# Forecasting Dutch inflation using machine learning methods

Robert–Paul Berben*    Rajni N. Rasiawan†    Jasper M. de Winter‡

January 7, 2025

**Abstract**

We extract tone-adjusted,

*Keywords:* Inflation forecasting, Big data, Machine learning, Random Forest, Ridge regression
*JEL Classification:* C22, C53, C55, E17, E31

---

*Economics and Research Division, De Nederlandsche Bank, Amsterdam, The Netherlands, r.p.berben@dnb.nl
†Economics and Research Division, De Nederlandsche Bank, Amsterdam, The Netherlands, r.n.rasiawan@dnb.nl
‡Economics and Research Division, De Nederlandsche Bank, Amsterdam, The Netherlands, j.m.de.winter@dnb.nl

# 1 Introduction

Inflation forecasts are important inputs for monetary policy decision making. Unfortunately, forecasting inflation turns out to be quite tricky, and very simple time-series models often provide the best forecasts. Quoting Faust and Wright (2013) "We find that [...], extremely simple inflation forecasts – that however take account of nowcasting and secular changes in the local mean inflation rate – are just about the best that are available."

In recent years, various papers have been published showing that the combination of "big data" and machine learning (ML) models can lead to more accurate inflation forecasts. Medeiros et al. (2021) consider several ML models and the FRED-MD database developed by McCracken and Ng (2016) to forecast U.S. consumer price index (CPI) inflation, and are able to outperform standard benchmarks, with the random forest (RF) model delivering the smallest errors. Naghi et al. (2024) elaborate on the work by Medeiros et al. (2021) by extending the sample from 2016 to 2022, considerably expanding the set of machine learning models, and also applying the methods to Canadian and U.K. datasets. They find that the claimed superiority of the RF model for forecasting U.S. inflation carry over to U.K. inflation, but much less so to Canadian data. Furthermore, from 2020 onward the RF model performs poorly on U.S. data. Similarly, Huang et al. (2024) use a large dataset of macroeconomic and financial predictors to forecast Chinese CPI inflation and producer price index (PPI) inflation. They find that penalized linear regression models outperform benchmarks. Maehashi and Shintani (2020) analyze, among other things, Japanese CPI and wholesale price index (WPI) inflation. They show that beyond the very short run, the RF model and boosted trees outperform the benchmark AR model. Das and Das (2024) conclude that for Indian CPI the RF model provides much more accurate forecasts than an ARIMA model, but only when the COVID-19 pandemic period is included. Using pre-pandemic data the ARIMA model turns out to be difficult to beat. Lenza et al. (2023) apply the quantile random forest (QRF) model to euro area data. QRF density forecasts appear competitive with those of a linear model and the ECB survey of professional forecasters. While inflation in many of the large economic blocks has been studied extensively, evidence for smaller countries –like the Netherlands– is limited. Kohlscheen (2022) employs the RF model to examine the drivers of CPI inflation in a panel of high-income countries, including the Netherlands. Results for the Netherlands separately are not provided however. Vedder and van de Winkel (2024), replicating Goulet Coulombe et al. (2022)

on Dutch data, demonstrate that many ML models can generate significantly better forecasts of CPI inflation than an AR model. The current paper aims to contribute to the literature by studying the usefulness for forecasting (sub-components of) HICP inflation in the Netherland, taking as a benchmark both simple time series models and the inflation projections of De Nederlandsche Bank (DNB), the central bank of the Netherlands.

In this paper we build on Medeiros et al. (2021) and use a large dataset and multiple ML models to forecast the year-on-year Dutch inflation rate. Our target variable is the harmonized index of consumer price (HICP), and three sub-components of the HICP, i.e.: services inflation, non-energy industrial goods (NEIG) inflation and core inflation (HICP excluding energy and food). We compare the forecast accuracy of various types of ML models for these four target variables by using 129 monthly series from January 1990 to January 2024. We compute 1- to 12-months ahead forecasts for the four target variables, using a 18 year rolling estimation window. We compute the first forecast in January 2010. We calculate the inflation rate as the year-on-year percentage change in the price index series, which is the measure typically used by policy makers. Following Goulet Coulombe et al. (2021), we generate forecasts of year-on-year inflation for $h$ periods ahead using two different approaches. The direct approach entails forecasting year-on-year percentage changes $h$ periods ahead directly (using non-seasonally adjusted data), using a single model. Alternatively, the so-called path average approach requires estimation of $h$ different models, each used to forecast month-on-month changes in the target index, $1, 2, \ldots, h$ periods ahead, respectively. These $h$ forecasts are subsequently cumulated. This method uses seasonally adjusted data. We follow Goulet Coulombe et al. (2022) and assess the marginal improvement in out-of-sample $R^2$ ($OOS - R^2$). We gauge the forecast accuracy relative to two benchmarks: a simple linear time series model and an institutional forecast, the inflation projections of De Nederlandsche Bank (DNB), the central bank of the Netherlands. Finally, we examine how our results change when we only include data from the period prior to the pandemic.

Our main results can be summarized as follows. **First**, it is possible to beat simple linear time series models in terms of root mean squared forecast errors (RMSFE). Gains range from negligible at short horizons to more than 40% at longer horizons, depending on the target variable. Cumulative sums of squared forecast errors (CSSFE) indicate that the gains increase towards the end of the sample. **Second**, The performance of each ML model depends on the inflation measure and the forecast horizon. For headline HICP inflation, ridge regression (shrinkage) provides the

3

most accurate forecasts across for most horizons. In contrast, factor models are often preferred for forecasting services inflation. The results for forecasting core inflation and NEIG inflation are more mixed. This conclusion generally holds true not only when considering the best-performing model for each horizon, but also when extending the analysis to the top five models selected for each horizon and target variable. **Third**, for headline HICP inflation, the lowest RMSFE is obtained using path average forecasts, whereas direct forecast are often more accurate for the other series. However, when comparing the $OOS - R^2$ of direct forecasts to those of path average forecasts, the differences are often not statistically significant. **Fourth**, even during the pre-pandemic period, it is possible to outperform simple linear time series models. For headline HICP inflation, the gains become somewhat larger, while for the other inflation series, the gains are slightly smaller. Furthermore, direct forecasts more often outperform path average forecasts, suggesting that in tranquil times, the additional flexibility offered by path average forecasts is less necessary. **Fifth**, some ML models provide more accurate inflation forecast than DNB's official inflation forecast, particularly for NEIG inflation. Hence, ML models can be an important addition to the central bank's toolkit.

**Relation to the literature**. Multiperiod-ahead times series forecasts can be computed in various ways. Marcellino et al. (2006) use a large dataset of U.S. macroeconomic time series and demonstrate that iterated AR($p$) forecasts tend to outperform direct AR($p$) forecasts, with the relative performance improving as the forecast horizon extends. Quaedvlieg (2021) reaches similar conclusions using tests for multi-horizon superior predictive ability. Kock and Teräsvirta (2016) examine multiperiod-ahead forecasting using nonlinear neural network prediction methods. They find that iterated and direct forecasts often exhibit similar performance, with their ranking depending on the dataset. Goulet Coulombe et al. (2021) provide a new perspective on this literature. They systemically compare direct forecasts to path average forecasts for several target variables. In a linear context, the two approaches should yield fairly similar results. However, for nonlinear models, path average forecasts offer greater flexibility, potentially leading to more accurate forecasts. They demonstrate that the preference for direct versus path average forecasts is often variable specific, although the latter is generally preferred for variables that strongly co-move with the business cycle. They find that for U.S. CPI inflation, direct forecasts are usually more accurate than path average forecasts. Conversely, we find the opposite for Dutch headline HICP inflation.

The surge in inflation following the pandemic posses a challenge to inflation forecasters. Bobe-

ica and Hartwig (2023) propose allowing the residuals of vector autoregression models of inflation to have a fat-tailed distribution, while Lenza and Primiceri (2022) advocate explicitly modeling the change in residual volatility in order to capture the extraordinary large innovations during the pandemic era. On the other hand, data must be sufficiently informative to discriminate among competing forecasting models. In this context, Hansen et al. (2011) have shown that their so-called model confidence set contains considerably fewer inflation forecasting models when computed over 1970–1983 compared to 1984–1996. The great inflation of the 1970s, related to the oil price shocks and the subsequent rapid disinflation of the early 1980s, proved to be very helpful in selection competitive forecasting models. We demonstrate that the benefits from using ML models for forecasting inflation depend on whether the pandemic and post-pandemic periods are included or excluded from the evaluation sample.

A better understanding of inflation dynamics at the component level is of relevance for policy makers. Barkan et al. (2023) introduce a hierarchical neural network model for predicting U.S. CPI components, while Joseph et al. (2024) forecast CPI inflation in the United Kingdom using a large set of monthly CPI items. We demonstrate that the dynamics of headline HICP inflation can differ significantly from those of services or NEIG inflation, necessitating different forecasting models. For services inflation, factor models clearly outperform other model types in both the pre-pandemic and full data samples. However, for headline HICP and NEIG inflation, the best-performing models exhibit more variation.

In the academic literature, simple time series models are often used as benchmarks. A more pertinent issue for policymakers and other users of inflation forecasts is whether newly developed models can outperform existing institutional forecasts. In other words, it is important to determine whether adding one or model ML models to the existing suite of models is beneficial. Lenza et al. (2023) show that QRF forecasts of euro area inflation perform on par with the published Eurosystem inflation projections. Yoon (2021) demonstrates that boosting and RF models can produce forecasts of Japanese GDP growth that are more accurate than those made by the IMF and the Bank of Japan. Araujo and Gaglianone (2023) show that ML models can, in numerous cases, outperform traditional econometric models for Brazilian inflation. We find that it is not easy to outperform DNB's official inflation forecast for the Netherlands. However, for NEIG inflation, shrinkage models do produce more accurate forecasts up to 12 months ahead.

The rest of the paper is organized as follows. In Section 2 we describe the data, the forecast

design and the forecast evaluation metrics. In Section 3 we introduce the ML models. The main results are in Section 4. Concluding remarks are in Section 5. Details on the dataset, additional charts and a complete set of all results are shown in Appendix A, B, C, and the online Appendix (here), respectively.

# 2   Data, forecast design and forecast evaluation

## 2.1   Data

We have compiled a dateset consisting of monthly observations of 129 Dutch and international macroeconomic time series. All series were downloaded on March $4^{th}$, 2024, and cover the period from January 1990 to January 2024. The target variables in our analysis are four inflation series: headline HICP (HICP) inflation, services inflation (HICPS), non-energy industrial goods inflation (HICPNEIG), and core inflation (headline HICP excluding energy and food, HICPMEF).Our list of time series can be subdivided into eleven groups, following McCracken and Ng (2016) and Medeiros et al. (2021), namely: (1) output & income, (2) labor market, (3) consumption, (4) orders & inventories,(5) money & credit, (6) interest & exchange rates, (7) commodity prices, (8) producer prices, (9) domestic prices, (10) price expectations, and (11) stock market. In addition to the monthly series in our database, we include as potential predictors the four principal components computed from this set of variables. Further details are provided in Section 3.5. We consider four lags of all variables, as well as four autoregressive terms of the dependent variable. Hence, the analysis contemplates 532 potential predictors. Table A.1 in Appendix A presents more details on the data series used, seasonal adjustments, and transformations of all series. In this paper, we consider two methods for constructing forecasts of year-on-year inflation: the direct method and the path average method. Details on both methods are presented in Section 2.4. The datasets used to compute these forecasts differ in how trending predictors are transformed to achieve stationarity. For path average forecasts, trending time series are transformed to stationarity by taking (log) first differences. For direct forecasts, year-on-year changes of the predictors are used.

