both main and interaction effects of features. The results of this study provide a valuable insight into how different features and their interactions affect the choice of forecast-model selection and give a more refined picture of the relationship between features and the choice of forecast-model which is particularly valuable for ongoing research in the field of feature-based time series analysis.

Keywords: FFORMS, time series, machine learning interpretability, black-box models, LIME

1 Introduction

The time series forecasting field has been evolving for a long time and has introduced a wide variety of forecasting methods. However, for a given time series the selection of an appropriate forecasting-method among many possibilities is not straight forward. This selection is one of the most difficult tasks as each method perform best for some but not all tasks. The features of time series are considered to be an important factor in identifying suitable forecasting models (Meade 2000; Makridakis & Hibon 2000). However, a description of relationship between the features and the performance of algorithms is rarely discussed in the field of forecasting.

There have been several recent studies on use of meta-learning approaches to automate the

forecast model selection based on the features computed from the time series (Shah 1997; Prudêncio & Ludermir 2004; Lemke & Gabrys 2010; Kück, Crone & Freitag 2016). Meta-learning approach provides a systematic guidance on model selection based on knowledge acquire from historical data set. The key idea is, forecast-model selection is posed as a supervised learning task. Each time series in the meta-data set is represented as a vector of features and labeled according to the “best” forecasting method(i.e. lowest MASE, etc.). Then a meta-learner is trained to identify suitable forecasting models. [Usually a machine learning algorithm is used.] With the era of big data, such an automated model selection process is necessary because the cost of invoking all possible forecasting method is prohibitive. However, these work suffer from the limitation of providing meaningful interpretations that can enhance understanding of relations between features and model outcomes which leads to the questions of

1. How features are related to the property being modeled?
2. How features interact with each other to identify the suitable forecasting method?
3. Which features contribute the most to classify a specific instance?

This results in less transparency and trustability of the model. To best of our knowledge, very limited efforts have been taken to understand how the models are making its decisions and what is really happening inside these complex model structures.

On the other hand, aside from the goal of developing automated forecast-model selection framework few researchers have made an attempt to provide a description of relationship between the features and the performance of algorithms (Schnaars 1984; Wang, Smith-Miles & Hyndman 2009; Lemke & Gabrys 2010; Petropoulos et al. 2014). However, these studies are limited by the scale of problem instances used, diversity of forecasting-methods used, quality of features considered, and modelling approached used to identify the relationship between features and forecast model performance. Most of these studies are typically restricted to simple statistical techniques such as simple linear models, decision trees, etc.

To fill this gap, this paper makes a first step towards providing a comprehensive explanation of the relationship between time series features and forecast-model selection using machine learning interpretability techniques. This paper builds on the method from our previous work ref, in which we introduced the FFORMS (Feature-based FOREcast Model Selection) framework. The random forest algorithm is used to model the relationship between features and “best” performing forecast-model. A large collection of time series is used to train the model. We use 35 features to capture different characteristics of time series. The reason for the choice of random

forest algorithm is its ability to model complex variable interactions. One noticeable significance of our approach is this can be parallized to for any given computing budget and time. Although the prediction accuracy of random forest algorithm is high, it is not easy to interpret what is happening inside the forest because of the two-step randomization. In this work we aim at providing a deeper understanding of the underlying mechanism and influence of features in forecast-model selection.

[What is the usefulness of analyzing the relationship between features and model performance?] Understanding the role of features is worthwhile even if producing an accurate and generalizable model is the only objective of the modelling. This is because the less transparency of the model may be distrusted regardless of their predictive performance. The methodology we propose here is a novel application of machine learning interpretability methods to visualize and explore the role of features in forecast-model selection.

This paper proceeds as follows. In [Section 2](#) we describe the application of FFORMS framework to M4competition data. The main contribution of this from our previous work Talagala, Hyndman & Athanasopoulos (2018) is we extend the FFORMS framework to model weekly, daily and hourly series. [Section 3](#) gives background on machine learning interpretability techniques that are used to identify role of features in forecast model selection. In [Section 4](#) we discuss the results. [Section 5](#) concludes.

2 FFORMS Application to M4 competition data

The FFORMS framework consists of two main components: i) *offline phase*, which includes the development of a classification model and ii) *online phase*, use the classification model developed in the offline phase to identify “best” forecast-model. We develop separate classifiers for yearly, monthly, quarterly, weekly, daily and hourly series.

2.1 FFORMS framework: offline phase

2.1.1 observed sample

We split the time series in the M4 competition into training set and test set. The time series in the training set are used as the set of observed time series. The time series in the test set are used to evaluate the classification models. Further, for yearly, quarterly and monthly time series in addition to the time series provided in the M4 competition we used the time series of M1 and

M3 competitions. Table 1 summarizes the number of time series in the observed sample and the test set in each frequency category.

Table 1: *Composition of the time series in the observed sample and the test set*

Frequency	Observed Sample			Test set
	M1	M3	M4	M4
Yearly	181	645	22000	1000
Quarterly	203	756	23000	1000
Monthly	617	1428	47000	1000
Weekly	-	-	259	100
Daily	-	-	4001	226
Hourly	-	-	350	64

2.1.2 simulated time series

As described in Talagala, Hyndman & Athanasopoulos (2018), we augment the reference set by adding multiple time series simulated based on each series in the M4 competition. We use several standard automatic forecasting algorithms to simulate multiple time series from each series. Table 2 shows the different automatic forecasting algorithms used under each frequency category. The automated ETS and ARIMA are implemented using `ets` and `auto.arima` functions available in the forecast package in R (Hyndman et al. 2018). The `stlf` function in the forecast package (Hyndman et al. 2018) is used to simulate multiple time series based on multiple seasonal decomposition approach. As shown in Table 2 we fit models to each time series in the M4 competition database from the corresponding algorithm and then simulate multiple time series from the selected models. Before simulating time series from daily and hourly series we convert the time series into multiple seasonal time series (msts) objects. For daily time series with length less 366 the frequency is set to 7 and if the time series is long enough to take more than a year ($\text{length} > 366$), the series is converted to a multiple seasonal time series objects with frequencies 7 and 365.25. For hourly series, if the series length is shorter than 168, frequency is set to 24, if the length of the series is greater than 168 and less than or equals to 8766 only daily and weekly seasonality are allowed setting the frequencies to 24 and 168. In this experiment the length of the simulated time series is set to be equal to: length of the training period specified in the M4 competition + length of the forecast horizon specified in the competition. For example, the series with id “Y13190” contains a training period of length 835. The length of the simulated series generated based on this series is equals to 841 (835+6).

As illustrated in Talagala, Hyndman & Athanasopoulos (2018), the observed time series and the simulated time series form the reference to build our classification algorithm. Once we create the reference set for random forest training we split each time series in the reference set into

Table 2: Automatic forecasting algorithms used to simulate time series

Algorithm	Y	Q	M	W	D	H
automated ETS	✓	✓	✓			
automated ARIMA	✓	✓	✓			
forecast based on multiple seasonal decomposition				✓	✓	✓

training period and test period.

2.1.3 Input: features

The FFORMS framework operates on the features of the time series. For each time series in the reference set features are calculated based on the training period of the time series.

Table 3: Time series features

	Feature	Description	Y	Q/M	W	D/H
1	T	length of time series	✓	✓	✓	✓
2	trend	strength of trend	✓	✓	✓	✓
3	seasonality 1	strength of seasonality corresponds to frequency 1	-	✓	✓	✓
4	seasonality 2	strength of seasonality corresponds to frequency 2	-	-	-	✓
5	linearity	linearity	✓	✓	✓	✓
6	curvature	curvature	✓	✓	✓	✓
7	spikiness	spikiness	✓	✓	✓	✓
8	e_acf1	first ACF value of remainder series	✓	✓	✓	✓
9	stability	stability	✓	✓	✓	✓
10	lumpiness	lumpiness	✓	✓	✓	✓
11	entropy	spectral entropy	✓	✓	✓	✓
12	hurst	Hurst exponent	✓	✓	✓	✓
13	nonlinearity	nonlinearity	✓	✓	✓	✓
14	alpha	ETS(A,A,N) $\hat{\alpha}$	✓	✓	✓	-
15	beta	ETS(A,A,N) $\hat{\beta}$	✓	✓	✓	-
16	hwalpha	ETS(A,A,A) $\hat{\alpha}$	-	✓	-	-
17	hwbeta	ETS(A,A,A) $\hat{\beta}$	-	✓	-	-
18	hwgamma	ETS(A,A,A) $\hat{\gamma}$	-	✓	-	-
19	ur_pp	test statistic based on Phillips-Perron test	✓	-	-	-
20	ur_kpss	test statistic based on KPSS test	✓	-	-	-
21	y_acf1	first ACF value of the original series	✓	✓	✓	✓
22	diff1y_acf1	first ACF value of the differenced series	✓	✓	✓	✓
23	diff2y_acf1	first ACF value of the twice-differenced series	✓	✓	✓	✓
24	y_acf5	sum of squares of first 5 ACF values of original series	✓	✓	✓	✓
25	diff1y_acf5	sum of squares of first 5 ACF values of differenced series	✓	✓	✓	✓
26	diff2y_acf5	sum of squares of first 5 ACF values of twice-differenced series	✓	✓	✓	✓
27	seas_acf1	autocorrelation coefficient at first seasonal lag	-	✓	✓	✓
28	sediff_acf1	first ACF value of seasonally-differenced series	-	✓	✓	✓
29	sediff_seacf1	ACF value at the first seasonal lag of seasonally-differenced series	-	✓	✓	✓
30	sediff_acf5	sum of squares of first 5 autocorrelation coefficients of seasonally-differenced series	-	✓	✓	✓
31	seas_pacf	partial autocorrelation coefficient at first seasonal lag	-	✓	✓	✓
32	lmres_acf1	first ACF value of residual series of linear trend model	✓	-	-	-
33	y_pacf5	sum of squares of first 5 PACF values of original series	✓	✓	✓	✓
34	diff1y_pacf5	sum of squares of first 5 PACF values of differenced series	✓	✓	✓	✓
35	diff2y_pacf5	sum of squares of first 5 PACF values of twice-differenced series	✓	✓	✓	✓

The description of the features calculated under each frequency category is shown in Table 3. A

comprehensive description of the features used in the experiment is given in Talagala, Hyndman & Athanasopoulos (2018).

2.1.4 Output: class-labels

In addition to the class labels used by Talagala, Hyndman & Athanasopoulos (2018) we include some more class labels when applying the FFORMS framework to the M4 competition time series. The description of class labels considered under each frequency is shown in Table 4. We fit the corresponding models outlined in Table 4 to each series in the reference set. The models are estimated using the training period for each series, and forecasts are produced for the test periods.

Table 4: *Class labels*

class label	Description	Y	Q/M	W	D/H
WN	white noise process	✓	✓	✓	✓
AR/MA/ARMA	AR, MA, ARMA processes	✓	✓	✓	-
ARIMA	ARIMA process	✓	✓	✓	-
SARIMA	seasonal ARIMA	✓	✓	✓	-
RWD	random walk with drift	✓	✓	✓	✓
RW	random walk	✓	✓	✓	✓
Theta	standard theta method	✓	✓	✓	✓
STL-AR		-	✓	✓	✓
ETS-notrendnoseasonal	ETS without trend and seasonal components	✓	✓	✓	-
ETStrendonly	ETS with trend component and without seasonal component	✓	✓	✓	-
ETSDampedtrend	ETS with damped trend component and without seasonal component	✓	✓	-	-
ETStrendseasonal	ETS with trend and seasonal components	-	✓	-	-
ETSDampedtrendseasonal	ETS with damped trend and seasonal components	-	✓	-	-
ETSseasonalonly	ETS with seasonal components and without trend component	-	✓	-	-
snaive	seasonal naive method	✓	✓	✓	✓
tbats	TBATS forecasting	-	✓	✓	✓
nn	neural network time series forecasts	✓	✓	✓	✓
mstlets		-	-	✓	✓
mstlarima		-	-	-	✓

The `auto.arima` and `ets` functions in the `forecast` package are used to identify the suitable (S)ARIMA and ETS models. In order to identify the “best” forecast-model for each time series in the reference set we combine the mean Absolute Scaled Error (MASE) and the symmetric Mean Absolute Percentage Error (MAPE) calculated over the test set. More specifically, for each series both forecast error measures MASE and sMAPE are calculated for each of the forecast models. Each of these is respectively standardized by the median MASE and median sMAPE calculated across the methods. The model with the lowest average value of the scaled MASE and scaled sMAPE is selected as the output class-label. Most of the labels given in Table 4 are self-explanatory labels. In STL-AR, mstlets, and mstlarima, first STL decomposition method applied to the time series and then seasonal naive method is used to forecast the seasonal component. Finally, AR, ETS and ARIMA models are used to forecast seasonally adjusted data

respectively.

2.1.5 Train a random forest classifier

A random forest with class priors is used to develop the classifier. We build separate random forest classifiers for yearly, quarterly, monthly, weekly, daily and hourly time series. The wrapper function called `build_rf` in the `seer` package enables the training of a random forest and returns class labels("best" forecast-model) for each time series.

2.2 FFORMS framework: online phase

The online phase of the algorithm involves generating point forecasts and 95% prediction intervals for the M4 competition data. First, the corresponding features are calculated based on the full length of the training period provided by the M4 competition. Second, point forecasts and 95% prediction intervals are calculated based on the predicted class labels, in this case forecast-models. Finally, all negative values are set to zero.

3 Peeking inside FFORMS

The main theme of this paper is to explore the nature of the relationship between features and forecast-model selection. We use both model-diagnostic approaches and machine learning interpretability approaches.

3.1 Model-diagnostics

Model-diagnostic is an important aspect in evaluating the accuracy of the model's predictions as well as the model's understanding of the nature of the relationship between features and predicted outcome.

3.1.1 Out-of-bag(OOB) error and uncertainty measure for each observation

It is argued in order to estimate the test error of a bagged model it is not necessary to perform cross-validation approach, because each tree is grown using different bootstrap samples from the training set and a part of training data is not used in the tree construction (Breiman (2001); Chen, Liaw & Breiman (2004)). In general, each bagged tree does not make use of around one third of observations to construct the decision tree. These observations are referred to as the out-of-bag(OOB) observations. Each tree is grown based on different bootstrap samples hence,

each tree has different set of OOB observations. These OOB samples can be used to calculate internal estimation of the test set error.

3.1.2 Representation of model in the data space and data in the model space (d-in-ms)

Wickham, Cook & Hofmann (2015) explains the importance of displaying the “model in the data space (m-in-ds)” and “data in the model space (d-in-ms)”. Displaying the data in the model space (d-in-ms) is the most commonly used approach for model-diagnostics. For example, plot of fitted values versus residuals (Wickham, Cook & Hofmann (2015)). D-in-MS is a visualisation of embedding high-dimensional data into a low-dimensional space generated from the model. Visualisation of D-in-MS do not help to gain an understanding of the nature of the relationship between features predicted outcome. In order to address this issue Wickham, Cook & Hofmann (2015) and Silva, Cook & Lee (2017) have highlighted the importance of visualisations of model in the data space. In the context of classification, representation of m-in-ds could be achieved by first, projecting the training data set into meaningful low-dimensional feature space and then visualize the complete prediction regions or their boundaries. In other words this can be considered as the visualization of predictor space in the context of the data space. See Wickham, Cook & Hofmann (2015) for visualisation method of this kind and Silva, Cook & Lee (2017) for comparable method for random forest algorithms.

3.2 Machine Learning Interpretability

In recent years, there have been a growing interest for interpretability of machine learning algorithms with European General Data Protection Regulation (GDPR) stipulates the explainability of all automatically made decision concerning individuals. We explore the role of features in two different perspectives: i) global explanation of feature contribution: overall role of features in the choice of different forecast model selection, and ii) local explanation of feature contribution: nature of the contributions features make for a prediction of a specific instance. We will introduce each of these ideas briefly below. Model-diagnostics tools are used.