Figure 1(a) illustrates the development of headline HICP inflation and its three sub-components. After a period of inflation around 2 percent, headline HICP inflation increased during the pandemic. Following the first case of COVID-19 in the Netherlands in February 2020, inflation began

a steep ascent in the summer of 2021, peaking at 17.1% in September 2022 before starting to decline. The main driver behind this increase was the significant rise in energy and food prices. This is evident from Figure 1(b), which shows the development of core inflation. The increase in core inflation during the pandemic was much less steep, although still noticeable. NEIG inflation, shown in Figure 1(c), is much more volatile than core inflation. The smaller increase in core inflation during the pandemic appears to be due to the relatively stable path of services inflation (Figure 1(d). The pass-through of higher energy and food prices in the industrial sector is much stronger than in the services sector. This is also clearly seen by examining some distributional characteristics in Figures 1(a)– 1(d). NEIG inflation has the highest coefficient of variation (1.8), followed by headline HICP inflation (0.9), core inflation (0.7) and services inflation (0.5).
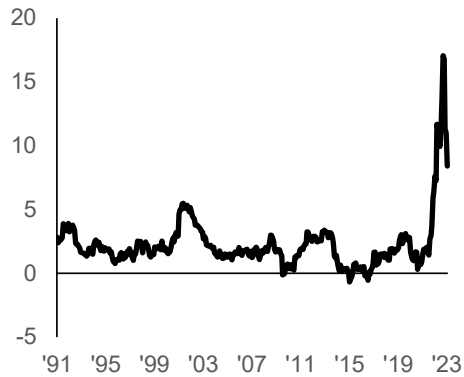
## 2.2 Forecast design

The forecasts are based on a rolling window with a fixed length of 216 months (18 years). This means that the number of forecasts depends on the forecast horizon. One advantage of a rolling-window framework is that it provides some protection against structural breaks or trends in the target variable. We re-calculate the seasonal inflation factors for each estimation window to avoid hindsight bias in these factors.

We employ a quasi real-time design, taking into account the data publication delays as of our download date (March $4^{th}$, 2024). However, we ignore the possibility of data revisions for the predictors, such as industrial production. The latter implies that we might overestimate the forecasting accuracy of the ML models. Unfortunately, a large real-time dataset for the Netherlands does not yet exist. Furthermore, Bernanke and Boivin (2003) have shown that the scope of the dataset appears to matter more for forecasting accuracy than the use of real-time (unrevised) data. Moreover, HICP inflation and its sub-components, our target variables, are not revised after the initial publication, which further mitigates the drawback of our quasi real-time design. Overall, it is very unlikely that the relative ranking of the ML models in terms of forecasting accuracy will change in a meaningful way when conducting a full-fledged real-time analysis. We evaluate the out-of-sample forecast performance of the ML models over the 1- to 12-month forecasting horizon ($h = 1, \ldots, 12$). Since the maximum horizon of the institutional forecast is only 10 months, the comparison between these forecasts and the ML models is limited to the 1- to 10-month horizon.
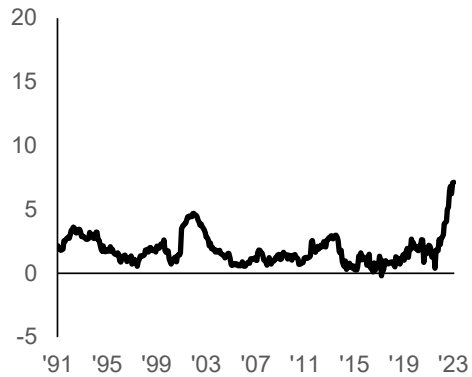
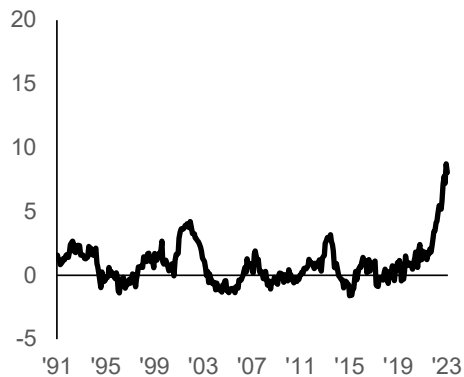## Figure 1: HICP inflation and its main sub-components

(a) HICP headline



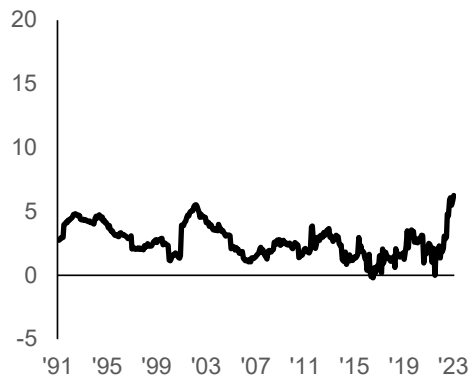| | |
|---|---|
| Minimum: | -0.7 |
| P25: | 1.4 |
| Median: | 1.8 |
| Mean: | 2.3 |
| P75: | 2.6 |
| Maximum: | 17.1 |
| Coeff. of variation: | 0.9 |

(b) HICP excluding food & energy



| | |
|---|---|
| Minimum: | -0.2 |
| P25: | 1.0 |
| Median: | 1.6 |
| Mean: | 1.8 |
| P75: | 2.4 |
| Maximum: | 7.1 |
| Coeff. of variation: | 0.6 |

(c) HICP non-energy industrial goods



| | |
|---|---|
| Minimum: | -1.6 |
| P25: | -0.2 |
| Median: | 0.6 |
| Mean: | 0.9 |
| P75: | 1.5 |
| Maximum: | 8.7 |
| Coeff. of variation: | 1.8 |

(d) HICP services



| | |
|---|---|
| Minimum: | -0.2 |
| P25: | 1.7 |
| Median: | 2.5 |
| Mean: | 2.7 |
| P75: | 3.5 |
| Maximum: | 6.3 |
| Coeff. of variation: | 0.5 |

## 2.3 Forecast evaluation

Following the literature, we measure the accuracy of the forecasts using the root mean squared forecast error (RMSFE). We compute the gain in RMSFE relative to benchmark models for each horizon separately.

The standard Diebold and Mariano (1995) (DM) test procedure is utilized to test the significance of the gains in relative predictive accuracy. We also implement the model confidence set (MCS), which selects the subset of best models at a given confidence level (Hansen et al., 2011).

To gain insight into the evolution of the forecast errors over time, we examine the forecast accuracy of the ML models over both the full sample and over the pre-pandemic sample. The full sample runs from January 2010 (2010M1) to December 2023 (2023M12), meaning that the first forecast is produced in January 2010 and the last (1-month ahead) forecast in December 2023. The pre-pandemic sample ends in January 2019, ensuring that even the 12-month ahead forecast refers to a date before the outbreak of the pandemic. Furthermore, we calculate the cumulated sum of squared forecast errors difference (CSSFED), which is defined as:

$$\text{CSSFED}_{M,h} = -\sum_{t=t_0}^{t_1}(e_{t,M,h}^2 - e_{t,BM,h}^2) \tag{1}$$

where $M$ represents a machine learning model and $BM$ represents a benchmark model. A CSS-FED *above* zero indicates that the forecasts of the machine learning model have a *lower* CSSFE up until that point in time, and are therefore more accurate than the benchmark model. Conversely, a CSSFED *below* zero indicates that the benchmark model has a lower forecast accuracy at that point in time. Additionally, a *decrease* in the CSSFED indicates that the model performance of the machine learning model is decreasing relative to the benchmark model, while an *increase* indicates the opposite.

Finally, we analyze the influence of several model features in the spirit of Goulet Coulombe et al. (2022). First, we compare the forecast accuracy of direct forecasts against path average forecasts. Second, we compare forecasts that include factors among the predictors with those that do not. Finally, we analyze the impact of using targeted predictors. For each model $M$, we

compute the pseudo-out-of-sample $R^2$ ($OOS - R^2$):

$$R_{t,h,M}^2 = 1 - \frac{\epsilon_{t,h,M}^2}{\frac{1}{T}\sum_{t=1}^{T}(y_{t+h} - \hat{y}_{t+h})} \tag{2}$$

where $\epsilon_{t,h,M}^2$ is the squared forecast error of model $M$ for horizon $h$ at time $t$. To assess the influence of feature $f$, we run Diebold-Mariano-style regressions for model $M_f$, which includes the feature, and model $M_{-f}$, which does not include the feature:

$$R_{t,h,M_f}^2 - R_{t,h,M_{-f}}^2 = \alpha_{M,f} + \nu_{t,h} \tag{3}$$

The main advantage of equation 3 is that the coefficient $\alpha_{M,f}$ can be interpreted as the marginal improvement in the $OOS - R^2$, and is not unit- or series-dependent. We then count the number of models for which the version with feature $f$ makes significantly more accurate forecasts than the version without feature $f$, and vice versa. We also calculate the median difference between the pseudo out-of-sample $R^2$ for models with and without feature $f$, separately for the various model types.

## 2.4   Direct and path average forecasts

Let $Y_t$ be (a sub-component of) the monthly HICP index, and let $x_t$ be a large macroeconomic dataset comprising of $N$ predictors, for $t = 1, \ldots, T$. $Y_t$ is not seasonally adjusted. Our target variable is $y_{t+h}$, the year-on-year percentage change in $Y_t$ $h$ periods into the future: $y_{t+h} = 100 * (Y_{t+h} - Y_{t+h-12})/Y_{t+h-12}$. Direct forecasts of $y_{t+h}$, denoted by $\hat{y}_{t+h}^{dir}$, can be obtained using the following prediction model:

$$y_{t+h}^{dir} = G_h(x_t) + \epsilon_{t+h} \tag{4}$$

where $G_h(\cdot)$ is a (potentially non-linear) mapping between the predictor variables $x_t$ and future inflation. Following Goulet Coulombe et al. (2021), another method to compute forecasts of year-on-year inflation $h$ periods into the future is by 'averaging' 1 to $h$ periods ahead forecasts of month-on-month changes in $Y_t$. Consider the following alternative prediction model:

$$y_{t+h}' = G_h(x_t') + \epsilon_{t+h}' \tag{5}$$

where $y'_{t+h} = \ln(Y'_{t+h}/Y'_{t+h-1})$ and $Y'_t$ is equal to $Y_t$ corrected for seasonal factors, $\zeta_t$, i.e. $Y_t = Y'_t + \zeta_t$. A path average forecast of year-on-year inflation, denoted by $\hat{y}^{pa}_{t+h}$, is then computed as:

$$\hat{Y}^{pa}_{t+h} = \exp(\ln(Y'_t) + \sum_{i=1}^{h} \hat{y}'_{t+i}) + \zeta_{t+h} \tag{6}$$

$$\hat{y}^{pa}_{t+h} = 100 * (\hat{Y}^{pa}_{t+h} - Y_{t+h-12})/Y_{t+h-12} \tag{7}$$

# 3  Models

We construct $h$-period ahead forecasts of the target variable $y_t$ using various models, distinguishing five 'types' of models: benchmark models, shrinkage models, tree models, ensemble models, and factor models. Model names are abbreviated according to the following convention: `[method].[model].[factor].[modelsel]`. `[method]` can take on two values, depending of the method used for constructing the forecasts. `[method]` equal to `Y` ('year-on-year') denotes direct forecasts. Path average forecasts are represented by `M`, as they are computed using forecasts of month-on-month changes in the price index. `factor` indicates whether factors are included in the set of predictors (`F`) or not (`NF`). For some models, tuning of the parameters is done in multiple ways. In those cases, `[modelsel]` can take on the values `BIC` (Bayesian Information Criterion) and `CV` (Cross Validation). If only a single method is used to tune the model, `[modelsel]` is left blank.