3.3 General Notation

Let $\mathcal{P} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ be the historical data set we use to train the classifier. Consider a p -dimensional feature vector $X = (X_1, X_2, \dots, X_p)$ and a dependent variable, best forecasting method for each series Y . Let \mathcal{G} be the unknown relationship between X and Y . Zhao & Hastie (n.d.) term this as “law of nature”. Inside the FFORMS framework, random forest algorithm

tries to learn this relationship using the historical data we provided. We denote the predicted function as g .

3.4 Global Interpretability Methods

Global interpretability evaluate the behavior of a model on entire data set. Global perspective of model interpretation helps users to understand the overall modeled relationship between features and the model outcome. For example, which features are contribute mostly to the predictive mechanism of the fitted model, complex interactions between features, etc. In the following subsections we provide a description of tools we use to explore the global perspective of the model.

3.5 Analysis of Feature contribution

Jiang & Owen (2002) explains variable importance under three different views: i) causality: change in the value of Y for a increase or decrease in the value of x , ii) contribution of X based on out-of-sample prediction accuracy and iii) face value of X on prediction function g , for example in linear regression model estimated coefficients of each predictor can be considered as a measure of variable importance. See Jiang & Owen (2002) for comparable face value interpretation for machine learning models. In this paper we use the first two notions of variable importance. Partial dependency functions and individual conditional expectation curves are used to explore the “causality” notion of variable importance while Mean decrease in Gini coefficient and Permutation-based variable importance are used to capture the second notion of variable importance-features contribution to the predictive accuracy(Zhao & Hastie (n.d.)). We will introduce each of these variable importance measures below.

3.5.1 Mean decrease in Gini coefficient

Mean decrease in Gini coefficient is a measure of how each feature contributes to the homogeneity of the nodes and leaves in the resulting random forest proposed by Breiman (2001).

3.5.2 Permutation-based variable importance measure

The permutation-based variable importance introduced by Breiman (2001) measures the the prediction strength of each feature. This measure is calculated based on the out-of-bag (OOB) observations. The calculation of variable importance is formalized as follow: Let $\tilde{\mathcal{B}}^{(k)}$ be the OOB sample for a tree k , with $k \in \{1, \dots, ntree\}$, where $ntree$ is the number of trees in the random

forest. Then the variable importance of variable X_j in k^{th} tree is:

$$VI^{(k)}(X_j) = \frac{\sum_{i \in \mathcal{B}^{(k)}} I(\gamma_i = \gamma_{i,\pi_j}^k)}{|\mathcal{B}^{(k)}|} - \frac{\sum_{i \in \mathcal{B}^{(k)}} I(\gamma_i = \gamma_i^k)}{|\mathcal{B}^{(k)}|},$$

where γ_i^k denotes the predicted class for the i^{th} observation before permuting the values of X_j and γ_{i,π_j}^k is the predicted class for the i^{th} observation after permuting the values of X_j . The overall variable importance score is calculated as:

$$VI(X_j) = \frac{\sum_1^{ntree} VI^{(t)}(x_j)}{ntree}.$$

Permutation-based variable importance measures provide a useful starting point for identifying relative influence of features on the predicted outcome. However, they provide a little indication of the nature of the relationship between the features and model outcome. To gain further insights into the role of features inside the FFORMS framework we use partial dependence plot (PDP) introduced by Friedman, Popescu, et al. (2008).

3.5.3 Partial dependence plot (PDP)

Partial dependence plot can be used to graphically examine how each feature is related to the model prediction while accounting for the average effect of other features in the model. Let X_s be the subset of feature we want examine partial dependencies for and X_c be the remaining set of features in X . Then g_s , the partial dependence function on X_s is defines as

$$g_s(X_s) = E_{x_c}[g(x_s, X_c)] = \int g(x_s, x_c) dP(x_c).$$

In practice, PDP can be estimated from a training data set as

$$\bar{g}_s(x_s) = \frac{1}{n} \sum_{i=1}^n g(x_s, X_{ic}),$$

where n is the number of observations in the training data set. Partial dependency curve can be created by plotting the pairs of $\{(x_s^k, \bar{g}_s(x_{sk}))\}_{k=1}^m$ defined on grid of points $\{x_{s1}, x_{s2}, \dots, x_{sm}\}$ based on X_s . FFORMS framework have treated the forecast-model selection problem as a classification problem. Hence, in this paper partial dependency functions displays the probability of certain class occurring given different values of feature X_s .

3.5.4 Variable importance measure based on PDP

Greenwell, Boehmke & McCarthy (2018) introduced a variable importance measure based on the partial dependency curves. The idea is to measure the “flatness” of partial dependence curves for each feature. A feature whose PDP curve is flat, relative to the other features, indicates that the feature does not have much influence on the predicted value as it changes while taking into account the average effect of the other features in the model. The flatness of the curve is measured using the standard deviation of the values $\{\bar{g}_s(x_{sk})\}_{k=1}^m$.

3.5.5 Individual Conditional Expectation (ICE) curves

While partial dependency curves are useful in understanding the estimated relationship between feature and predicted outcome in the presence of substantial interaction between features, it can be misleading. Goldstein et al. (2015) proposed the Individual Conditional Expectation (ICE) curves to address this issue. Instead of averaging $g(x_s, X_{iC})$ over all observations in the training data, ICE plots the individual response curves by plotting the pairs $\{(x_s^k, g(x_{sk}, X_{iC}))\}_{k=1}^m$ defined on grid of points $\{x_{s1}, x_{s2}, \dots, x_{sm}\}$ based on X_s . In other words partial dependency curve is simply the average of all the ICE curves.

3.5.6 Variable importance measure based on ICE curves

This method is similar to the PDP-based VI scores above, but are based on measuring the “flatness” of the individual conditional expectation curves. We calculated standard deviations of each ICE curve. We then computed a ICE based variable importance score – simply the average of all the standard deviations. A higher value indicates a higher degree of interactivity with other features.

3.6 Assessment of Interaction Effect

Friedman’s H-statistic (Friedman, Popescu, et al. (2008)) is use to test the presence of interaction between all possible pair of features. This statistic is computed based on the partial dependence functions. For two way interaction between two specific variable x_j and x_k , Friedman’s H-statistic is defined as follow,

$$H_{jk}^2 = \sum_{i=1}^n [\bar{g}_s(x_{ij}, x_{jk}) - \bar{g}_s(x_{ij}) - \bar{g}_s(x_{ik})]^2 / \sum_{i=1}^n \bar{g}_s^2(x_{ij}, x_{jk}).$$

The Friedman’s H-statistic measures the fraction of variance of two-variable partial dependency,

$\bar{g}_s(x_{ij}, x_{jk})$ not captured by sum of the respective individual partial dependencies, $\bar{g}_s(x_{ij}) + \bar{g}_s(x_{jk})$. In addition to Friedman's H-statistic we also use the PDP of two variables to visualize the interaction effect.

Note that the, PD plots, ICE curves and PD-, ICE-associated measures and Friedman's H-statistic are computationally intensive to compute, especially when there are large number of observations in the training set. Hence, in our experiments ICE and PDP-based variable importance are measured based on the subset of randomly selected training examples.

3.7 Local Interpretable Model-agnostic Explanations (LIME)

Global interpretations help us to understand entire modeled relationship. Local interpretations help us to understand the predictions of the model for a single instance or a group of similar instances. In other words this allows users to zoom into a particular instance or a subset and explore how different features affect the resulting prediction. We use Local Interpretable Model-agnostic Explanations (LIME) approach introduced by Ribeiro, Singh & Guestrin (2016) for explaining individual predictions which relies on the assumption that "every complex model is linear on a local scale". This is accomplished by locally approximating the complex black-box model with a simple interpretable model. Ribeiro, Singh & Guestrin (2016) highlighted features that are globally important may not be in the local context and vice versa. The algorithm steps can be summarized as follow:

1. Select an observation of interest which we need to have explanations for its black-box prediction.
2. Create a permuted data set based on the selected observation. Permuted data set is created by making slight modifications to the features of selected observations.
3. Obtain similarity scores by calculating distance between permuted data and selected observation.
4. Obtain predicted outcomes for all permuted data using the black-box model.
5. Select m number of features best describing the black-box model outcome. This can be accomplished by applying feature selection algorithms such as ridge regression, lasso, ridge regression, etc.
6. Fit a simple linear model to the permuted data based on m selected features and similarity scores in step 3 as weights and complex model prediction outcomes in step 4 as response variable.

7. Use the estimated coefficients of simple linear model to explain the local behaviour corresponds to the selected observation in step 1.

An alternative for explaining local behaviour of complex models is proposed by Lundberg & Lee ([2017](#)) based on game theory named “Shapley values”.

4 Results

4.1 Yearly data

Model diagnostic: FFORMS framework, Yearly series

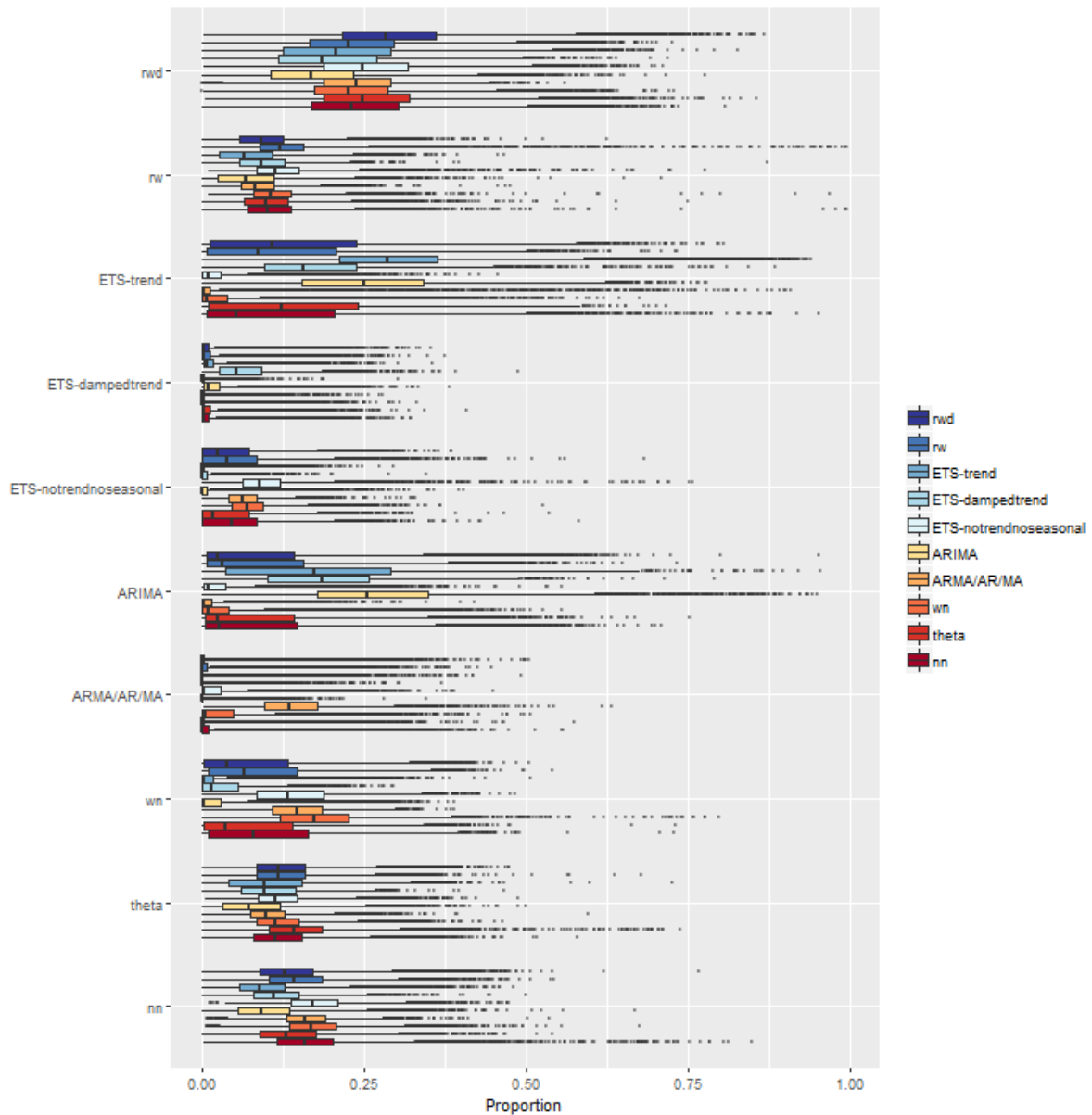


Figure 1: Distribution of proportion of times each yearly time series was assigned to each class based on OOB sample. Each row represent the predicted class label and colours of boxplots corresponds to class label of “best” forecasting method. There are ten rows in the plot corresponds to each predicted class represented by Y-axis. X-axis denotes the proportion of times the time series was predicted into each class. On each row, distribution of correctly classified class dominated the top, indicating a fairly good classification of the model fitted.

Figure 1 shows the distribution of proportion of times each observation (in our case each time series) was assigned to each class based on OOB sample. Each row represent the **predicted**

class label and colours of boxplots corresponds to **true class label**. The proportion 1 indicates, that the time series was always predicted to the corresponding class and 0 being never. This is an alternative way of visualizing the vote-matrix information in the random forest model. The other way of representing vote matrix involves ternary plot (Sutherland et al. (2000)) and jittered side-by-side dotplot (Ehrlinger 2015; Silva, Cook & Lee 2017). To overcome the problem of overlapping data points due to the scale of the training data set, similarity of classes and relatively large number of class labels, boxplot diagrams are used. This Figure 1 helps to evaluate the model performance in the data space (model-in-the-data-space) (Silva, Cook & Lee (2017)). On each row of Figure 1 the distribution of proportions corresponds to the time series in which the predicted class label and true class label are the same dominates the top indicating a fairly good classification of the model. In addition to that, in each row the distributions corresponds to the classes similar to the properties of the predicted class label also dominate the others. For example, within ETS-trend predicted class, the distributions correspond to the true class labels, ETS-damped trend, ARIMA, were also assigned with high probability and less values were assigned to ARMA/AR/MA, White noise process and ETS (ANN)/ ETS(MNN). This confirms that our FFORMS framework successfully learnt the similarities and dissimilarities between the classes itself. On average, random walk with drift have a high chance of getting selected with yearly time series. The results of M3-competition also concluded random walk with drift perform well for yearly time series. One reason for this is as shown in Kang, Hyndman, Li, et al. (2018) on average yearly series of M1, M3 and M4 are generally trended.

Figure 2 shows the contribution of features to the predictive mechanism of the FFORMS framework with respect to variable importance scores. Permutation-based variable importance and Gini feature importance measure are used to evaluate the overall feature importance. Moreover, most important features for each class is identified based on three measures: i) permutation-based variable importance, ii) partial dependence functions based variable importance and iii) ICE-curves based variable importance measure. The one that shows the highest importance is ranked 25, the second best is ranked 24, and so on. Finally, for each category, an average rank for each feature is computed based on the mean value of all rankings across all the feature importance metrics considered. The features, strength of trend and test statistic of Phillips-Perron(PP) unit root test, linearity, first autocorrelation coefficient of the differenced series, beta and lmres_acf1 are appear to be most important features in each class. Spikiness appear to be an important feature in the overall classification, except ARIMA category even though it does not appear to be among top five features, spikiness has been assigned a relatively high importance among all categories. These results indicates on average the features related to trend,

nonstationarity, overall shape of the trend (linear(measured by linearity), damped(measured by beta), exponential(measured by curvature)) and randomness (from spikiness, and $lmres_acf1$) are important for the choice of yearly time series forecasting methods. Further, y_acf1 is appear to be important in random walk with drift class and ARMA/AR/MA class. The length of time series (N) is assigned a high importance in random walk with drift, ETS-dampedtrend and neural-network class compared to others. Further exploration of partial dependency plots revealed neural network approach is likely to be selected in forecasting time series with long history of observations and probability of selecting random walk decreases as the length of the time series increases. Further, first correlation coefficient of the twice-differenced series is appear to be most important in ARIMA class as this category contains the higher order differenced series. Hurst exponent and entropy appear to be equally important in stationary classes. Within ETS-damped trend category beta and curvature ranked as important features. On the other hand, sum of squares of first five autocorrelation coefficients of the twice-difference series and lumpiness show lowest contribution across many classes.

Figure 3 shows the partial dependency curves, and associated confidence intervals of the top-three features that get selected most in each class. The three features shows a non-linear relationship with predicted outcome within each class. Probability of selecting ETS-trend, ARIMA, ETS-without seasonal and trend component and neural network models increases steadily as ur_pp increases. As expected probability of selecting stationary models decreases as the test statistic of Phillip-Perrion test increases and this probability remains zero after the value of 0 of ur_pp . Random walk with drift, ETS-trend, ETS-damped trend, ARIMA shows a increasing relationship with trend, whereas random walk, ETS-without trend and seasonal components, and stationary models shows a monotonically decreasing relationship as trend increases. The theta class shows parabolic relationship with trend. It is interesting to observe that probability of selecting neural network models decreases with very high trend value. The reason could be very clear highly trended series are more likely to select ETS-trend, ETS-dampedtrend and ARIMA models which lesson the chance of neural network models for them. The wide confidence bands around the partial dependency functions of linearity indicated the higher variability of ICE curves. Narrow confidence band corresponds to the random walk with drift indicate all individual ICE curves rise sharply around the value linearity = 0 and remained stable beyond value linearity=5. Random walk class depicted the mirror image of random walk with drift class. For ARMA/AR/MA class all the ICE curves increase sharply around 0 and decline steadily after that. Similar relationship appeared in white noise class with wide confidence bands whereas ARIMA and neural network show an opposite relationship.

Figure 4 shows the heat maps of relative strength of all possible pairwise interactions for each class. The relative strength of two-way interaction between features were determined using formulae developed by Friedman, Popescu, et al. (2008), which is implemented in the `iml` (Molnar, Casalicchio & Bischl (2018),) package in R. The test statistic of Phillip-Perron test, strength of trend, and linearity show a weak interaction with other features in all the class. However, except ETS-trend class, trend and `ur_pp` shows an high level of interactivity. These two features are appear to be among top 5 in all the classes according to the variable importance measures. Further, narrow confidence bands corresponds to these features in the partial dependency plots also confirms the less interactivity. In almost all the cases partial correlation and auto-correlation based features are heavily interacting. However, the first correlation coefficient of the difference series do not interact with other features heavily in the case of ARIMA class. Further, almost all pair of features appear to be interacting with in neural network category. ?? and ?? communicate more information about the nature of two-way feature interaction effect.

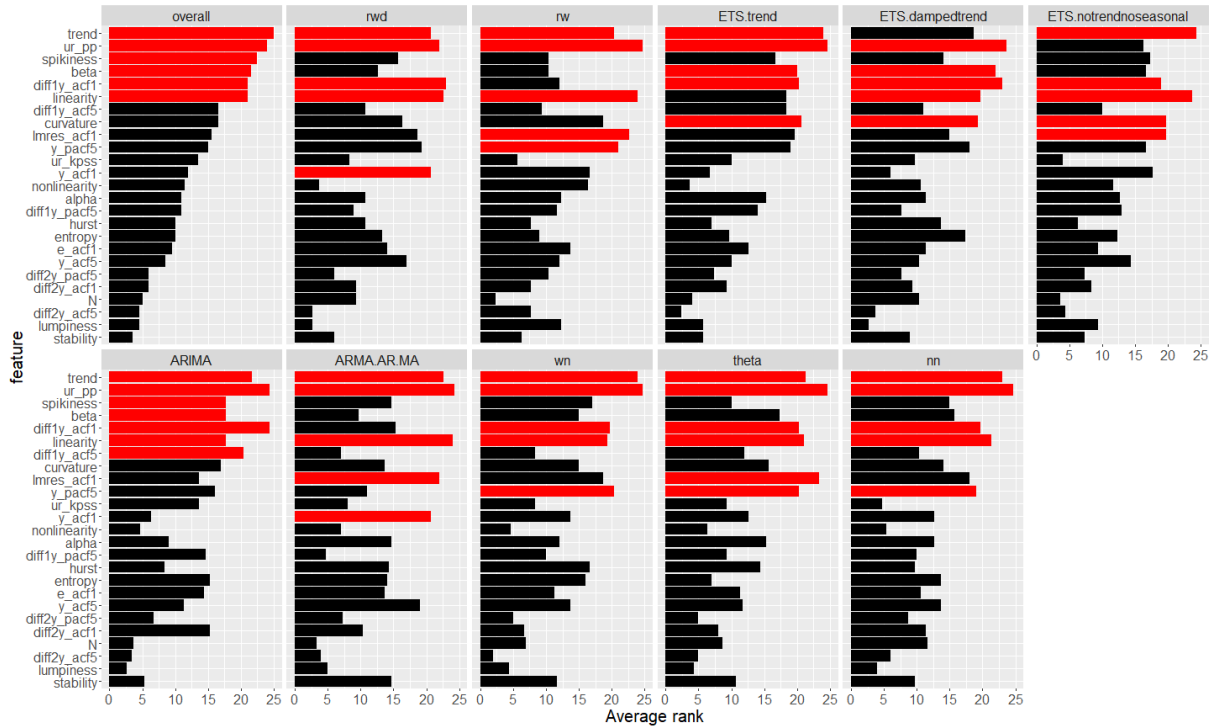


Figure 2: Feature importance plot for yearly series. Permutation-based VI measure and mean decrease in gini coefficients are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on three measures: i) permutation-based variable importance, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.

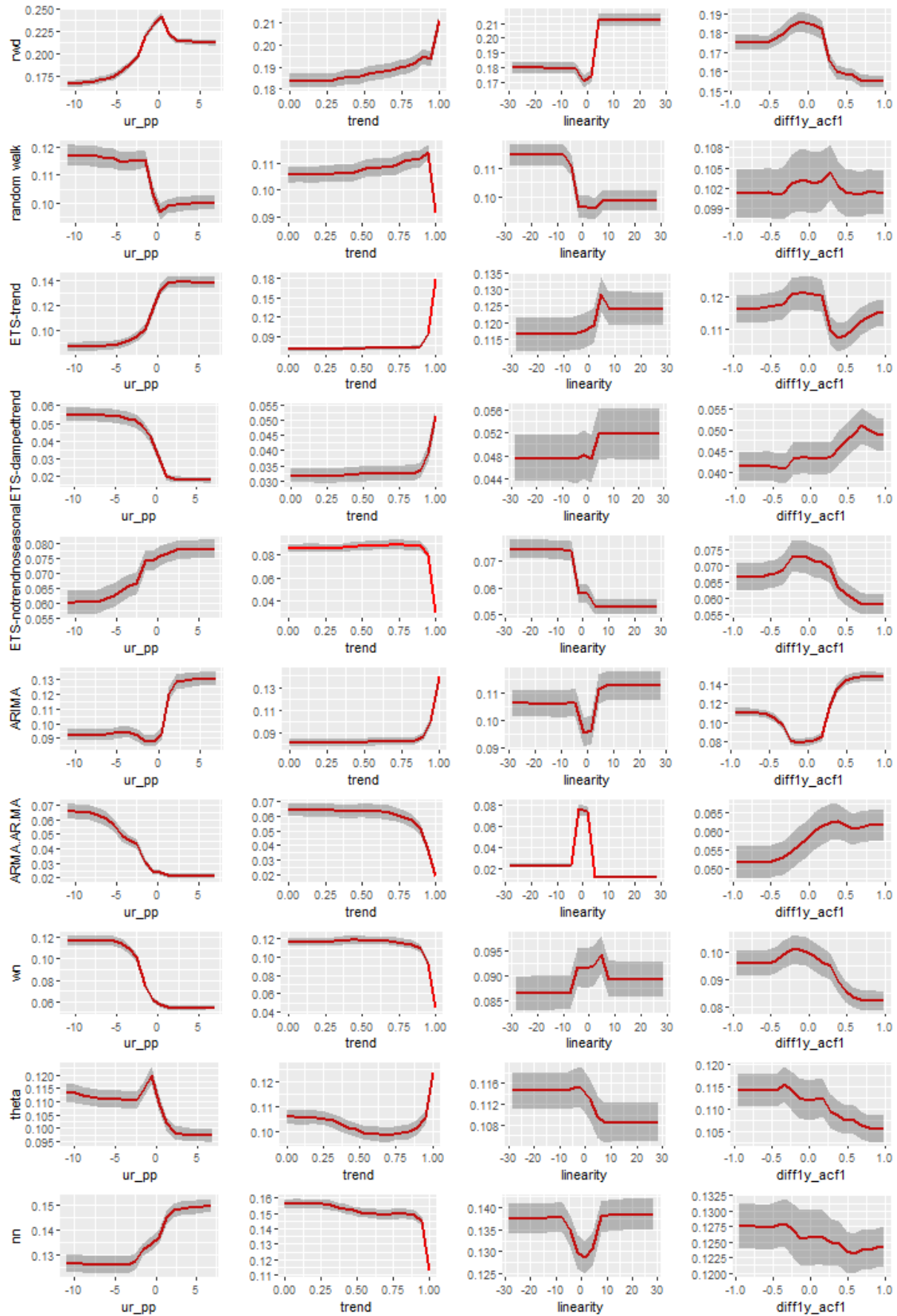


Figure 3: Partial dependence plots for the top-three features get selected most within each class. The shading shows the 95% confidence intervals. Y-axis denotes the probability of belong to corresponding class.



Figure 4: Heat maps of relative strength of all possible pairwise interactions calculated based on Friedman's H-statistic

4.2 Quarterly and Monthly data

Figure 5 shows oob error-based classification for quarterly and monthly data. According to Figure 5 the corresponding distributions depicted similar patterns both quarterly and monthly data. For quarterly and monthly data, the same set of features and class-labels are used to train the model. Hence, this consistency between the results of quarterly and monthly series would provide evidence in support of the validity and trustability of the model. On each row of ETS models a single distribution corresponds the correct classification class dominates the top. However, on average probability of selecting ETS-models for forecasting quarterly and monthly data is relatively low compared to the random walk with drift, stalar, tbats and neural network classes. Except the time series labeled as ARMA/AR/MA all other quarterly time series have a very low chance of being to ARMA/AR/MA class. Further, all distributions corresponds to the tbats row located further away from zero. This indicates all time series select tbats model at least once from the individual trees in the forest. Except few outliers, distributions within neural network category also show a slight upward deviation from zero. However, the upper boundary of these distributions do not surpass the upper boundaries of dominating box plots in the random walk with drift class and SARIMA class. Further, within stlar, tbats, theta and neural network classes all distributions level at similar proportionalities. These types of similarities in the distributions indicates the the appropriateness of using combination forecasting. Further, these information are useful in identifying potential time series models for combination forecast and improve the existing combination approaches proposed in the M4-competition. In addition to that the similarities and diversities observed in the boxplots indicates the neighborhood of cases in their respective instance space.

Figure 6 and Figure 7 shows variable importance plots for quarterly and monthly data respectively. For both quarterly and monthly data strength of seasonality, trend, linearity and spikiness are appear to be most important features across all categories. Even though the lumpiness does not appear as an top five feature within classes it is appear to be an important feature in the overall classification process and a relatively high rank is assigned within many classes. In the case of yearly data low variable importance is assigned to both stability and length of the series. However, in quarterly and monthly data high variable importance is assigned to length of the series and stability. One notable difference between quarterly series and monthly series is, for monthly data length of the series is ranked among top specially in random walk with drift, random walk, ETS with seasonal and trend component, ETS-seasonal, SARIMA and ARIMA classes. Hence, PDP of N for each class of monthly series are also shown in Figure 6 and Figure 7.

In addition to the strength of seasonality, the models available for handling seasonal components (snaive, SARIMA, all ETS models with seasonal component) assigned a high importance to the additional features related to seasonality such as ACF, PACF-based features related to seasonal lag or seasonally differenced series. Furthermore, as expected features calculated based on parameter estimated of ETS(A, A, A) have been ranked as important for the choice of ETS with damped trend and seasonal component and ETS with trend and seasonal component.

Figure 8 and Figure 9 show the partial dependency functions of the features that get selected most often in the top. Except for random walk, partial dependency curves of seasonality and trend show a similar behaviour for both quarterly and monthly data. Hence, for seasonality and trend, the partial dependency curves computed based quarterly are presented except for random walk. Probability of selecting a model with a parameter to handle the seasonal effect (snaive, all ETS models with seasonal component, SARIMA, tbats, theta, stlar) increases as the seasonality increases. Further for classes rwd, all ETS model with trend component, SARIMA, ARIMA, tbats and theta probability of being selected increases as the strength of trend increases. On the other hand, opposite relationships were observed for snaive and ETS-seasonal which accounted seasonality only. This confirms the idea that the choice of model selection consistent with the expected relationship. According to the FFORMS framework trained on quarterly data probability of selecting random walk models remains stable up to 0.85 and drops sharply afterwards, whereas the FFORMS framework trained on monthly data indicates probability of selecting random walk models increases as trend increases. This could be due to the interaction effect of trend with other features. As shown in Figure 10 and Figure 11 trend shows a high interactivity with other features within random walk class. The characterization of the relationship with linearity differ between quarterly and monthly data for snaive, random walk, ETS-DTS, SARIMA, stlar, tbats, white noise and neural network class. Figure 10 and Figure 11 show the heat maps of relative strength of all possible pairwise interactions calculated based on Friedman's H-statistic for quarterly and monthly data respectively. Within all classes seasonality show less interactivity with other features, whereas lumpiness show high interactivity with other features. In addition to that for quarterly diff1y_pacf5 shows a very high interactivity with other features and for monthly data within many classes linearity, spikiness, y_pacf5, and ACF/PACF-based features related to seasonal lags and beta shows an interactivity with all features. For quarterly data y_acf1 and y_acf5 show higher interactivity within random walk with drift, ETS-seasonal and SARIMA while for monthly all classes show a high interactivity between y_acf1 and y_acf5. Further within each class, a subset of ACF/PACF-based features shows some interactivity. In general interactivity between features related to correlation structure of a time

series and overall shape (spikiness, linearity, curvature, etc) lead to the choice of forecast-model selection.



Figure 5: Distribution of proportion of times each time series was assigned to each class based on OOB sample. Each row represent the predicted class label and colours of boxplots corresponds to class label of “best” forecasting method. X-axis denotes the proportion of times the time series was predicted into each class. Results for quarterly and monthly series are shown in Figures A and B respectively.