In the following, we briefly describe the models that we use in our analysis.

## 3.1  Benchmark model

The first benchmark is the random walk model. Direct forecast using the random walk mmodel (`Y.RW`) matches the no-change forecast in Atkeson and Ohanian (2001), i.e. past year's inflation is used as a forecast: $y_{t+h} = y_t$. In path average forecasts, on the other hand, the no-change forecast (`M.RW`) is past month's month-on-month inflation, similar to Stock and Watson (1999): $y'_{t+h} = y'_t$. The second benchmark extends the two random walk models with a 'drift' parameter: $y_{t+h} = \beta_0 + y_t$ (`Y.RWD`) and $y'_{t+h} = \beta_0 + y'_t$ (`M.RWD`), respectively. The third benchmark is the autoregressive model of order $p$, where $p$ is determined by the BIC, with a maximum of 4 lags. The AR($p$) is used both to produce direct forecasts $y_{t+h} = \beta_0 + \beta_1 y_t + \cdots \beta_p y_{t+p-1}$ (`Y.AR`) and

path average forecasts $y'_{t+h} = \beta_0 + \beta_1 y'_t + \cdots \beta_p y'_{t+p-1}$ (M.AR), respectively. The parameters are estimated by OLS. Considering all forecast horizons (12) and target variables (4), M.RWD has the lowest RMSFE in 23 out of the 48 cases, and is hence selected as the main benchmark model in the remainder of the paper.

## 3.2 Shrinkage models

We estimate several shrinkage estimators where $G_h(x_t) = \beta_h x_t$. All methods minimize the objective function:

$$\sum_{t=1}^{T} = \left\{ (y_{t+h} - \beta_h x_t)^2 + \lambda J(\beta_h) \right\} \tag{8}$$

where $\lambda$ is the hyperparameter determining the degree of regularization. The methods differ in terms of the specification of the penalty term $J(\beta_h)$. We choose $\lambda$ either by BIC or CV.

**LASSO** (LAS) The least absolute shrinkage and selection operator (LASSO) was introduced by Tibshirani (1996) and corresponds to the penalty term given by $J(\beta_h) = \sum_{i=1}^{N} |\beta_{h,i}|$. The penalty term of LASSO is the $L_1$ norm, which shrinks parameters of irrelevant predictors to zero. To achieve consistent model selection, Zou (2006) proposed adaLASSO (ALAS). adaLASSO is similar to LASSO but includes weighting parameters $\omega_i$ obtained from a first-step estimation. In this paper, we use LASSO for this purpose, with the penalty given by $J(\beta_h) = \sum_{i=1}^{N} \omega_i |\beta_{h,i}|$.

**Ridge Regression** (RR) was proposed by Hoerl and Kennard (1970) and assumes an $L_2$ norm penalty $J(\beta_h) = \sum_{i=1}^{N} \beta_{h,i}^2$. The parameters of less-relevant predictors can become very small, but unlike LASSO, will rarely be exactly zero.

**Elastic Net** (EN) was designed to make the most out of LASSO and Ridge Regression, and includes these models as special cases. The penalty term of Elastic Net is given by $J(\beta_h) = \omega \sum_{i=1}^{N} |\beta_{h,i}| + (1 - \omega) \sum_{i=1}^{N} \beta_{h,i}^2$, where $\omega$ is an additional tuning parameter setting the relative importance of the $L_1$ and $L_2$ penalty, respectively. We fix $\omega$ at 0.5.

## 3.3 Tree models

A regression tree with $M$ terminal nodes can be written as:

$$y_{t+h} = \sum_{i=1}^{M} \beta_{h,i} 1_{\{x_t \in R_i\}} + \epsilon_{t+h} \tag{9}$$

where $1_{\{\cdot\}}$ is an indicator function, $R_i$ is a partition of the space of $x_t$ and $\beta_{h,i}$ is the sample average of $y_{t+h}$ given $x_t \in R_i$. Estimation of (9) entails finding the best tree structure to minimize $\sum_{t=1}^{T} \epsilon_{t+h}^2$. A strength of regression trees is its capability to deal with non-linearity and interaction terms among predictors. A weakness, though, is that due to over-fitting the out-of-sample forecasting properties can be very poor. To deal with the issue of over-fitting, we consider three types of ensemble methods.

**Random Forest** The Random Forest (RF) model was introduced by Breiman (2001). The method is based on bootstrap aggregation (bagging) of randomly constructed regression trees. While the forecast of a regression tree in each bootstrap sample may suffer from over-fitting, averaging forecasts of bootstrap samples diminishes the variation and yields a more stable forecast. Ideally, regression trees of different bootstrap samples should not be highly correlated, since otherwise averaging may not be effective in lowering the variance of the forecast. In the Random Forest model, a dropout procedure is used to de-correlate the regression trees of the bootstrap samples. For estimation we use R-package `ranger`, with default settings and 500 trees.

**Boosted Tree** In Boosted Trees (BTREE) multiple regression trees are constructed similarly to bootstrap aggregation to overcome over-fitting. The algorithm starts with an initial regression tree:

$$f_0(x_t) = \sum_{i=1}^{M} \beta_{h,i} 1_{\{x_t \in R_i\}} \tag{10}$$

Then the model is updated in an iterative fashion from $k-1$ to $k$ according to the following rule:

$$f_k(x_t) = f_{k-1}(x_t) + \eta \sum_{i=1}^{M_k} \beta_{h,k,i} 1_{\{x_t \in R_{k,i}\}} \tag{11}$$

where $\eta$ is a learning rate (that we set at 0.05) and the regression tree in step $k$ is estimated from the residual from step $k-1$, $y_{t+h} - f_{k-1}(x_t)$. For estimation we use R-package `xgboost`, with the maximum depth of a tree equal to 4 and the maximum number of boosting iterations set to 1000.

**BART** The BART model is a sum-of-trees model introduced by Chipman et al. (2010). The Bayesian approach addresses the problem of over-fitting by using prior distributions to regularize the fit of each individual tree. Consequently, each tree explains only a small fraction of the variation in the target variable. To compute BART forecasts, we follow recommendations of Prüser (2019). We set the number of trees to 200, and $\alpha$ and $\beta$, which jointly control the depth of the trees,

to 0.1 and 1, respectively. The hyperparameters $k$ and $q$, which together determine the tightness of the prior of the values of the terminal nodes, are set at 2 and 0.9.

Borup et al. (2023) argue that "RF applied in high dimensions without an initial weeding out of irrelevant predictors may fail to reach its full potential". We follow their advise and add `TRF` and `TBART` to the list of tree models. In these models, we first select around 30 relevant targeting predictors, applying soft thresholding in combination with the LASSO regularizer, as suggested by Bai and Ng (2008).

## 3.4 Ensemble models

**Random subset regression** Unlike traditional regression methods that select a single best subset of predictors, the idea of Complete Subset Regression (Elliott et al. (2013)) is to generate a large number (K) of forecasts based on different subsets of $x_t$, denoted $x_t^i$. The final forecast $\hat{y}_{t+h}$ is then computed by taking a simple average of the individual forecasts $\hat{y}_{t+h}^i$:

$$
\begin{aligned}
y_{t+h}^i &= \beta_h x_t^i + \epsilon_{t+h}, \qquad i = 1 \ldots K & (12) \\
\hat{y}_{t+h} &= \sum_{i=1}^{M} \hat{y}_{t+h}^i / M & (13)
\end{aligned}
$$

In this paper, we average over 1000 Random Subset Regressions (`RSR`), cf. Boot and Nibbering (2019). In each regression, the number of predictors, in addition to the four lags of the target variable, is randomly selected, with a maximum of four. We also investigate the performance of a targeted version of Random Subset Regression (`TRSR`), following Bai and Ng (2008) and Kotchoni et al. (2019).

**Bagging** In bagging (`BAGG`), bootstrap samples of the original predictor variables and the target variable are repeatedly generated. We construct $K = 100$ samples using the block bootstrap method, with a fixed block length of five months. For each bootstrap sample, we select around 10 relevant predictors using soft thresholding. Then, a regression is applied to compute a forecast $y_{t+h}^i$. The final forecast from bagging is constructed as the simple average of the forecasts from the individual bootstrap samples.

## 3.5  Factor models

Following Stock and Watson (1999), we compute $k$ common factors $F_t^k$ as the first $k$ principle components of all predictor variables $x_t$. The $h$-period ahead forecast is constructed by running a principal components regression of the form:

$$y_{t+h} = \beta_h F_t^k + \epsilon_{t+h} \qquad (14)$$

where $\beta_h$ is a $(1 \times k)$ vector of coefficients. To select the number of factors, we consider the information criteria of Bai and Ng (2002). In the full sample, their $IC_1$, $IC_2$ and $IC_3$ criteria suggests that 2, 2, and 9 factors, respectively, are appropriate. We settle for four factors. In addition to this plain vanilla factor model (FACT), we also include a factor model with targeted predictors (TFACT) and a boosted factor model (BFACT). In the targeted factor model, we follow Bai and Ng (2008) and select around 30 relevant predictors using soft thresholding in combination with the LASSO regularizer. In the boosted factor, we adopt the boosting algorithm as in Bai and Ng (2009) to the select the factors and the number of lags in the model. The maximum number of factors and lags is set at 10 and 4, respectively.

# 4  Results

This section presents the results of the out-of-sample forecasting experiments. In Subsection 4.1, we graphically display the RMSFEs for all target variables, models, and forecast horizons. In Subsection 4.2, we identify the best-performing model for each inflation measure and forecast horizon, analyze which model type most frequently emerges as the best performer (even when considering the top-5 models), and map out which specification choices contribute to forecasting accuracy (following Goulet Coulombe et al. (2021)). Finally, Subsection 4.3 examines whether ML models can outperform DNB's official inflation forecast. Detailed results, including MCS $p$-values, are provided in the online Appendix (here).

## 4.1  ML models against simple benchmark model

Figures 2 and 3 display the RMSFE relative to M.RWD for all 61 models, over the full sample and the pre-pandemic sample, respectively. The dots represent the relative RMSFEs of the individual

models. To visualize the central tendency of the RMSFE distribution, the bold vertical line indicates the median across all models. The right and left boundaries of the boxplot represent the 75th and 25th percentiles, respectively, summarizing the spread of the distribution. The rightmost whisker is positioned at the smaller of the maximum RMSFE value and the 75th percentile plus 1.5 times the interquartile range (IQR). The leftmost whisker is positioned at the larger of the minimum RMSFE value and the 25th percentile minus 1.5 times the IQR.

The main message from Figure 2 is that, when analyzed over the full sample, it is quite challenging to outperform a simple benchmark model (`M.RWD`) at the one-month forecasting horizon. However, statistical models show substantial improvement over the benchmark for longer forecasting horizons. The improvements in forecast accuracy over the benchmark model vary widely between models, ranging from minimal to as much as 48%, depending on the target variable and the forecasting horizon. There is a steadily increasing gain in forecasting accuracy of the statistical models for forecasting services inflation.

Figure 3 shows that, in the pre-pandemic sample, the median forecast across all models is unable to beat the simple benchmark for most HICP measures, except for headline HICP and services inflation. This seems to be related to the stability of inflation during this period. Thus, choosing the right model is essential. Remarkably, services inflation is the only HICP index that shows the same pattern in the pre-pandemic sample as in the full sample: a marked increase in the relative forecast accuracy when the forecast horizon is larger.

The differences between the full sample and the pre-pandemic sample are particularly notable for headline HICP. While most models are outperformed by the benchmark model (`M.RWD`) in the full sample, in the pre-pandemic sample, the median gain is increases with the forecasting horizon. This is most likely related to the high volatility in food and energy prices in the post-pandemic period. The fluctuations in food and energy prices directly impacted HICP inflation, whilst the other inflation series were not or only indirectly affected via second round-effects by the high food en energy inflation. The machine learning models seem ill equipped to pick up these movements[1]. Finally it is interesting to note that the spread of the forecasts appears to have widened since the pandemic, as shown by the larger boxplots in Figure 2 compared to Figure 3. This suggests that the surge in inflation since the pandemic has helped to distinguish between more and less accurate
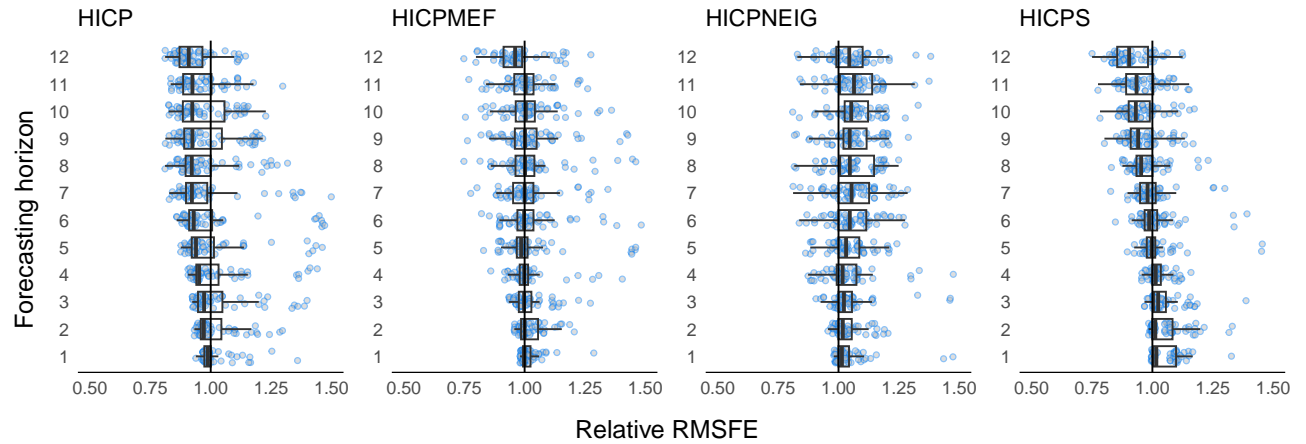
---

[1] This might also be caused by the selection of indicators in our database. There is a lack of long monthly time-series on energy components and food prices for the Netherlands.

Figure 2: Relative RMSFE, full sample (2010M1–2023M12)[†]

Figure 3: Relative RMSFE, pre-pandemic sample (2010M1–2019M12)[†]

forecasting models.

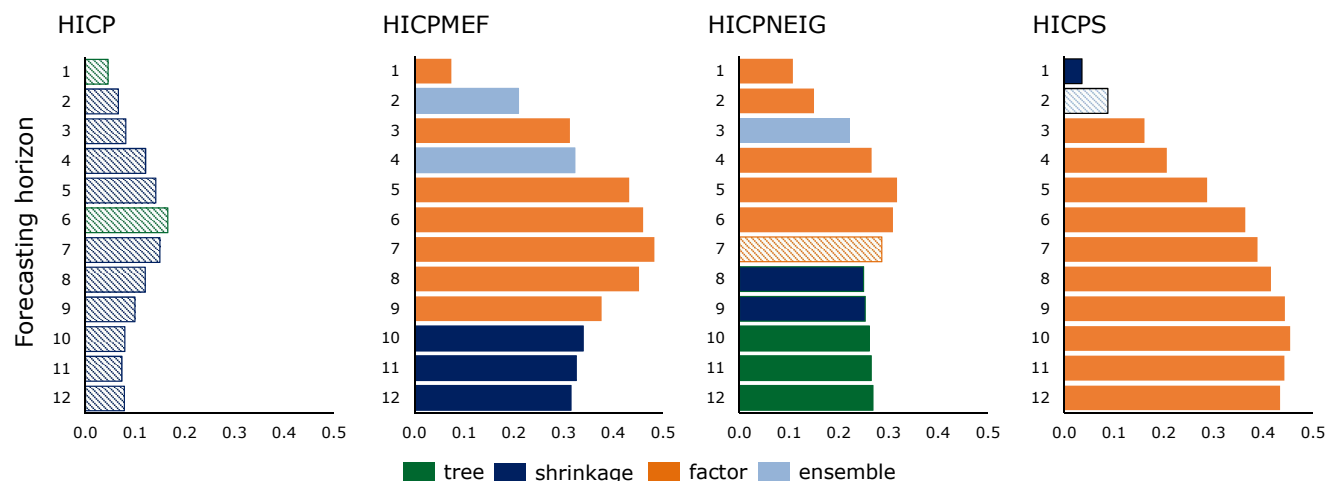## 4.2 The best ML model type for forecasting Dutch inflation

To gain more insight into the best-performing model we analyze two measures. Section 4.2.1 documents the RMSFE of the best-performing modeltype relative to the benchmark model. Section 4.2.2 examines the best-performing model sets. The idea here is that the difference between the best-performing modeltype and the runner up(s) can be quite small in practice.

### 4.2.1 The model with the lowest RMSFE

Figures 4 and 5 present the RMSFE and the model type of the best-performing model for each horizon, both for the full sample and the pre-pandemic sample. We distinguish four model types:shrinkage models (dark blue), tree models (green), ensemble models (light blue), and factor models (orange). The figures show the relative RMSFEs of the best-performing model against the benchmark M.RWD. A gain of 0.3 means the best model has an RMSFE that is 30% lower than the benchmark model. The figures also show the statistical significance of the gains in forecast accuracy. A solid bar indicates the gain with respect to the benchmark M.RWD model is significant at the 10% level. Non-significant entries are indicated by striped bars. We use the Diebold and Mariano test to test for statistical significance for the horizons 1 to 12 months.
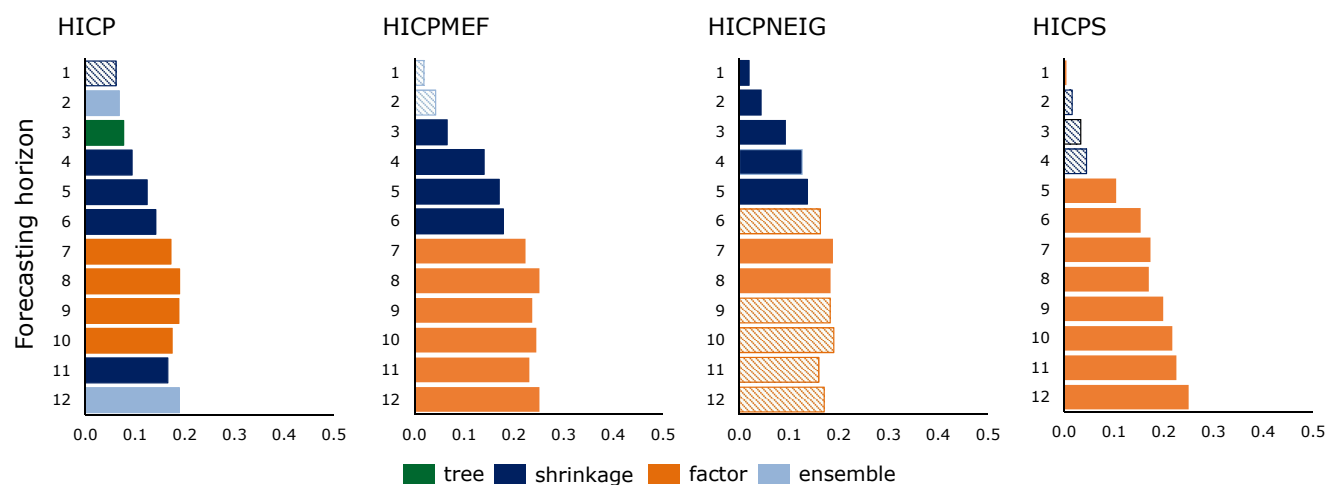
The main takeaways can be summarized as follows. First, measured over the full sample, none of the models are capable of systematically outperforming the headline HICP inflation forecasts of the benchmark model, as indicated by the striped bars for headline HICP in Figure 4. Although the average RMSFE of the best shrinkage model is approximately 10% lower on average, these differences are not significant. A possible explanation is that, since the volatility of HICP inflation is largely due to volatility in the energy and food components, our set of predictors is not sufficiently rich in terms of energy and food prices related series. Also, (announced) changes in taxes may play a role. Second, for core inflation, NEIG inflation, and services inflation machine learning models can significantly outperform the benchmark model for all horizons. Third, factor models are clearly overrepresented among the best-performing models when forecasting HICP sub-components over the full sample, although the dominance of the factor model is less clear for core inflation and NEIG inflation for forecasting over horizons longer than 7 months.

18

Figure 4: Gain in relative RMSFE best model, full sample (2010M1–2023M12)[†]

† [(RMSFE indicator model)/(RMSFE M.RWD)-1]×100. Solid bars indicate 10% significance level for Diebold and Mariano test. Striped bars indicate non-significance. HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.



Figure 5: Gain in relative RMSFE best model, pre-pandemic sample (2010M1–2019M12)[†]

† [(RMSFE indicator model)/(RMSFE M.RWD)-1]×100. Solid bars indicate 10% significance level for Diebold and Mariano test. Striped bars indicate non-significance. HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

For core inflation, the best shrinkage model beats the benchmark with the largest margin, and for NEIG inflation, the best shrinkage model and tree model beat the benchmark model. Fourth, Figure 5 shows that in the pre-pandemic sample, the best shrinkage model clearly dominates the benchmark model for forecasting over horizons up to 6 months, although not for all HICP sub-components. For longer forecasting horizons, factor models are the most accurate, with gains against the benchmark between 20% and 30%. Interestingly, the gain in forecasting accuracy on the 1- and 2-month forecasting horizons is either quite small or not statistically significant for HICP inflation and most sub-components. Fifth, the gains from using a machine learning model instead of a simple benchmark model are more pronounced for the full sample. This result indicates that the machine learning models are especially useful when policymakers most need them. Sixth, the gains in forecasting accuracy are largest for forecasting services inflation and increase with the forecast horizon. The gains are smallest for headline HICP, followed by NEIG inflation.

Figure 6 and 7 show the forecasting performance of all 62 models analyzed against the benchmark model to provide a deeper insight into which models within the model types perform best. The figures display heatmaps of the forecasting performance. The heatmap colors range from a gain in forecasting accuracy against the benchmark model ($RMSFE_{ML} - RMSFE_{M.RWD}$) of 5% (light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the benchmark model is *at least* 5% *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano test.