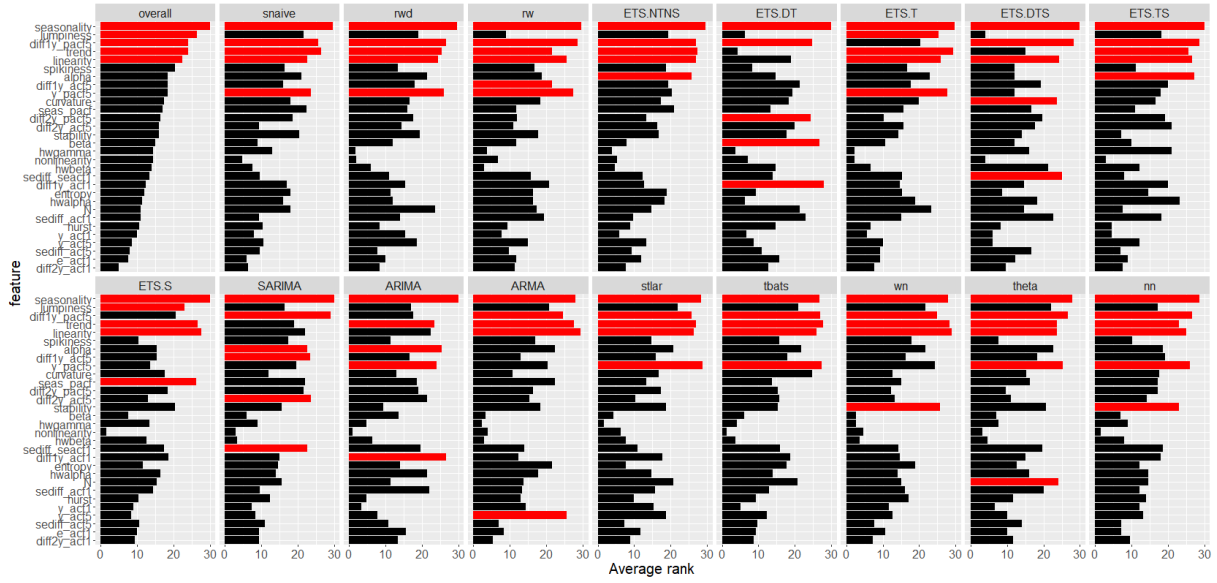


Figure 6: Feature importance plot for quarterly data. Permutation-based VI measure and mean decrease in Gini coefficients are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on the three measures: permutation-based VI, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.



Figure 7: Feature importance plot for monthly data. Permutation-based VI measure and mean decrease in Gini coefficients are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on the three measures: permutation-based VI, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.

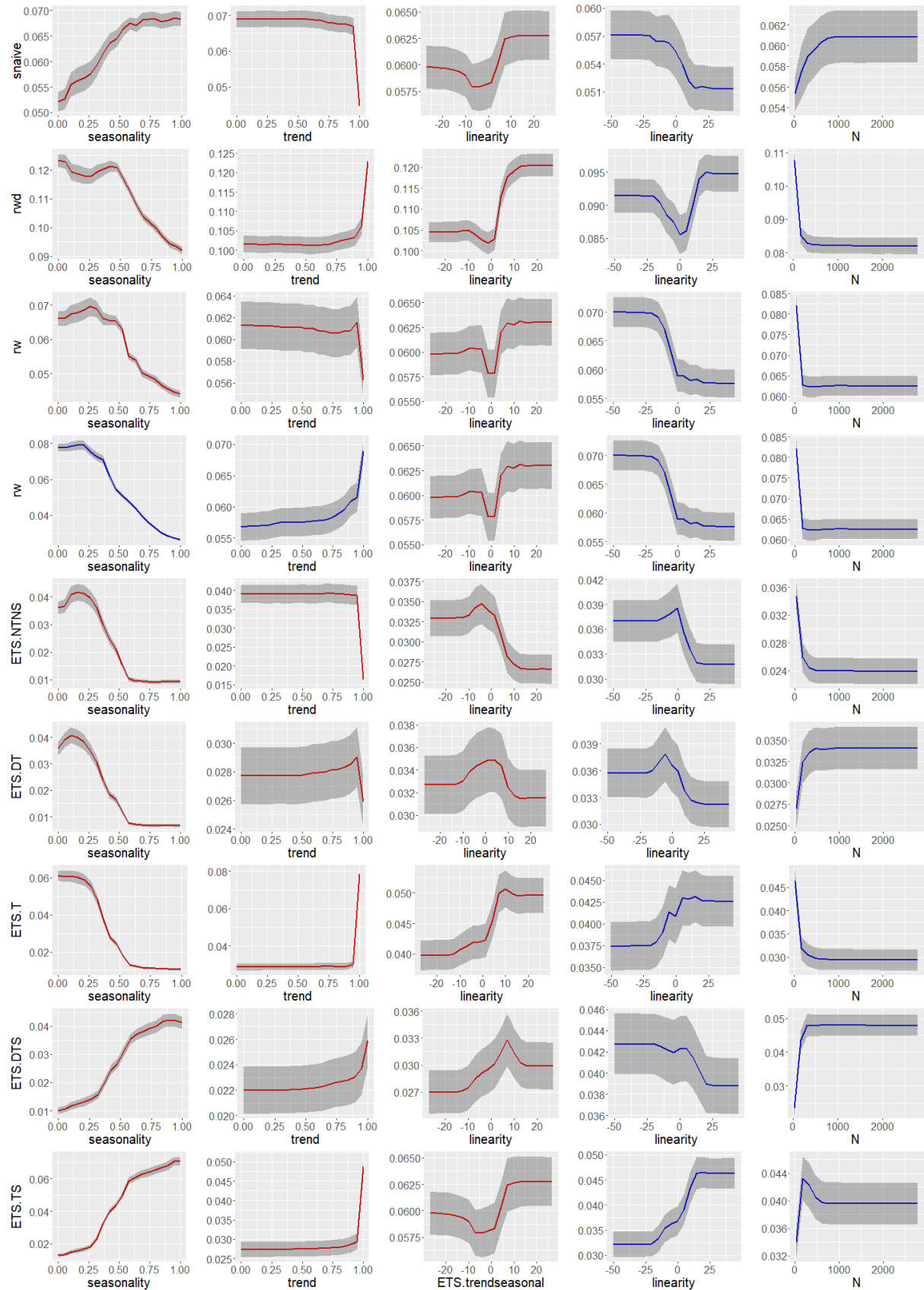


Figure 8: Partial dependence plots for the top-three features get selected most within each class inside quarterly and monthly FFORMS frameworks. Additionally, N is included to observe the effects stated in the literature. The shading shows the 95% confidence interval. Y-axis denotes the probability of belonging to corresponding class. Red colour is for PDP drawn based on quarterly data and blue colour is for the PDP drawn based on monthly data.

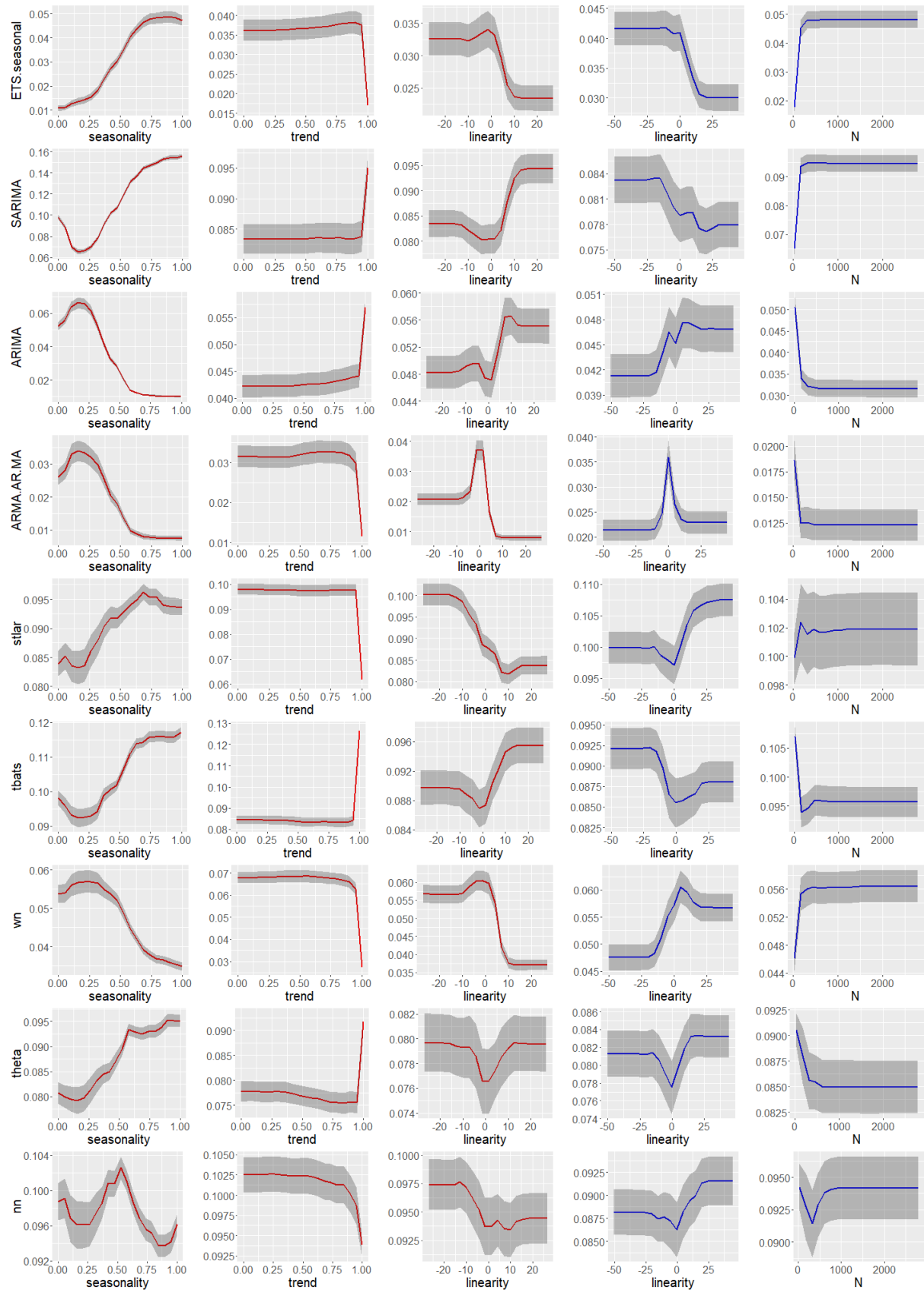
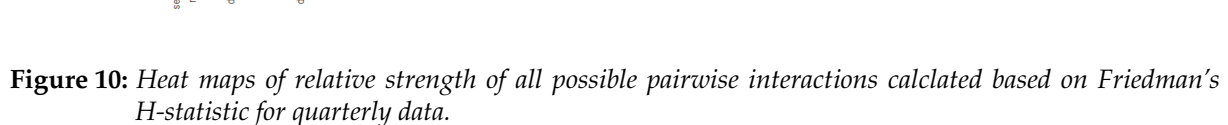


Figure 9: Partial dependence plots for the top-three features get selected most within each class inside quarterly and monthly FFORMS frameworks. Additionally, N is included to observe the effects stated in the literature. The shading shows the 95% confidence interval. Y-axis denotes the probability of belonging to corresponding class. Red colour is for PDP drawn based on quarterly data and blue colour is for the PDP drawn based on monthly data.



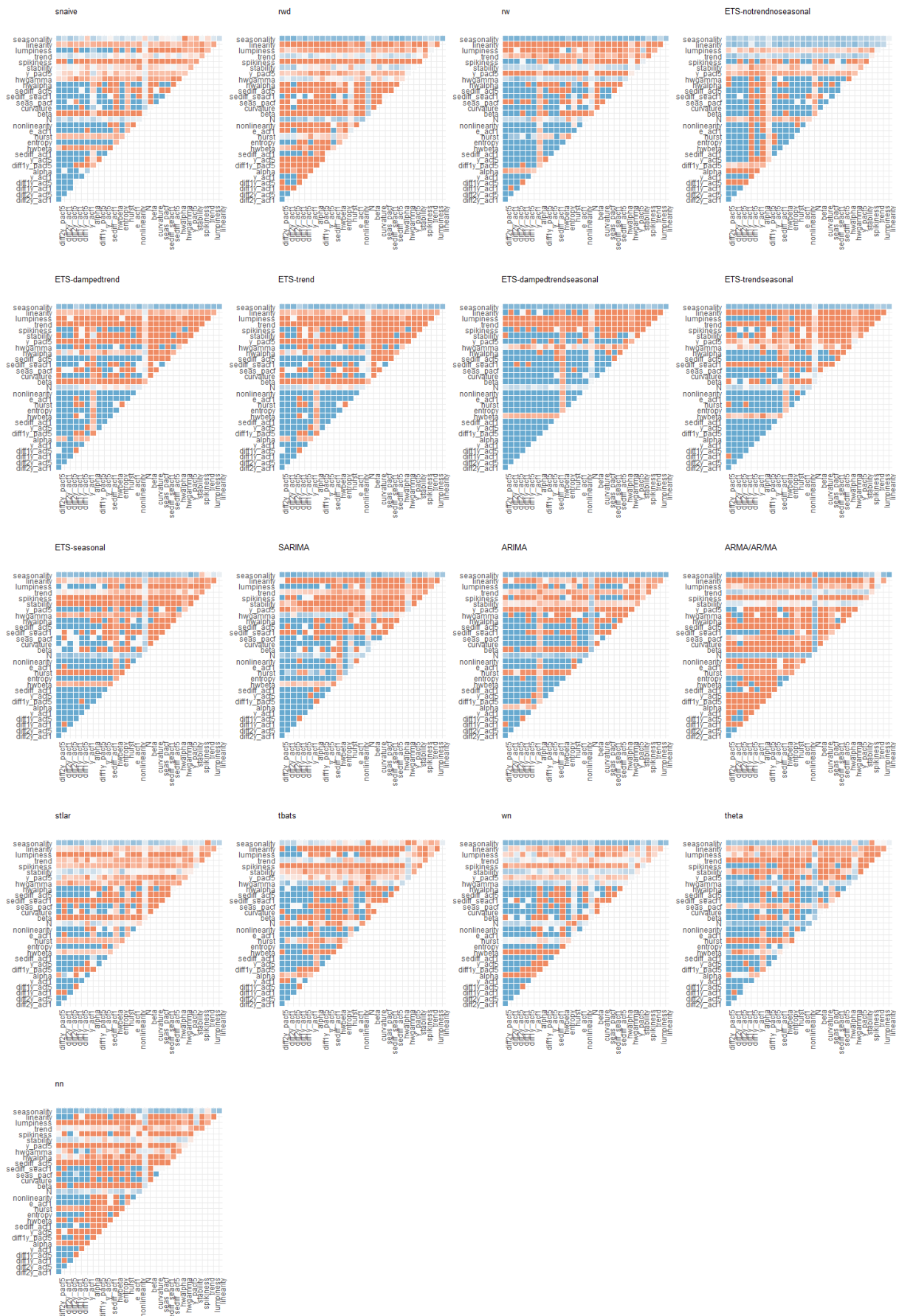


Figure 11: Heat maps of relative strength of all possible pairwise interactions calculated based on Friedman's H -statistic for monthly data.

4.3 Weekly

Figure 12 proportion of times each time series was classified to each class. Unlike, yearly, quarterly and monthly data theta models have low chance of being selected to forecast weekly data. The random walk with drift, tbats models and theta have higher chance of being selected. Except ARMA/AR/MA class the distribution corresponds to the true class label dominates others. ARMA/AR/MA class shows some unusual behaviour within some categories due to class imbalance ratio, ARMA/AR/MA class contains fewer number of observations in the training set.