For the full sample (Figure 6), we observe that no models consistently outperform the benchmark model for all forecasting horizons and inflation measures. We do find evidence of pockets of strong performance. The plain-vanilla factor model (Y.FACT.F), the targeted factor model (Y.TFACT), and the boosted factor model (M.BFACT.F) are among the most consistently outperforming models for forecasting core inflation and services inflation. As concluded earlier, in the full sample, no model outperforms the benchmark model when forecasting headline HICP inflation. Within the class of shrinkage models ridge regressions (M.RR.NF.BIC, M.RR.F.CV and M.RR.F.BIC) clearly outperform the other shrinkage models. Within the ensemble models the best strategy is to use Bagging (M.BAGG.NF and M.BAGG.F). Although also tree based methods outperform the benchmark model for core inflation and services inflation, the gains of the best tree-based models against the benchmark are much smaller than the relative gains for other model types. Strikingly, the non-linear tree-based methods do not systematically perform better than the

linear shrinkage and factor models.

In the pre-pandemic sample (Figure 7), the gains of the machine learning models are much smaller than in the full sample, with the exception of headline HICP inflation. These outcomes are in line with the results in the previous section. Also in the pre-pandemic sample, we observe the most consistent outperformance for shrinkage and factor models. The shrinkage models have the most consistent performance across the different inflation measures. Again, the ridge regressions perform remarkably well (`M.RR.NF.BIC`, `M.RR.F.CV` and `M.RR.F.BIC`).

In summary, we find further evidence for our earlier claim in Section 4.1 that the forecasting power for headline HICP inflation is obscured by the very volatile –and hard to forecast– food and energy inflation. When comparing the forecasting accuracy of machine learning models, it is important to look at different HICP sub-components, sample periods, and individual forecasting horizons to unlock the full forecasting potential of these models. We find that shrinkage models have the most consistent performance across inflatin measures, both before and after the COVID-19 crisis. Overall, ridge regression is the model that most consistently outperforms the benchmark model. Strikingly, in the very volatile pandemic period the non-linear tree-based methods do *not* systematically perform better than the linear shrinkage and factor models.

### 4.2.2 A broader view: the best model sets for forecasting inflation

Instead of focusing only on the best-performing model, we now expand the scope to include the five best-performing models. We define these top-5 models as those with the highest *p*-values in the model confidence set (MCS) as introduced by Hansen et al. (2011). The intuition is that these five models are almost as good at forecasting inflation as the best-performing model in Section 4.2.1 and could also be considered when choosing a model type. However, in this study, the total number of models happens to be unevenly distributed across model types, increasing the chance that certain model types are selected as in set op top-5 models. Specifically, 32% of all models are shrinkage models, 29% are tree-based models, 19% are ensemble models, 10% are factor models, and 10% are benchmark models.
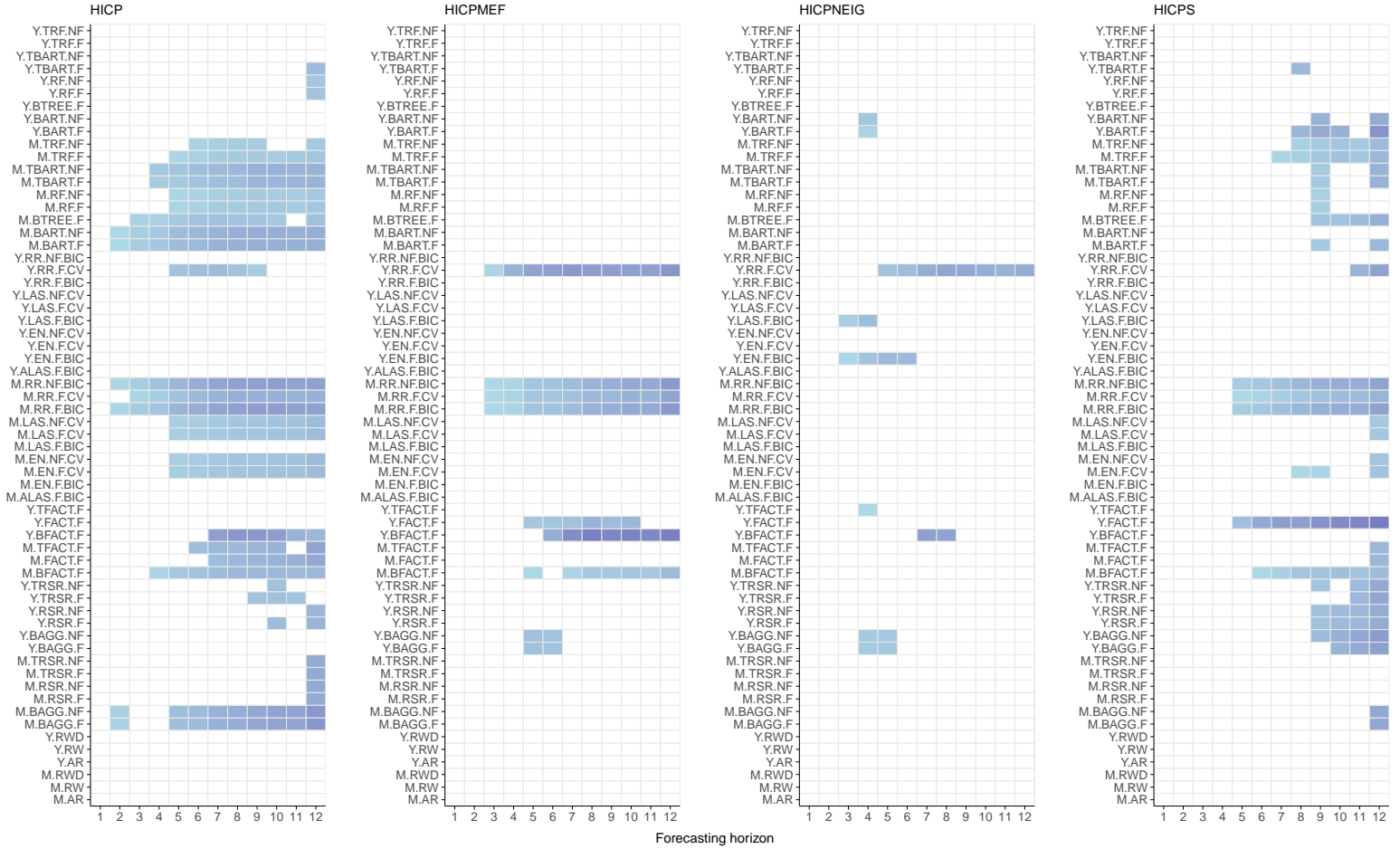
Figures 8 and 9 present the outcomes of the MCS in a specific form, showing the over/under representation of the different model types in the top-5. Since the total number of models is unevenly distributed across model types, we determine the over/under representation by dividing the number of models of a specific type in the top-5 by the share of this model type in the total

21

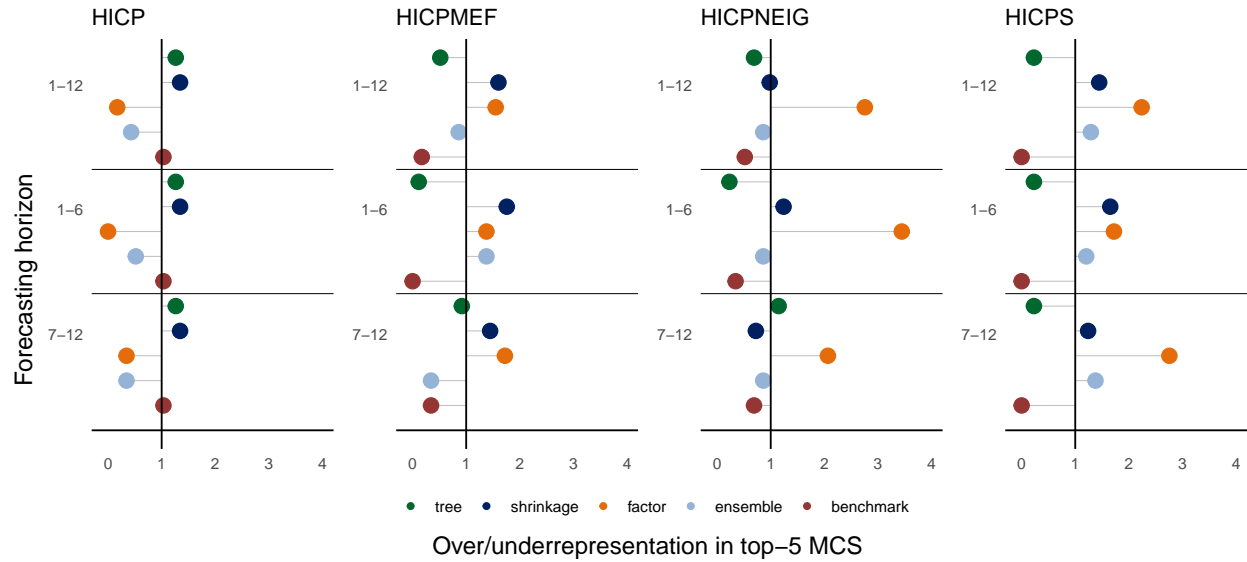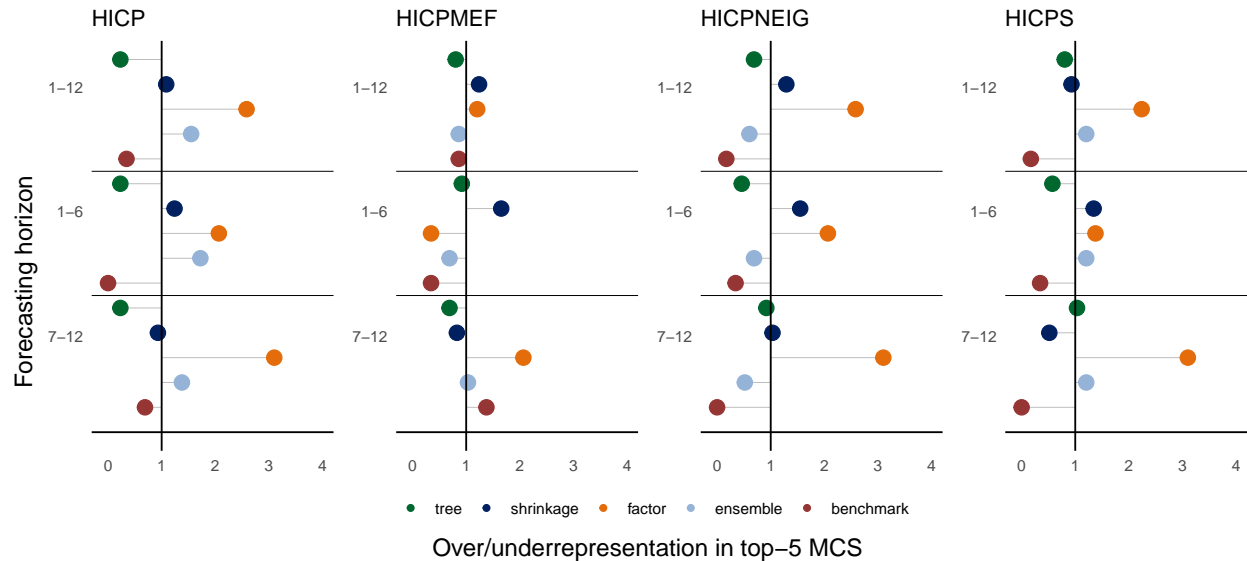Figure 6: Heatmap relative RMSFE, full sample (2010M1-2023M1) [†]

† Heatmap colors range from a gain in forecasting accuracy against the benchmark model of 5%(light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the benchmark model is *more* than 5% lower *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano-test.