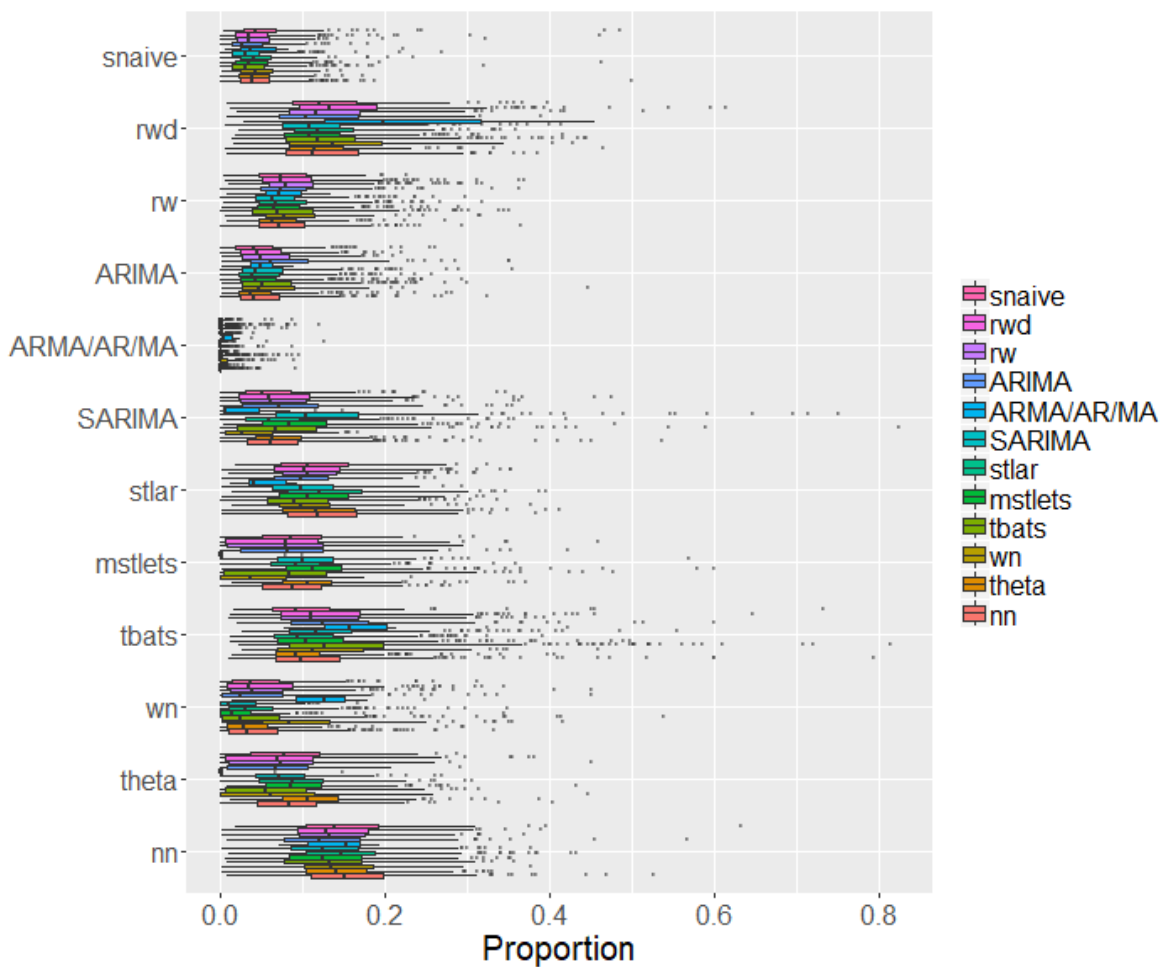


Figure 12: Distribution of proportion of times each time series was assigned to each class based on OOB sample. Each row represent the predicted class label and colours of boxplots corresponds to class label of “best” forecasting method. X-axis denotes the proportion of times the time series was predicted into each class.

According to the results of Figure 13 spikiness, linearity, trend, strength of seasonality, stability and lumpiness has been assigned a high importance. This is similar to the results of yearly, quarterly and monthly data. The length of series has been selected among top 5 by mstlets, tbats, theta and neural network models. According to the results of Figure 14 for mstlets models probability of being selected increases as the length of series increases while the opposite relationship was observed at neural network models. The tbats models show a non-monotonic relationship with the length of series. It is surprising to observe that for mstlets models probability of being selected decreases as seasonality increases. This could be due to the interaction effect of seasonality with others. For weekly data, the number of pairs showing a particular degree of interaction strength is relatively low compared to other yearly, quarterly and monthly data. In tbats class, trend, linearity, spikiness, seasonality and first autocorrelation coefficient of the seasonally differenced series interact heavily with other features.

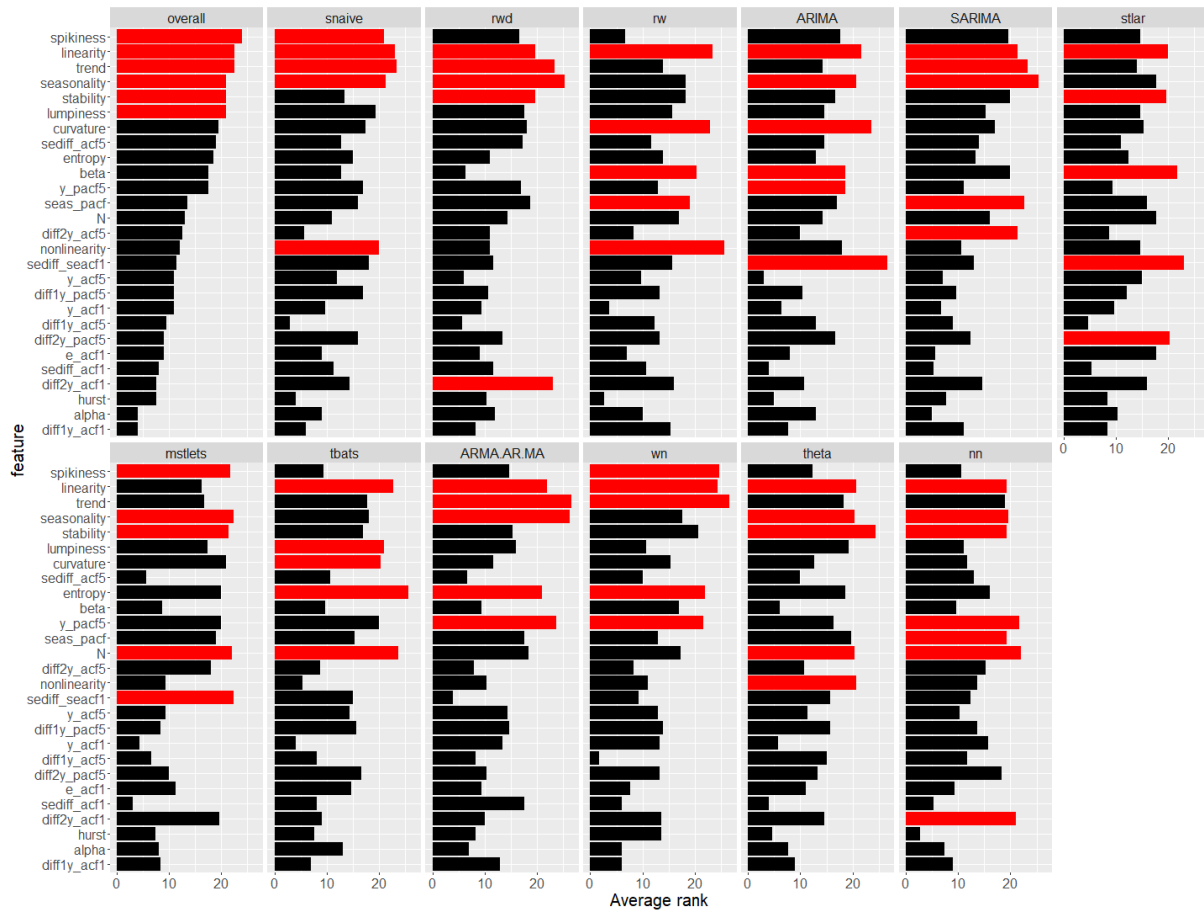


Figure 13: Feature importance plot for weekly data. Permutation-based VI measure and mean decrease in Gini coefficient are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on the three measures: permutation-based VI, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.

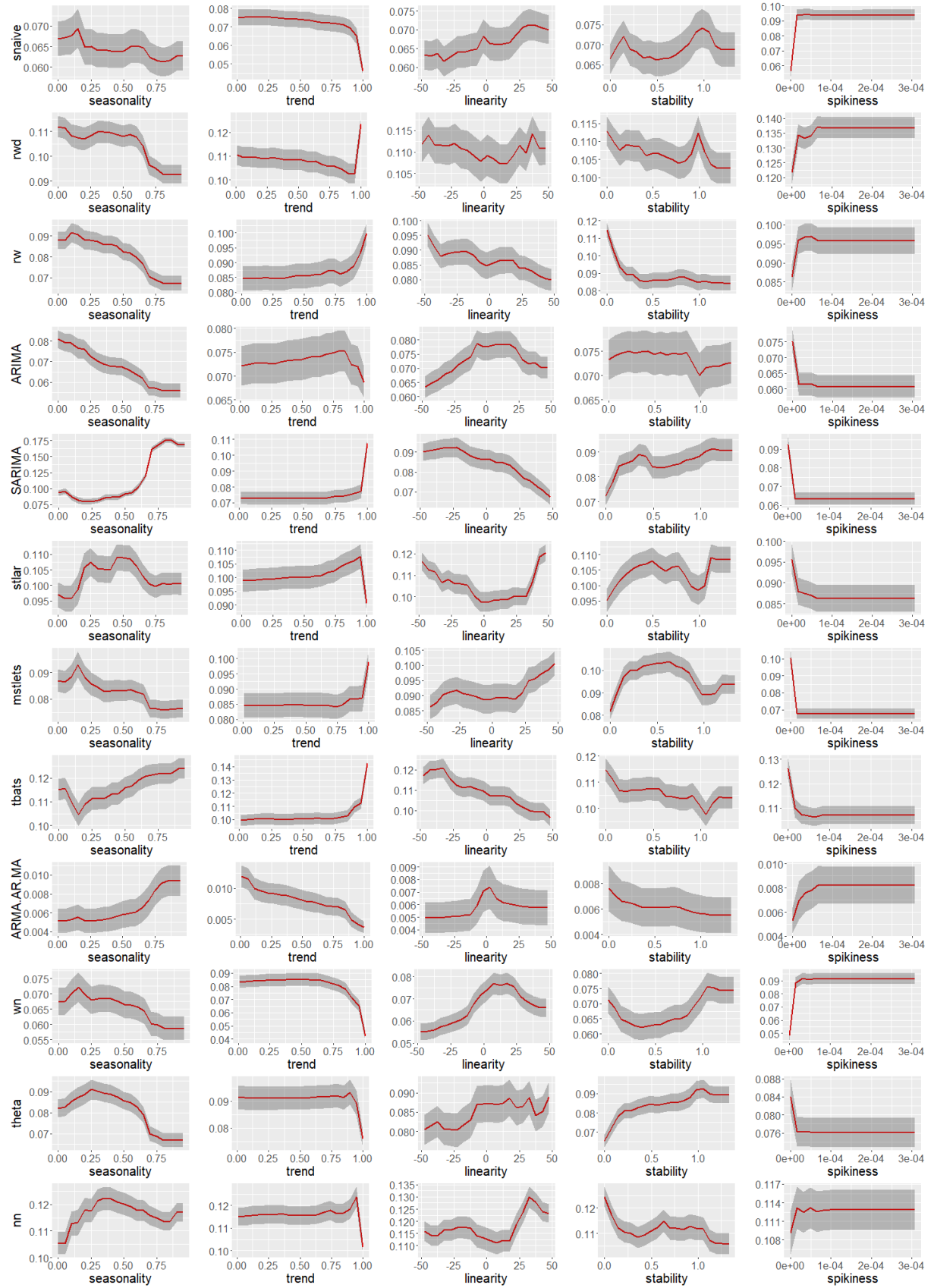


Figure 14: Partial dependence plots for the top ranked features from variable importance measures (weekly series). The shading shows the 95% confidence intervals. Y-axis denotes the probability of belong to corresponding class.



4.4 Daily and Hourly data

According to [Figure 16](#) the distributions corresponds to observations that have been correctly classified dominate the top for daily data. However, within daily series there are few observations that have been incorrectly classified to tbats class with very high probabilities. In general, neural network models have higher chance of being selected for daily time series. Overall, for hourly series random walk with drift models, tbats and neural network models have high chance of being selected. Furthermore, it is important to note that all hourly series have been assigned non-zero probability of getting selected to neural network class. Variable importance graph for daily and hourly data are shown in [Figure 17](#) and [Figure 18](#) respectively. The most important features for determining suitable forecast-models for daily time series are, strength of seasonality corresponds to the weekly seasonality (7), stability, trend, lumpiness and linearity. Strength of seasonality corresponds to the annual seasonality (365.25) is appeared to be important in determining the selection of theta models. Furthermore, length of the series is important in determining random walk, random walk with drift, mstlarima, mstlets, stlar, theta and nn classes. [Figure 19](#) shows the partial dependency plots of the top 3 features from the FFORMS framework. According to the results of [Figure 19](#) shorter series tends to select random walk with drift models while probability of selecting snaive, mstlarima and mstlets models increases as the length of series increases. Neural network models shows a non-monotonic relationship with length of the series (N). The theta models tends to be selected for series with high annual seasonality but very low weekly seasonality. According to [Figure 18](#), strength of daily seasonality (period=24) appear to be more important than than the strength of weekly seasonality (period=168). Furthermore, entropy, linearity, sum of squares of first 5 coefficients of PACF, curvature, trend, spikiness and stability were found to be the most important features in determining best forecasting method for hourly time series. Only snaive category ranked it among top 5 for hourly time series. The strength of seasonality corresponds to weekly seasonality also seems to be one of the most important feature for the classes snaive, random walk, mstlarima, and tbats. According to [Figure 20](#) probability of selecting random walk, random walk with drift, theta model and white noise process decreased with higher strength of daily seasonality, while the opposite relationship hold for other classes. On the other hand, probability of selecting random walk model increased with the increase in strength of weekly seasonality.

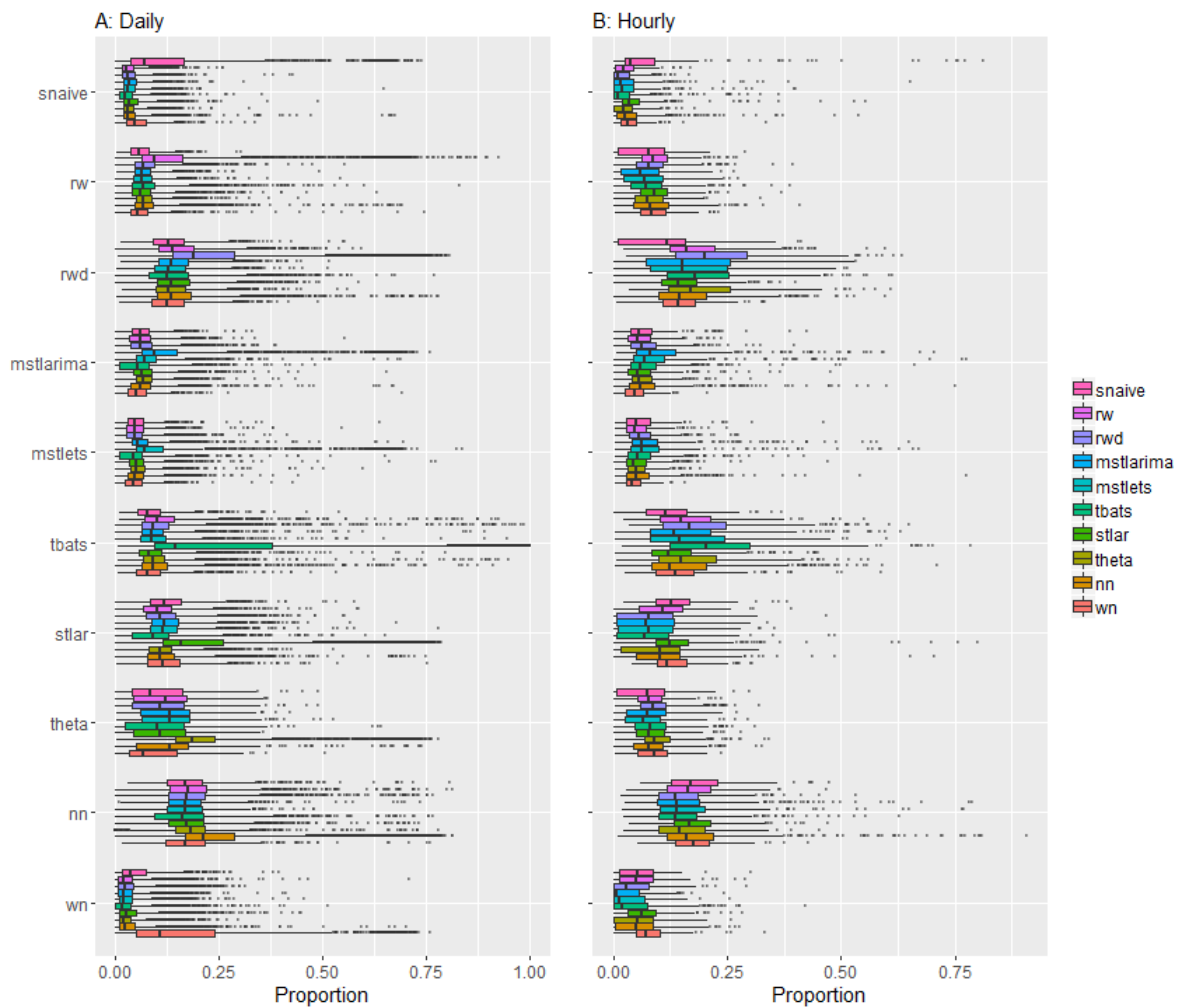


Figure 16: Distribution of proportion of times each daily time series was assigned to each class in the forest. Each row represents the predicted class label and colour of boxplots corresponds to true class label. There are ten rows in the plot corresponds to each predicted class represented by Y-axis. X-axis denotes the proportion of times a time series is classified in each class. On each row, distribution of correctly classified class dominated the top, indicating a fairly good classification of the model fitted.