Figure 7: Heatmap relative RMSFE, pre-pandemic sample (2010M1-2019M12) [†]

[†] Heatmap colors range from a gain in forecasting accuracy against the benchmark model of 5%(light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the benchmark model is *more* than 5% lower *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano-test.

Figure 8: Over & underrepresentation in top-5 MCS against random draw, full sample (2010M1–2023M12)[†]

Figure 9: Over & underrepresentation in top-5 MCS against random draw, pre-pandemic sample (2010M1–2019M12)[†]

number of models. If our measure equals 1, it means the number of models of a certain model type in the top-5 is in line with would result from a random selection from all models. If our measure is larger than 1, the model type is overrepresented in the top-5, indicating relatively high forecasting accuracy for this type of model. If our measure is below 1, it is underrepresented in the top-5. For conciseness, we present the outcomes averaged over the 1–12, 1–6, and 7–12 month horizons. Section B in the Appendix provides the measure for all individual forecasting horizons.

The figures show some interesting results. Averaged over the 1- to 12-month horizon, the figures reveal that the factor model type is overrepresented –often with a large margin– for all HICP sub-components and both samples, with the exception of headline HICP inflation over the full sample, where it is underrepresented. We already demonstrated that for headline HICP inflation in the full sample even the best model is indistinguishable from the benchmark model (see Figure 4). In line with this result, Figures 8 shows that the benchmark model is neither over- nor underrepresented for this inflation measure and sample. Second, averaged over the 1- to 12-month horizon, shrinkage methods are the clear runner-up to the factor models. Although the overrepresentation is usually much smaller than for the factor models, their representation is in line with or slightly above or below the equal representation line (the vertical line at an x-axis value of 1). This indicates that shrinkage models are a tough competitor for factor models in terms of forecasting accuracy. The number of ensemble models among the top-5 models is usually not much different from taking a random selection from all models. This indicates that the forecasting quality of the ensemble models does not really stand out. Strikingly, tree-type models are the worst model type using our measure. Tree models are *under*represented for all HICP measures and both samples.

### 4.2.3   A closer look at model specification

The previous two subsections presented our findings on the forecasting performance of model *types*. This section examines our outcomes through the lens of differences in model *specification*. As described in Section 3, models can differ across several dimensions. In this section, we look at three specification choices: (1) the transformation of the target variable, where we compare the relative performance of path average forecasts and direct forecasts, (2) the inclusion or exclusion of statistical factors in the set of predictors, and (3) the use of all predictors versus using only targeted predictors. The list of modeling choices could be extended further, but some of them are specific to the model type. For example, for some models, the tuning of the parameters is done in

several ways, such as via the BIC (Bayesian Information Criterion) or CV (Cross Validation). We do not consider these differences here.
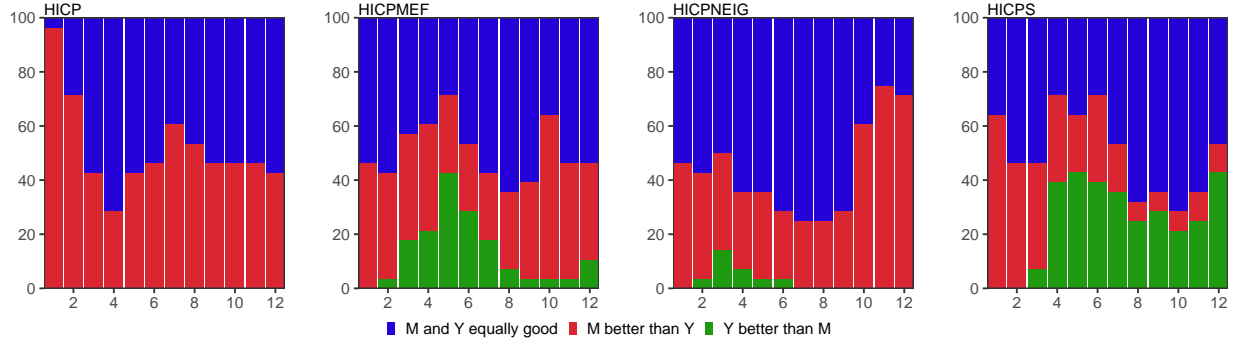
We run regressions inspired by Goulet Coulombe et al. (2021) and Goulet Coulombe et al. (2022); see equation (3) in Section 2. Figure 10 shows the outcome of our comparison of path average versus direct forecasts for the full sample. The blue bars show the percentage of models for which the pseudo-out-of-sample $R^2$ of the path average forecasts (M) is *not* significantly different from the direct forecasts (Y). On the other hand, the red bars show the percentage of models for which the path average forecasts are significantly better than the direct forecasts at the 10% level. Lastly, the green bars refer to cases for which the direct forecasts are significantly better than the path average forecasts, again at the 10% level. Figure 11 tests whether it is beneficial to add statistical factors to the set of predictors. Finally, Figure 12 examines whether using targeted predictors improves forecasting accuracy or not[2].

The main message from Figure 10 is that it is a 'no-regret' policy to forecast headline HICP inflation using path average forecasts. The blue (no difference between path average and direct forecast) and red bars (path average forecast better than direct forecast) sum to 100%, indicating that direct forecasts are never better. The same holds true for almost all forecasting horizons when forecasting NEIG inflation. The results are more ambiguous when forecasting core inflation and services inflation. The advantage of path average forecasts is their ability to model short-term volatility, leading to more stable long-term forecasts. On the other hand, direct forecasts may avoid cumulating errors, which are inherent in path average forecasts. The lower volatility of core inflation and services inflation lessen the need for using path average forecast as 'insurance' against short-term volatility. For the full sample, both the case for including factors in the set of predictors (Figure 11) and for using targeted predictors (Figure 12) is strong. For almost all horizons and measures of inflation, including factors doesn't hurt and often improves forecasting accuracy. The same conclusion holds for the use of targeted predictors.

To gain more insight into the quantitative differences between direct and path average forecasts, Figure 13 shows the gain in pseudo-out-of-sample $R^2$ of path average forecasts over direct forecasts, broken down by model type. For headline HICP inflation and NEIG inflation, the gains are larger at longer horizons, particularly for shrinkage models. Ensemble models perform worse
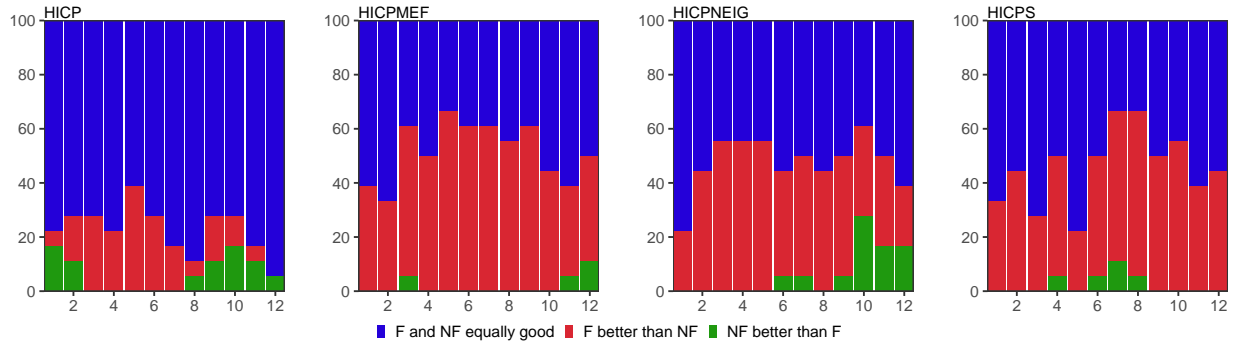
---

[2] Factor models are excluded from the comparison in Figure 11. Shrinkage models always use the complete set of predictors and are excluded from Figure 12.

Figure 10: Test of direct versus path average forecast, full sample (2010M1-2023M12)[†]



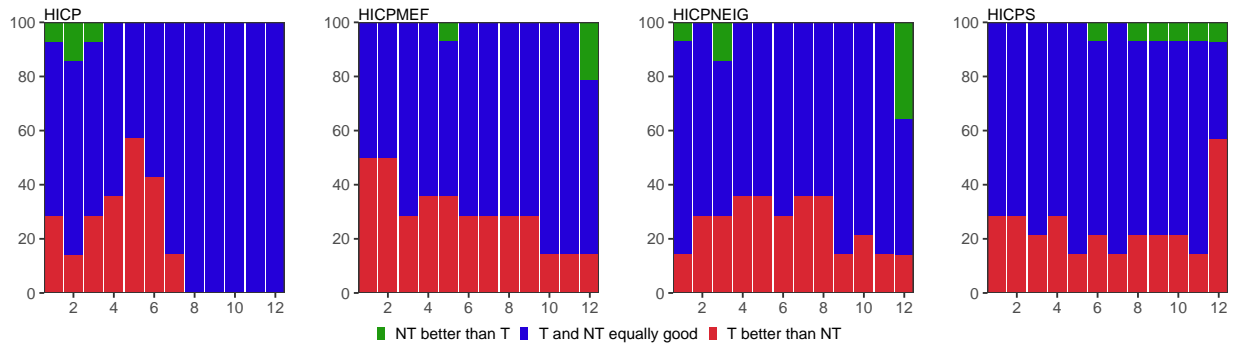M and Y equally good   M better than Y   Y better than M

[†] M: path average forecast, Y: direct forecast, HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure 11: Test of factor versus no factor forecast, full sample (2010M1-2023M12)[†]



F and NF equally good   F better than NF   NF better than F

[†] F: factor augmented forecast (excluding factor models), NF: forecast without factors (excluding factor models), HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure 12: Test of targeting versus no targeting of predictors, full sample (2010M1-2023M12)[‡]



NT better than T   T and NT equally good   T better than NT

[†] T: targeted predictors (excluding shrinkage models), NT: predictors not targeted (excluding shrinkage models), HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

when using path average forecasts. However, the weight given to the latter result should be low, because in the previous two sections, ensemble models hardly featured among the best-performing models and they were often underrepresented in the top-5 model sets. Depending on the forecasting horizon, factor models are sometimes better when using the direct forecasting approach.

The gains in the pseudo-out-of-sample $R^2$ from adding factors and targeted predictors are much smaller, especially from adding factors.

Section C in the Appendix presents the results for the pre-pandemic sample. Qualitatively, the results are broadly similar to those for the full sample. For example, the use of path average forecasts also tends to raise the pseudo-out-of-sample $R^2$. A key exception is that for factor models, it is clearly better to use *direct* forecasts for all inflation measures, apart from headline HICP inflation. This also holds for ensemble models, but given their weak forecasting performance it is inadvisable to use this model type anyway.

The main message from the pseudo-out-of-sample $R^2$ analysis is that, in general, path average forecasts are the preferred option for inflation forecasts in our dataset. Exceptions are factor models, where the picture is more mixed, especially in the pre-pandemic sample. The gains from following this strategy are sizable. Including factors and using targeted predictors is generally a 'no-regret' policy, although the gains are relatively small compared to the much larger gains from using the path average forecasts.
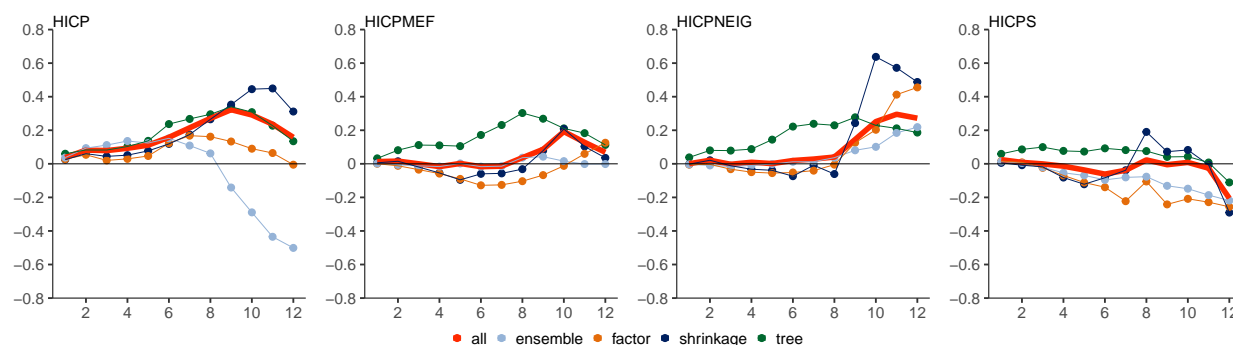
## 4.3   ML models versus institutional forecasts

In this section, we examine whether ML models can not only produce more accurate forecasts of Dutch inflation than simple benchmark models, but also outperform an established institutional forecast, DNB's official inflation forecast.

Since the introduction of the euro in 1999, DNB has been a member of the Eurosystem. The Eurosystem publishes macroeconomic projections four times a year[3]. Darracq Pariès et al. (2021) presents a recent assessment of the modeling toolbox currently in use within the Eurosystem, while Conrad and Enders (2024) investigate the accuracy of the inflation projections. DNB contributes to the quarterly projection exercises by providing forecasts of the Dutch annual inflation rate, ranging from 1- to 10-months into the future. These forecasts, known as the Narrow Inflation
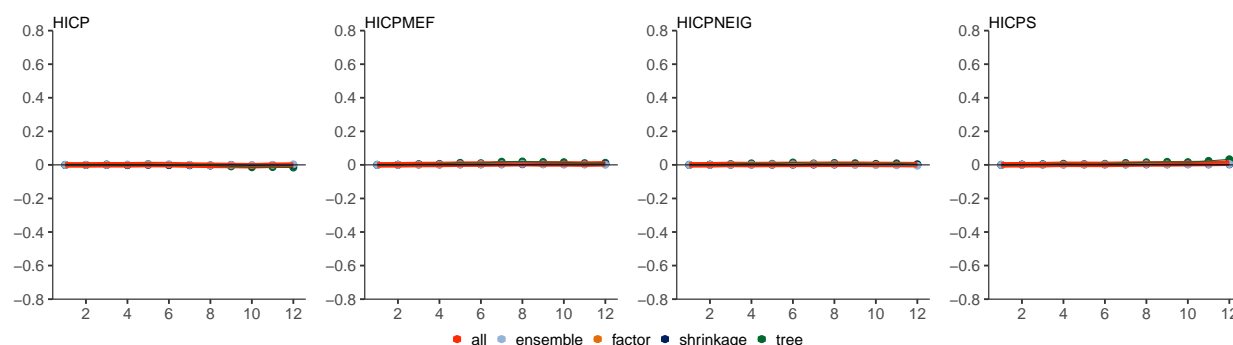
---

[3]   See https://www.ecb.europa.eu/press/projections/html/index.en.html and https://www.ecb.europa.eu/pub/pdf/other/staffprojectionsguide201607.en.pdf

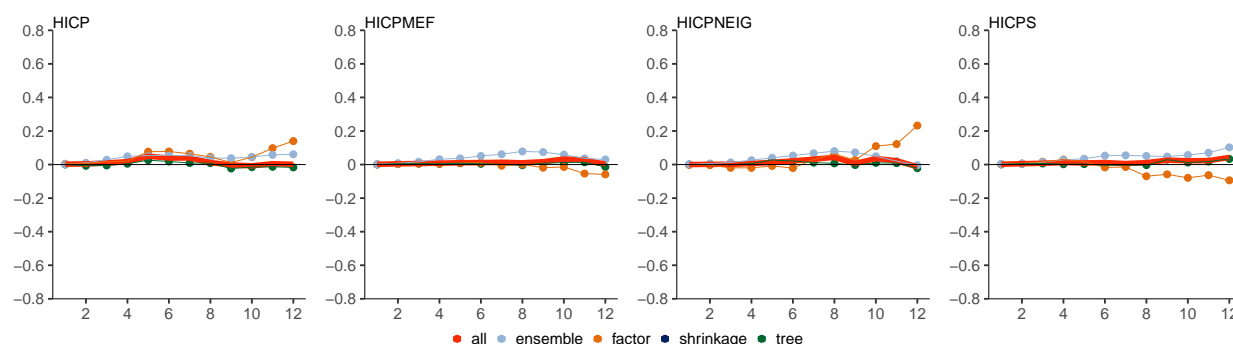Figure 13: Gain in OOS-R2 of path average forecast compared to direct forecast, full sample (2010M1-2023M12)†



† Gain/decrease in OOS-R2 of using path average forecasts, HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure 14: Gain in OOS-R2 of factor versus no factor forecast, full sample (2010M1-2023M12)†



† Gain/decrease in OOS-R2 of using factors (excluding factor models), HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure 15: Gain in OOS-R2 of targeting versus no targeting of indicators, full sample (2010M1-2023M12)†



† Gain/decrease using targeted factors or indicators (excluding shrinkage models), HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Projection Exercise (NIPE), will be referred to as NIPE throughout the remainder of the paper.
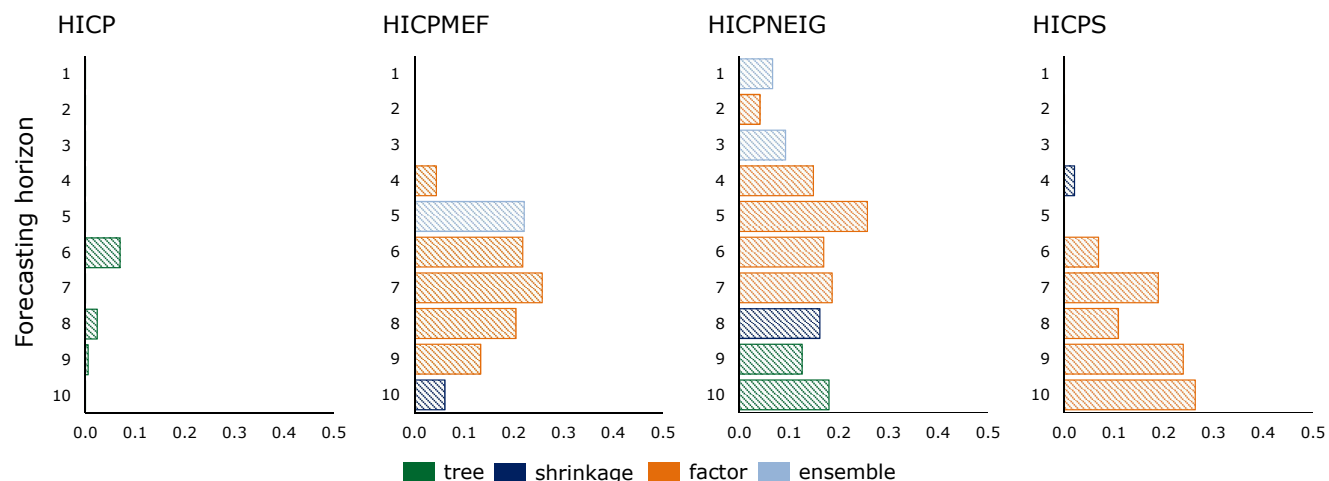
DNB's NIPE forecasts are based on a suite-of-models approach. This suite includes linear models for forecasting key components of HICP inflation, such as NEIG, food and services inflation, building on Den Reijer and Vlaar (2006). Additionally, SARIMA models are used to separately forecast more than 200 individual (COICOP) price sub-components of the HICP. Informal model averaging is employed to arrive at the final forecasts. The forecasts also consider announced government measures, such as changes in the VAT rate or energy taxes. Finally, the model-based forecasts are sometimes adjusted using expert-judgment to reflect relevant off-model information.

DNB makes inflation forecasts in four specific months of the year: February, May, August, and November. This contrasts with the ML model forecasts analyzed in Section 4, which are made every month. To fairly compare the forecast quality between NIPE forecasts and the ML forecasts, we will only use the forecasts from the four months when DNB forecasts are available. This means, for example, that we will only have 1-month ahead forecasts for March, June, September, and December, and 2-months ahead forecasts for April, July, October, and January, and so on.

Figures 16 and 17 present the RMSFE of the best model for each horizon. The models are again categorized as shrinkage models (dark blue), tree models (green), ensemble models (light blue), or factor models (orange) for the full sample (2010M1–2023M12) and the pre-pandemic sample (2010M1–2019M12). The figures show the relative RMSFEs of the best model compared to DNB's NIPE forecast. A gain of 0.3 means the best model has a RMSFE that is 30% lower than the NIPE forecast. The figures also indicate the statistical significance of the differences. A solid bar indicates the difference from the NIPE forecast is significant at the 10% level. We use the Diebold and Mariano test to assess the statistical significance for the horizons of 1 to 10 months. Non-significant entries are indicated by striped bars. If no bar appears, it means that none of the ML models can outperform the NIPE forecasts.
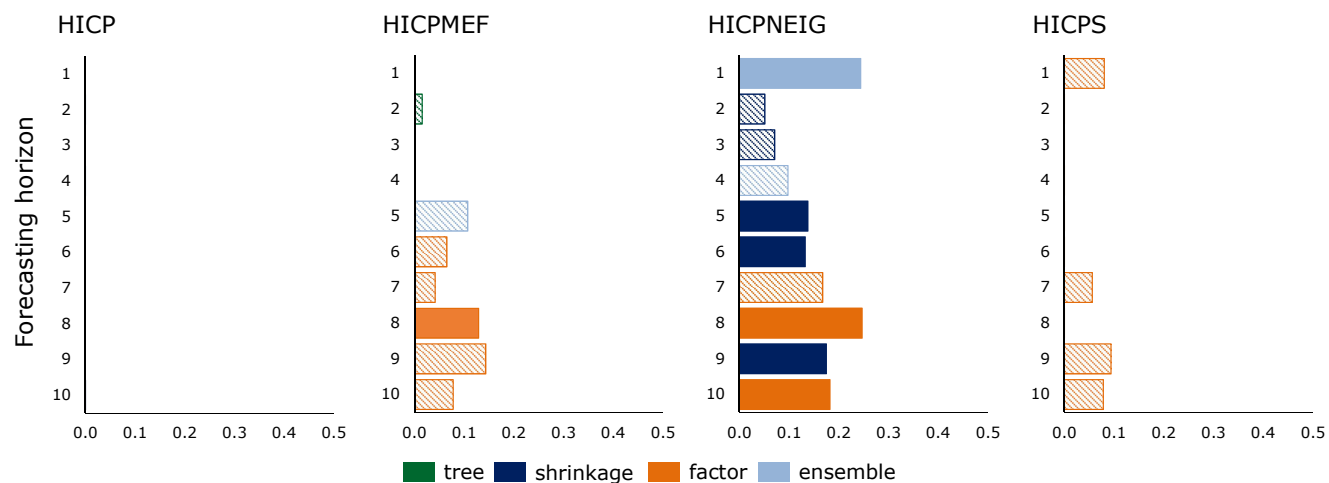
In the full sample (Figure 16), for some forecast horizons, the RMSFE of the best-performing ML model is substantially lower compared to the NIPE forecast, but the gain in forecast accuracy is in none of cases significant at the 10% level. For headline HICP inflation one possible explanation is that its fluctuations are largely due to changes in the energy and food components, and that our set of predictors may not be sufficiently rich in terms of energy and food price-related series. Additionally, (announced) changes in taxes may play a role for any of the inflation mea-

Figure 16: Gain in relative RMSFE best model, full sample (2010M1–2023M12)†



† [(RMSFE ML model)/(RMSFE NIPE)-1]×100. Solid bars indicate 10% significance level for Diebold and Mariano test. Striped bars indicate non-significance. HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure 17: Gain in relative RMSFE best model, pre-pandemic sample (2010M1–2019M12)†



† [(RMSFE ML model)/(RMSFE NIPE)-1]×100. Solid bars indicate 10% significance level for Diebold and Mariano test. Striped bars indicate non-significance. HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

sures. Furthermore, the result suggests that expert judgment, which is incorporated in the NIPE forecast but *not* in the ML models, *can* be crucial in times of volatile inflation. Finally, the lack of significance may be due to the high volatility of inflation during the pandemic, which makes the denominator in the Diebold and Mariano test very large.