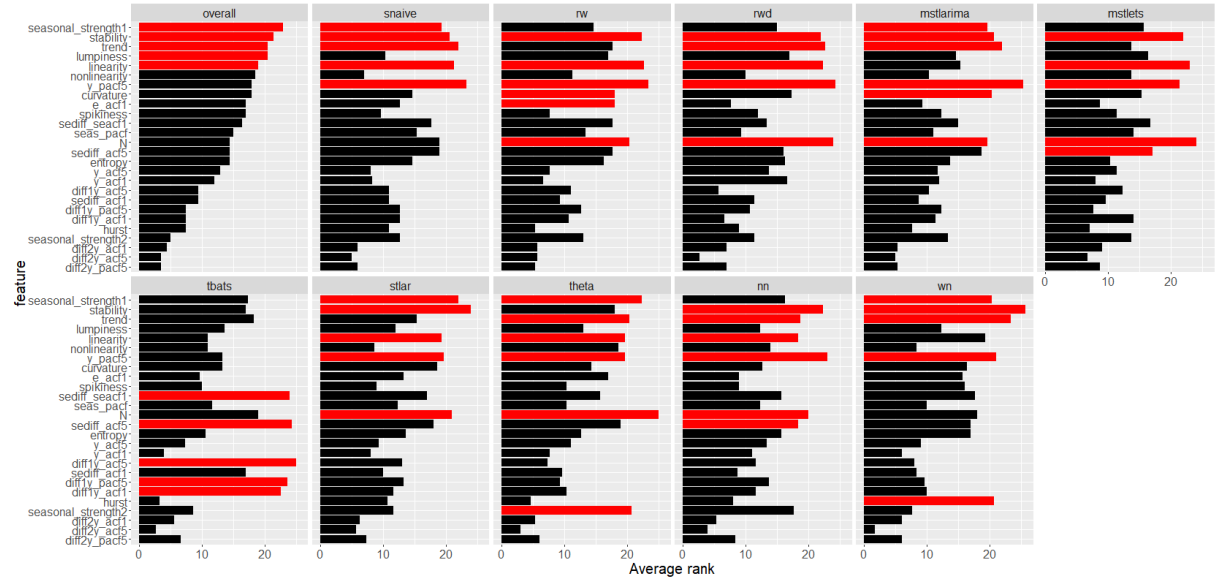


Figure 17: Feature importance plot for daily data. Permutation-based VI measure and mean decrease in Gini coefficients are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on the three measures: permutation-based VI, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.

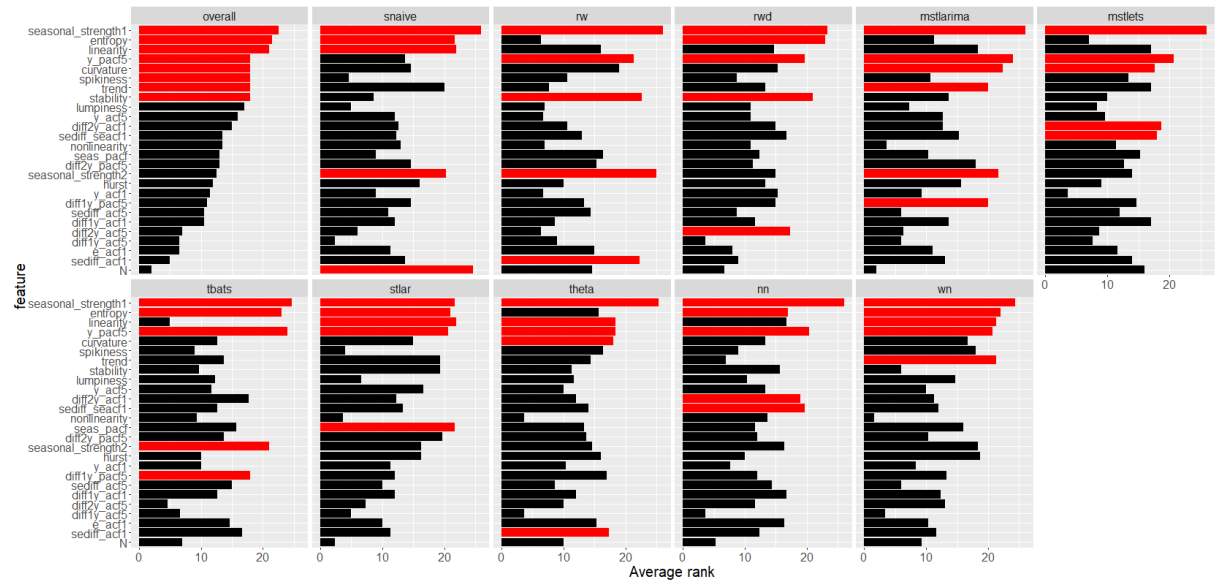


Figure 18: Feature importance plot hourly series. Permutation-based VI measure and mean decrease in Gini coefficients are used to evaluate overall feature importance. Class-specific feature importance is evaluated based on the three measures: permutation-based VI, PD-based VI measure, and ICE-based VI measure. Longer bars indicate more important features. Top 5 features are highlighted in red.

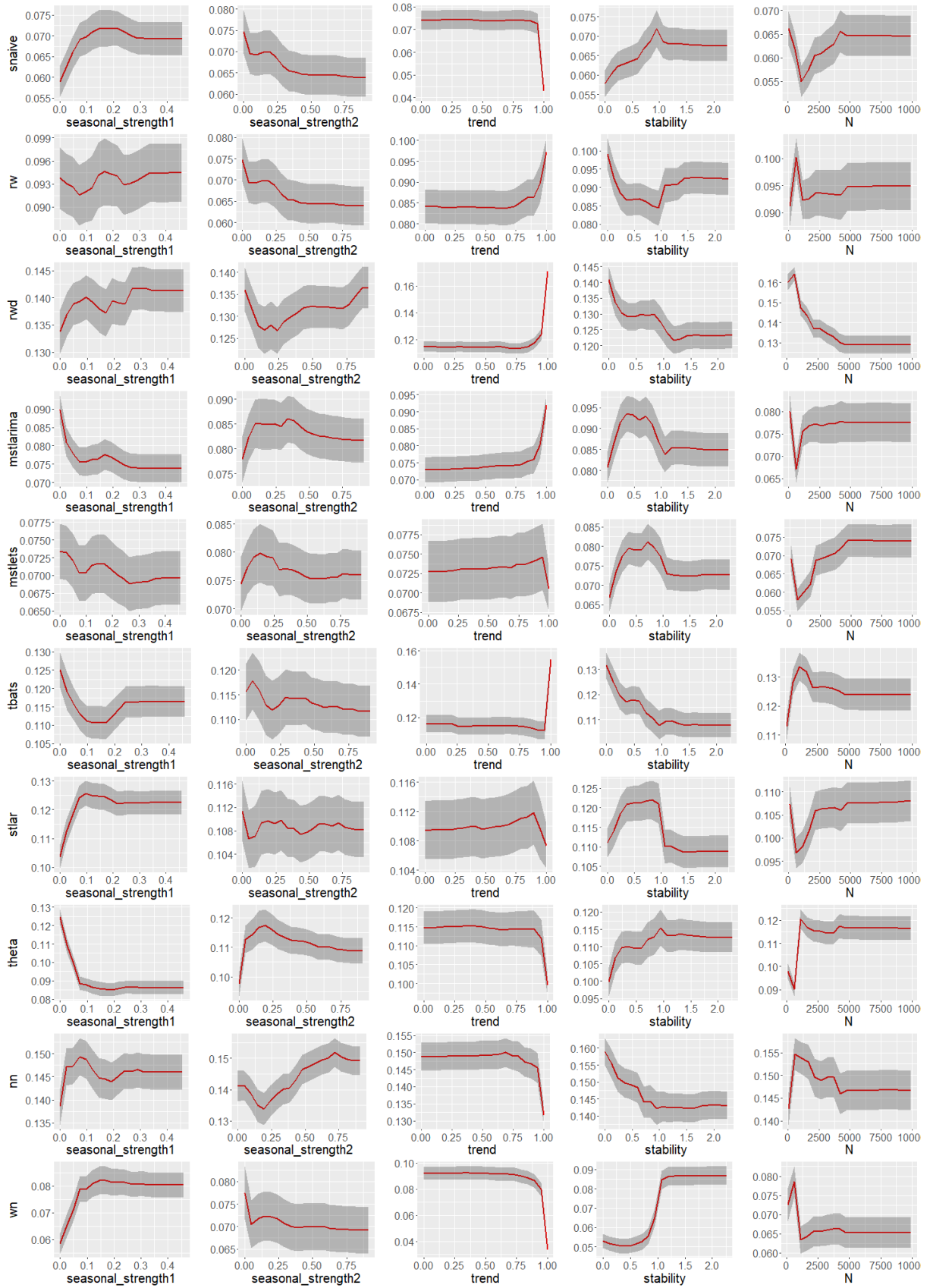


Figure 19: Partial dependence plots for the top ranked features from variable importance measures(daily series). The shading shows the 95% confidence intervals. Y-axis denotes the probability of belong to corresponding class.

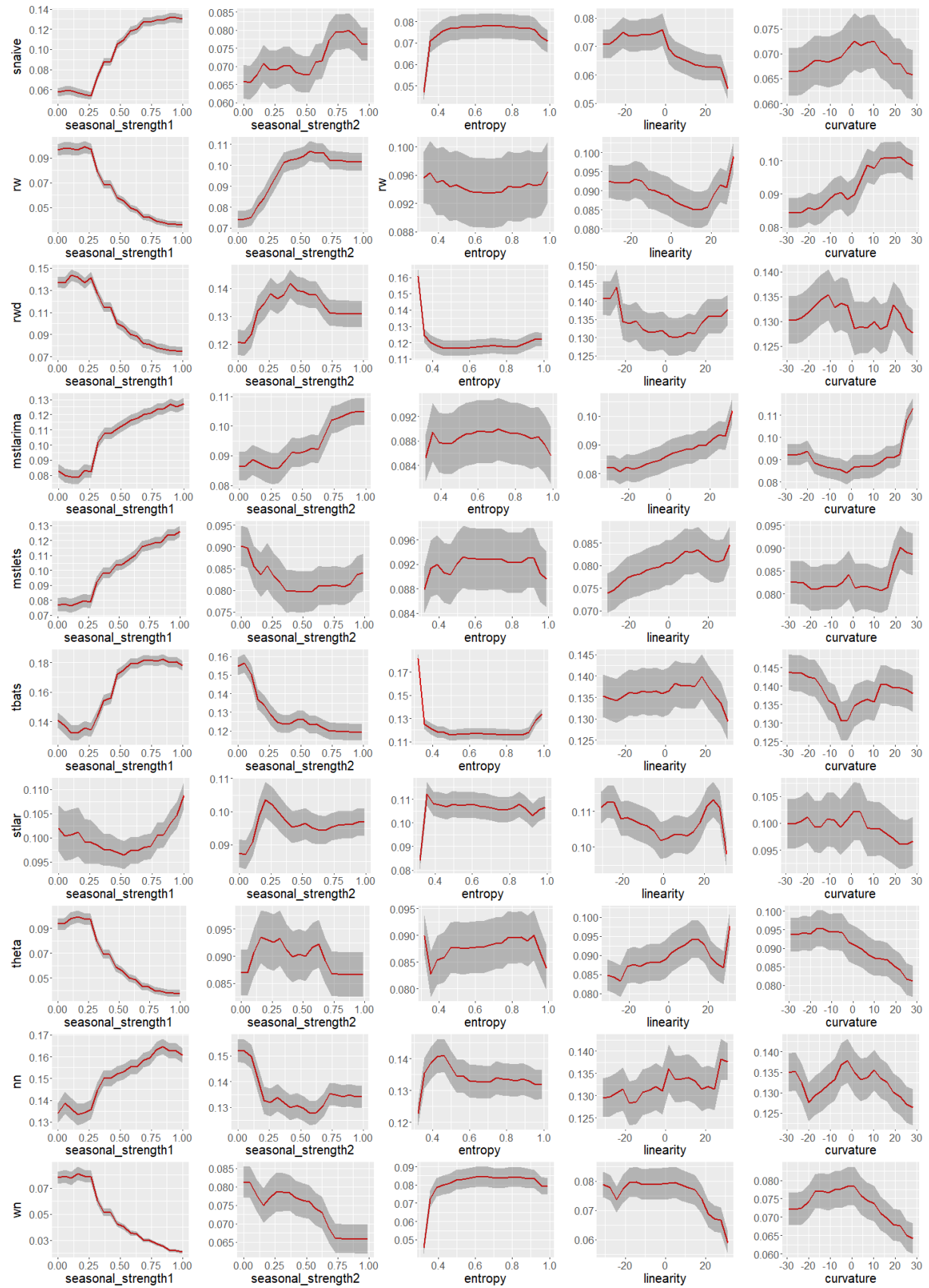


Figure 20: Partial dependence plots for the top ranked features from variable importance measures(daily series). The shading shows the 95% confidence intervals. Y-axis denotes the probability of belong to corresponding class.

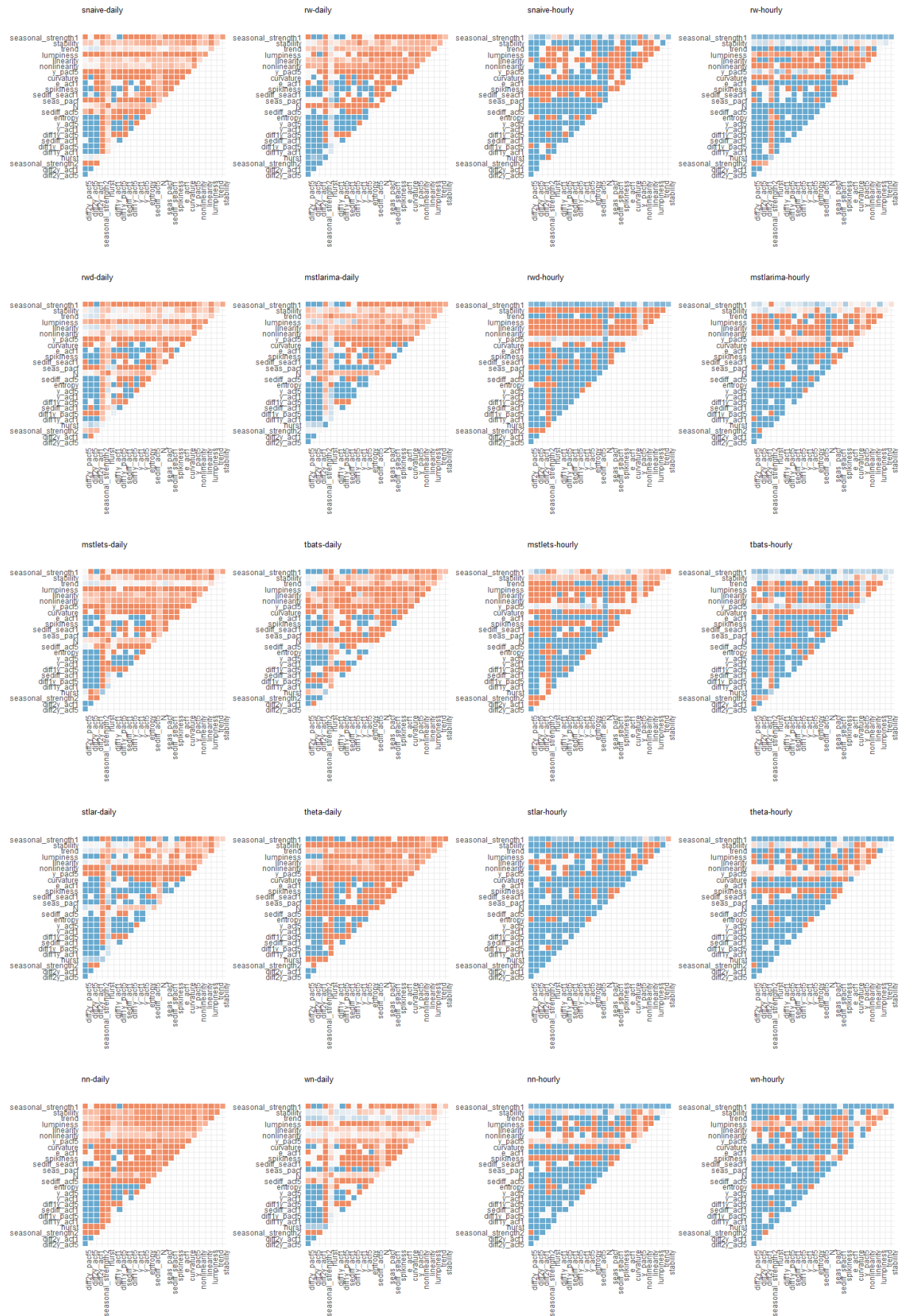


Figure 21: Heat maps of relative strength of all possible pairwise interactions calculated based on Friedman's H -statistic for daily data.

4.5 Local Interpretable Model-agnostic Explanations

Figure 22 shows the feature contribution for the instances highlighted in the PCA-space. The low linearity, with a trend value of less than 0.616, low values for Phillip-Perron test-statistic, Hurst exponent, autocorrelation coefficients of the time series and spikiness of the series causes the FFORMS framework to classify the first series with ARMA/AR/MA. On the other hand, LIME indicated the low value of first autocorrelation coefficients of the original series and the residual series of linear regression model contradicts to this decision. It is interesting to note that even though 2, 4 and 5 series located close proximity to each other in the PCA space their varying degree of trend, spikiness, first autocorrelation coefficient of the difference series, entropy, led them to classify as ETS-trend, neural-network and random walk respectively. From this approach we can gain insight into the local neighbourhood characteristics which lead to the choice of a particular neighbourhood over alternative destinations.

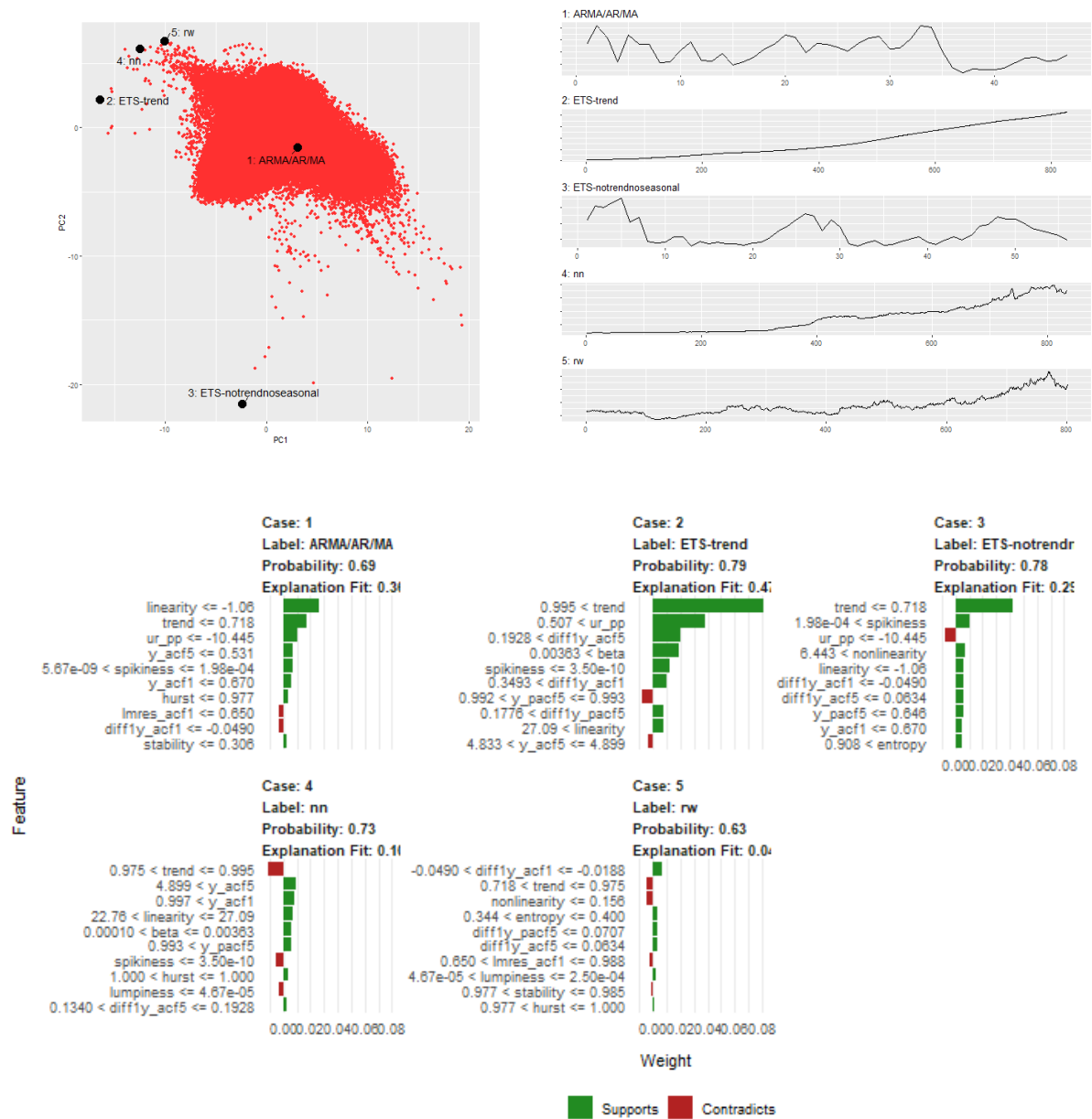


Figure 22: Local interpretable Model-agnostic explanations for six selected yearly time series. Features denoted with green colour are supporting features for an outcome label and length of the bar is proportional to the weight of a feature.

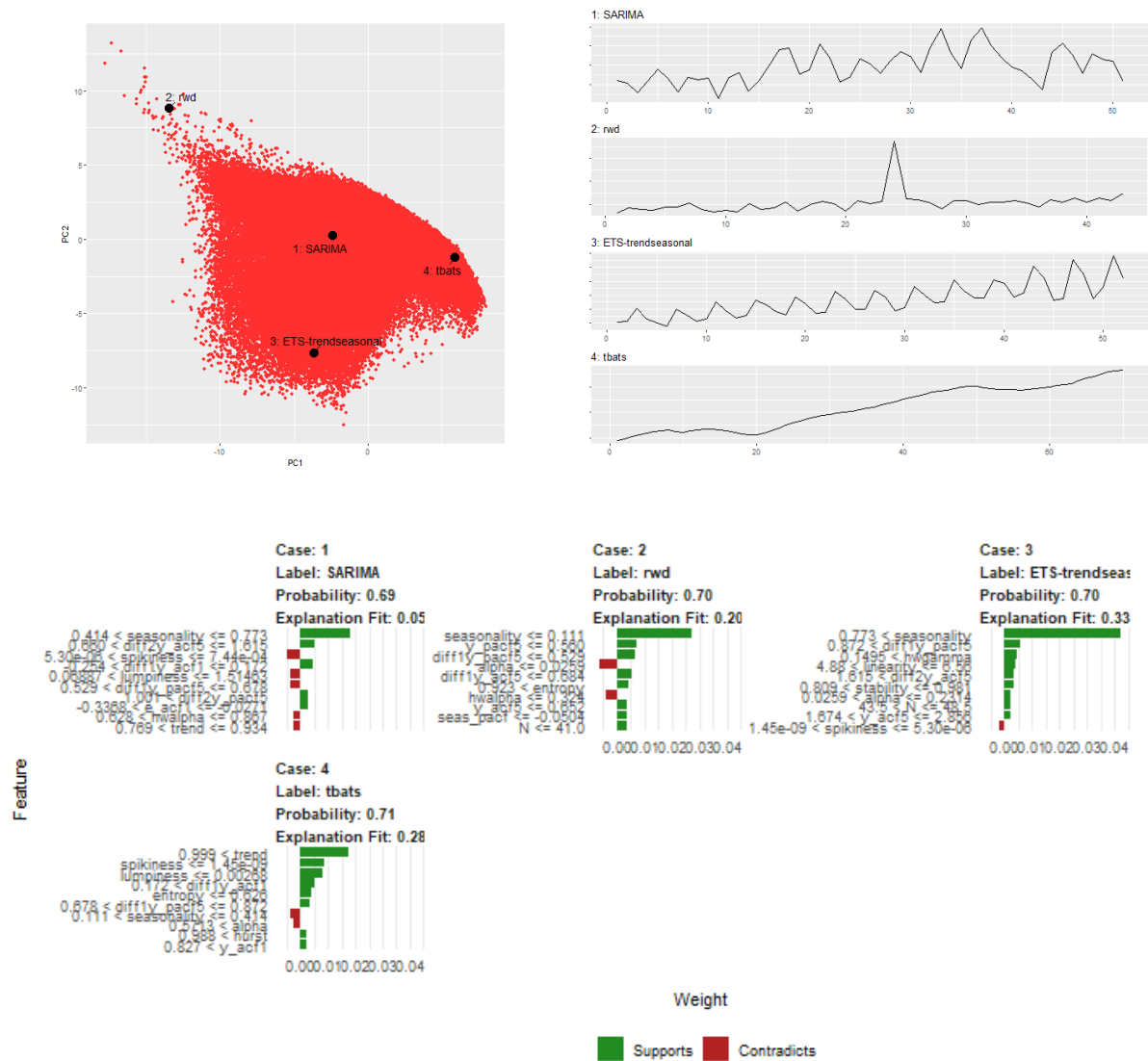


Figure 23: Local interpretable Model-agnostic explanations for six selected quarterly time series. Features denoted with green colour are supporting features for an outcome label and length of the bar is proportional to the weight of a feature.

5 Discussion and Conclusions

Forecast model selection is both time and computer cost intensive process. Consequently, the application of machine-learning approaches to predict suitable forecasting model from large number of potentially relevant time series features is a topic growing popularity in the field of time series forecasting. Recently we introduced a computationally efficient framework for individual forecast models selection based on the features computed from the time series. We called this framework FFORMS: Feature-based FOfRecast Model-Selection. The M4-competition submission based on FFORMS framework was placed eighth under prediction interval category. In this paper we used model-agnostic machine learning interpretability tools to explore what was happening under the hood of FFORMS framework and to gain an understanding of what features led to the choices of FFORMS framework. On the other hand explaining predictions is an important aspect in getting humans trust and use the proposed framework effectively, if the explanations are faithful and intelligible. Humans usually have prior knowledge about the application domain, which they can use to accept (trust) or reject prediction if they understand the reasons behind it.

We explored the role of features in two different perspectives: i) individual effect of feature, and ii) interaction effect of features. Overall, the features strength of trend, strength of seasonality, linearity, spikiness and curvature among the top 10 within each frequency category. Lemke & Gabrys (2010) also pointed out features related to nonstationarity and seasonality of a series are important factors for choosing a forecasting method. Features that frequently appear can be considered as more relevant than those that tend to appear less frequently. Partial dependency plots are used to visualize the learned relationship between features and the model predictions. The displayed relationships confirm to domain knowledge expectations. However, since several number of features are used to build the framework with comparable contributions, and thus, all individual contributions are small. According to the results of daily and hourly data we also observed neural network modelling model was appropriate for forecasting high frequency data. In response to the results of the M3-competition this has been pointed out by many commentators (Makridakis & Hibon 2000). Further, our results show that the performance of various methods depends upon the length of the time series. Short time series tends to select simple methods such as random walk models, naive, etc. ETS models with both trend and seasonal components, SARIMA models, mstl models tend to provide accurate forecasts with longer time series as these are more parameterized models.

As FFORMS framework is developed on top of random forest algorithm takes into account

every possible interactions. We used Friedman's H-statistic to identify most important two-way interactions. It was apparent from the heat matrices of Friedman's H-Statistic presented that a substantial interaction effect exist between the features. The strength trend showed less interactivity in yearly series data, reflecting that these features are more important on their own. The features involve in interaction and their strength of interaction effect differ across the different frequency categories (yearly, quarterly, monthly, weekly, daily and hourly) as well as forecast-models (random walk, ETS models, etc.). However, it is interesting to note that in each frequency category all or subset of ACF/PACF-based features interact each other. This confirm that information regarding correlation structure of the time series is an essential information for the choice of model selection. [Figure 24](#) - [Figure 28](#) show partial dependency plot for most frequently appeared interacting feature combination in each of the frequency categories. According to figures [Figure 24](#) - [Figure 28](#) the unique pattern of interactivity exist within each class are useful for separating one from another.

Exploration of conditions learnt by the FFORMS framework also support practitioners to make a good educated guess on suitable forecast-model for a given problem. Further the results of this study is useful in identifying new ways to improve forecasting accuracy by capturing different features of time series.

Appendix

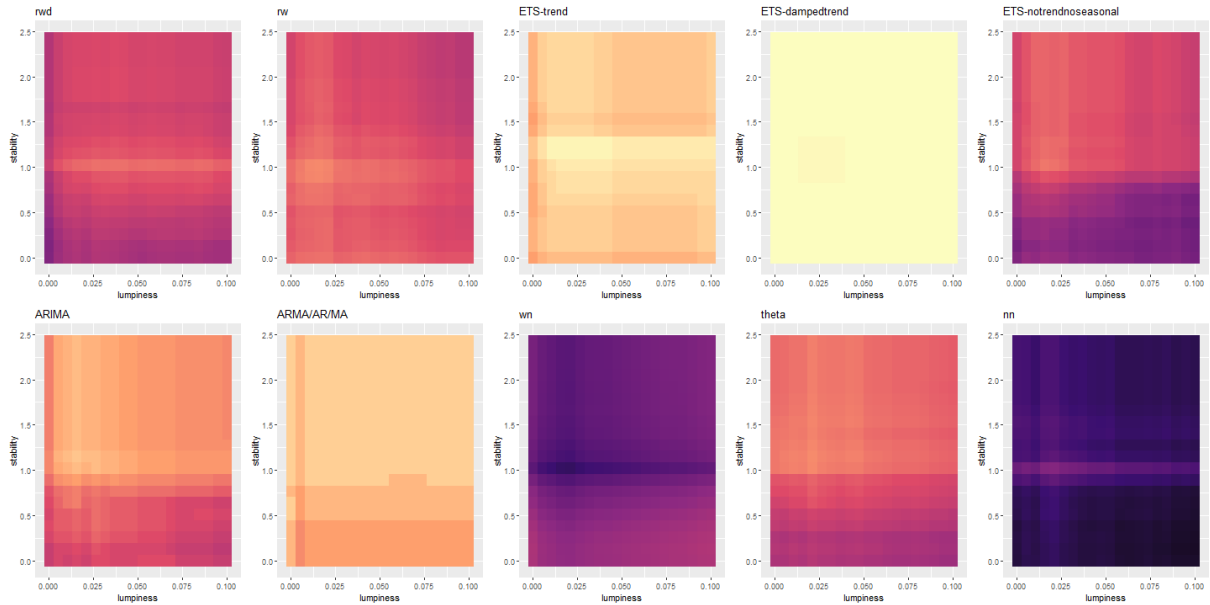


Figure 24: Partial dependence plot of model selection probability and the interaction of stability and lumpiness for yearly data.

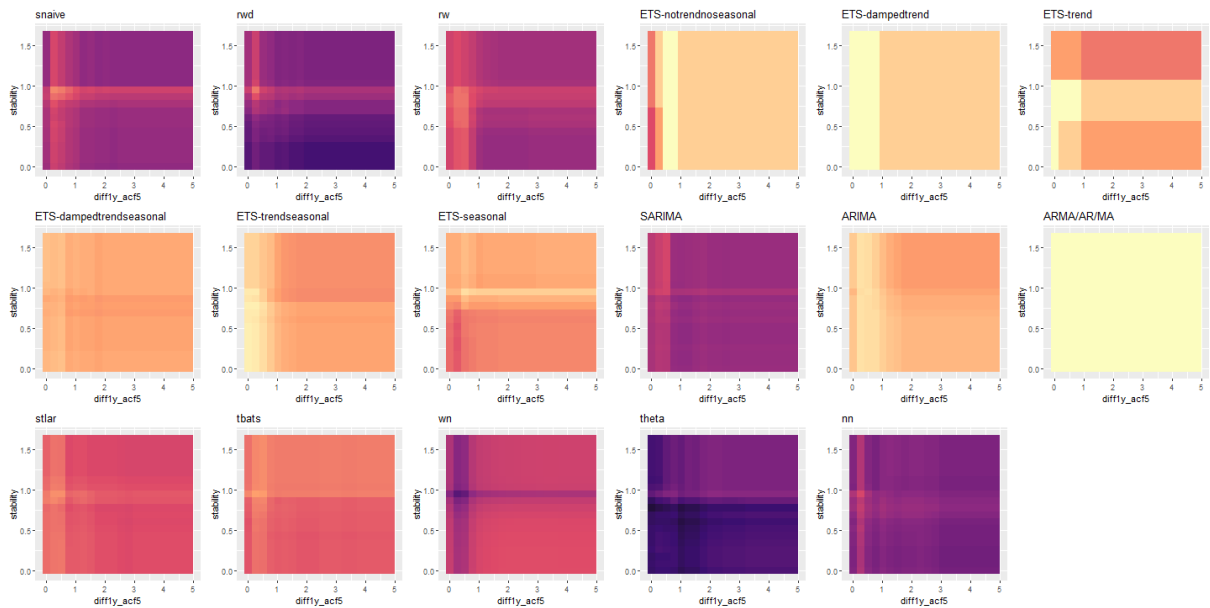


Figure 25: Partial dependence plot of model selection probability and the interaction of stability and diff1y_acf5 for quarterly data

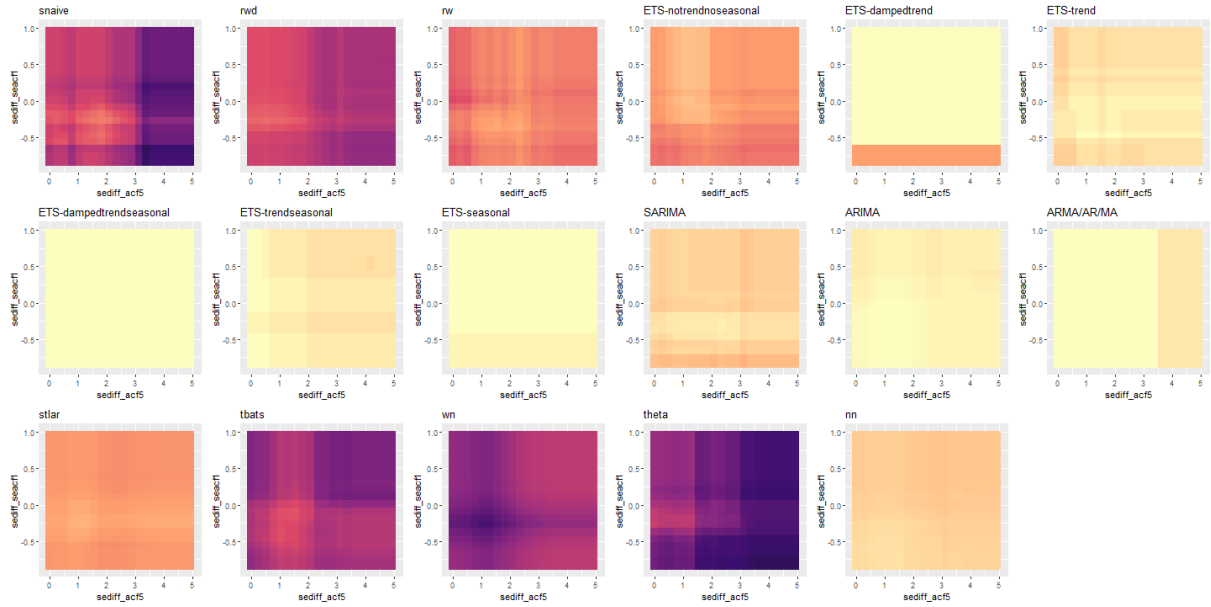


Figure 26: Partial dependence plot of model selection probability and the interaction of `sediff_seacf1` and `sediff_acf5` for monthly data

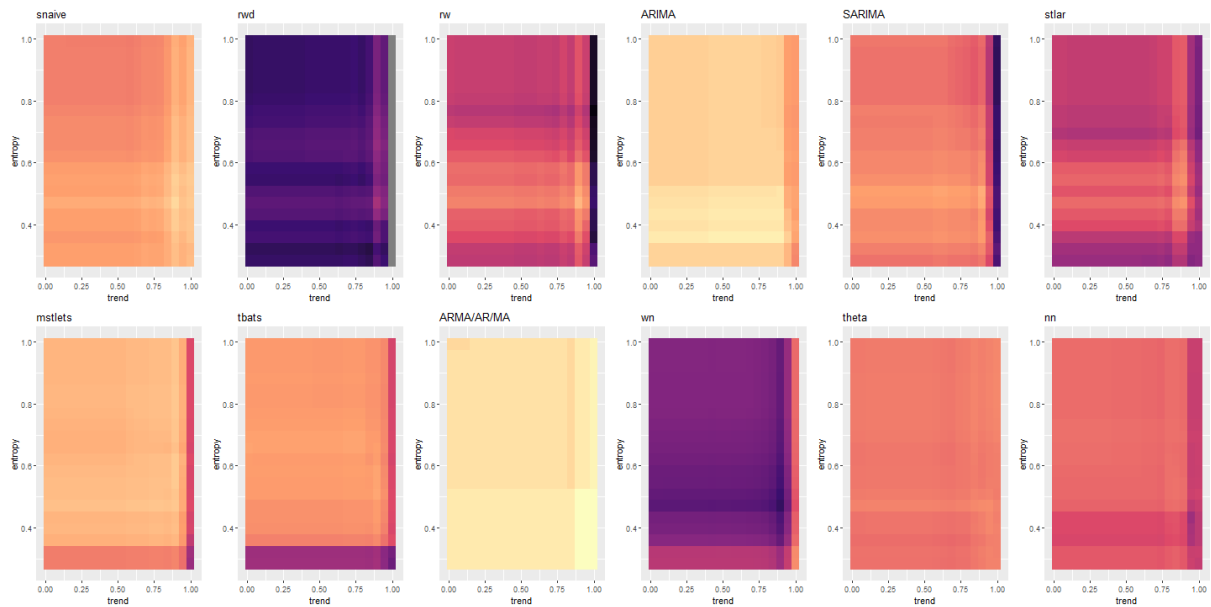


Figure 27: Partial dependence plot of model selection probability and the interaction of `trend` and `entropy` for weekly data.

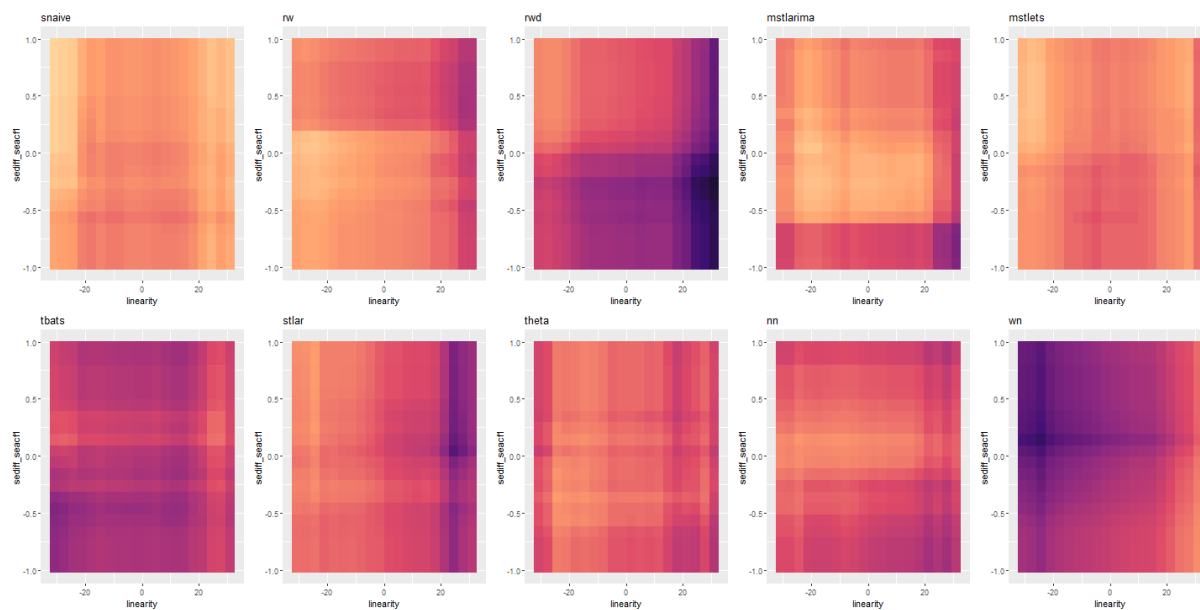


Figure 28: *Partial dependence plot of model selection probability and the interaction of `sediff_seacf1` and `linearity` for hourly data.*

References

- Breiman, L (2001). Random forests. *Machine Learning* **45**(1), 5–32.
- Chen, C, A Liaw & L Breiman (2004). *Using random forest to learn imbalanced data*. Tech. rep. University of California, Berkeley. <http://statistics.berkeley.edu/sites/default/files/tech-reports/666.pdf>.
- Ehrlinger, J (2015). ggRandomForests: Visually Exploring a Random Forest for Regression. *arXiv preprint arXiv:1501.07196*.
- Friedman, JH, BE Popescu, et al. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics* **2**(3), 916–954.
- Goldstein, A, A Kapelner, J Bleich & E Pitkin (2015). Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics* **24**(1), 44–65.
- Greenwell, BM, BC Boehmke & AJ McCarthy (May 2018). A Simple and Effective Model-Based Variable Importance Measure.
- Hyndman, R, G Athanasopoulos, C Bergmeir, G Caceres, L Chhay, M O’Hara-Wild, F Petropoulos, S Razbash, E Wang & F Yasmeeen (2018). *forecast: Forecasting functions for time series and linear models*. R package version 8.3. <http://pkg.robjhyndman.com/forecast>.
- Jiang, T & AB Owen (2002). Quasi-regression for visualization and interpretation of black box functions. *Technical Report, Stanford University*.
- Kang, Y, RJ Hyndman, F Li, et al. (2018). *Efficient generation of time series with diverse and controllable characteristics*. Tech. rep. Monash University, Department of Econometrics and Business Statistics.
- Kück, M, SF Crone & M Freitag (2016). Meta-learning with neural networks and landmarking for forecasting model selection an empirical evaluation of different feature sets applied to industry data. In: *Neural Networks (IJCNN), 2016 International Joint Conference on*. IEEE, pp.1499–1506.
- Lemke, C & B Gabrys (2010). Meta-learning for time series forecasting and forecast combination. *Neurocomputing* **73**(10), 2006–2016.
- Lundberg, SM & SI Lee (2017). A unified approach to interpreting model predictions. In: *Advances in Neural Information Processing Systems*, pp.4765–4774.
- Makridakis, S & M Hibon (2000). The M3-Competition: results, conclusions and implications. *International journal of forecasting* **16**(4), 451–476.

- Meade, N (2000). Evidence for the selection of forecasting methods. *Journal of forecasting* **19**(6), 515–535.
- Molnar, C, G Casalicchio & B Bischl (2018). Iml: An R package for interpretable machine learning. *The Journal of Open Source Software* **3**(786), 10–21105.
- Petropoulos, F, S Makridakis, V Assimakopoulos & K Nikolopoulos (2014). ‘Horses for Courses’ in demand forecasting. *European Journal of Operational Research* **237**(1), 152–163.
- Prudêncio, RB & TB Ludermir (2004). Meta-learning approaches to selecting time series models. *Neurocomputing* **61**, 121–137.
- Ribeiro, MT, S Singh & C Guestrin (2016). Why should I trust you?: Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD’16)*. ACM, New York, NY, USA. ACM, pp.1135–1144. <https://doi.org/10.1145/2939672.2939778>.
- Schnaars, SP (1984). Situational factors affecting forecast accuracy. *Journal of marketing research*, 290–297.
- Shah, C (1997). Model selection in univariate time series forecasting using discriminant analysis. *International Journal of Forecasting* **13**(4), 489–500.
- Silva, N da, D Cook & EK Lee (2017). Interactive graphics for visually diagnosing forest classifiers in R. *arXiv preprint arXiv:1704.02502*.
- Sutherland, P, A Rossini, T Lumley, N Lewin-Koh, J Dickerson, Z Cox & D Cook (2000). Orca: A visualization toolkit for high-dimensional data. *Journal of Computational and Graphical Statistics* **9**(3), 509–529.
- Talagala, TS, RJ Hyndman & G Athanasopoulos (2018). Meta-learning how to forecast time series. *Technical Report 6/18, Monash University*.
- Wang, X, K Smith-Miles & RJ Hyndman (2009). Rule induction for forecasting method selection: meta-learning the characteristics of univariate time series. *Neurocomputing* **72**(10), 2581–2594.
- Wickham, H, D Cook & H Hofmann (2015). Visualizing statistical models: Removing the blind-fold. *Statistical Analysis and Data Mining: The ASA Data Science Journal* **8**(4), 203–225.
- Zhao, Q & T Hastie (n.d.). Causal Interpretations of black-box models ().