Strikingly, the picture does not change much if we exclude the volatile pandemic period (Figure 17) and focus on the pre-pandemic sample. Even then, the ML models struggle to outperform the NIPE forecast. This might be because both the NIPE model and the ML models perform much better in less volatile times. Additionally, the NIPE model can incorporate pre-announced government measures (such as tax increases) and anticipated increases in rents wages (relevant for services inflation)and gas prices. The only exception is NEIG inflation. For this HICP sub-component, the best-performing ML models can significantly outperform the NIPE forecast for some horizons, especially for medium (5- and 6-months) and long horizons (8- to 12-months). Shrinkage models and factor models are most often among the best-performing models.

Figures 18 and 19 illustrate the forecasting performance of all 62 models compared to the NIPE forecasts, offering deeper insights into which models within each model type perform best. The figures display heatmaps of the forecasting performance. The heatmap colors range from a gain in forecasting accuracy against the benchmark model ($RMSFE_{ML} - RMSFE_{NIPE}$) of 5% (light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the NIPE is *at least* 5% *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano test.
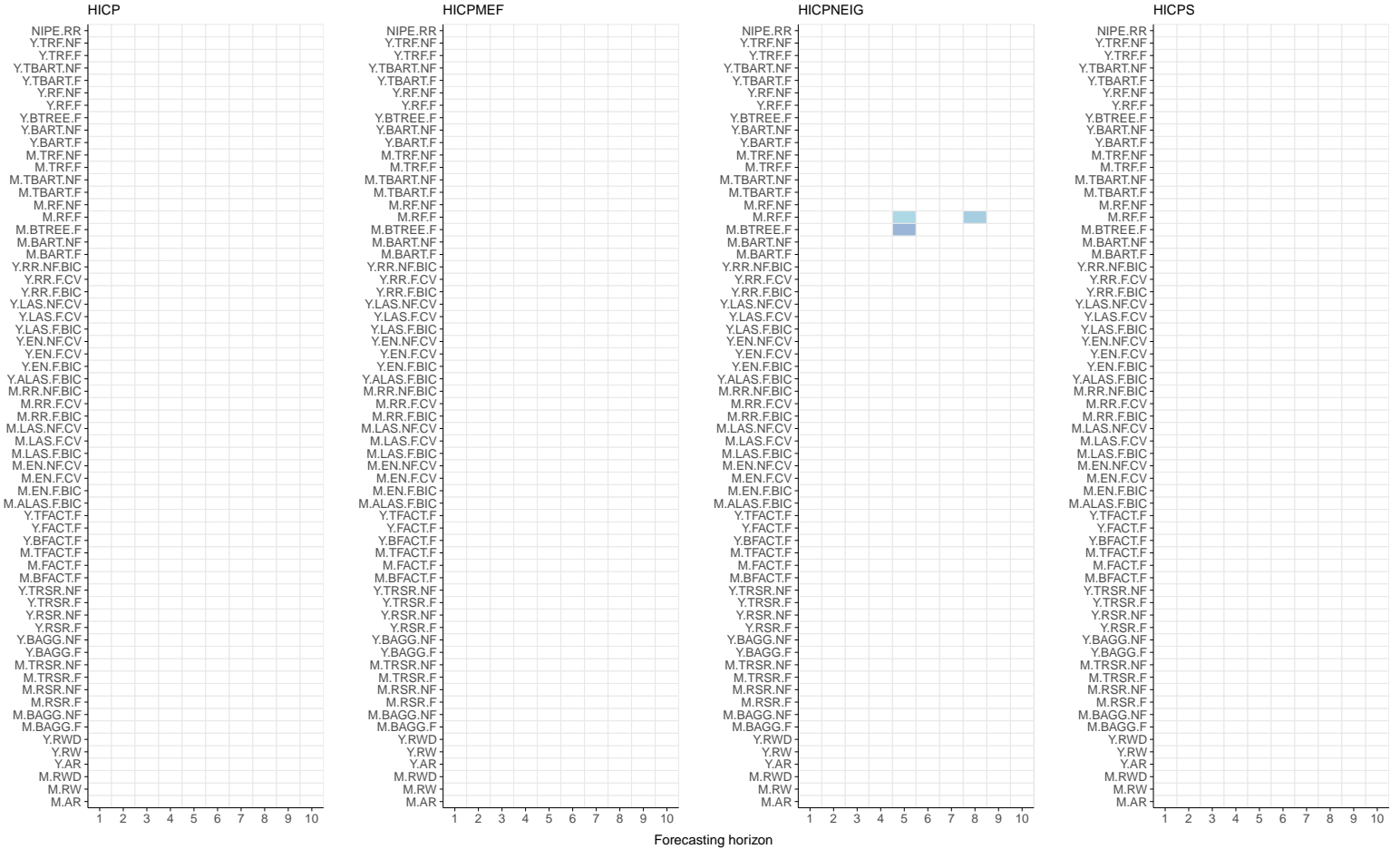
Unsurprisingly, Figure 18 shows an almost completely white canvas. Only two models outperform the NIPE forecasts for NEIG inflation, but the gains are small[4]. The same holds for the pre-pandemic sample. The only measure for which we observe more consistent gains for some models is NEIG inflation. Among all 62 models, the best performance is observed for the ridge regression models (`Y.RR.F.CV`, `M.RR.NF.BIC`, `M.RR.F.CV` and `M.RR.F.BIC`), as these models outperform the NIPE forecasts across multiple horizons[5].

To better understand the performance of the ridge regression models for NEIG inflation over time, Figure 20 illustrates the evolution of realized NEIG inflation (left-hand axis) and the evolu-

---

[4] These models were not the best-performing models. That is why the gains compared to the NIPE in Figure 16 are not statistically significant
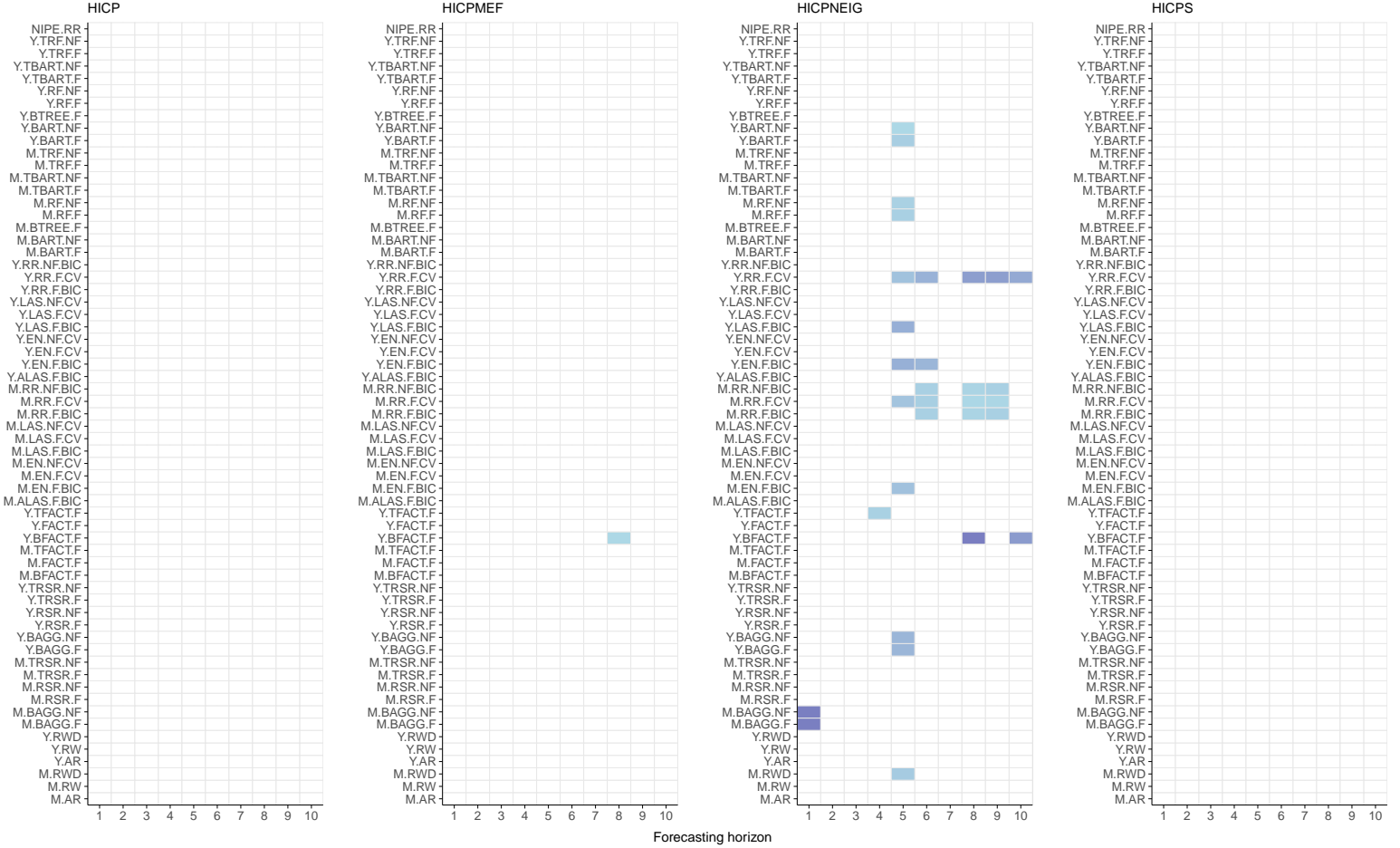
[5] We tested for average superior predictive ability, as described by Quaedvlieg (2021). The ridge models have - on average - statistically significant lower mean Square Forecast Error compared to the NIPE over short horizons (1- to 5-months), long horizons (6- to 10-months) and all horizons (1- to 10-months).

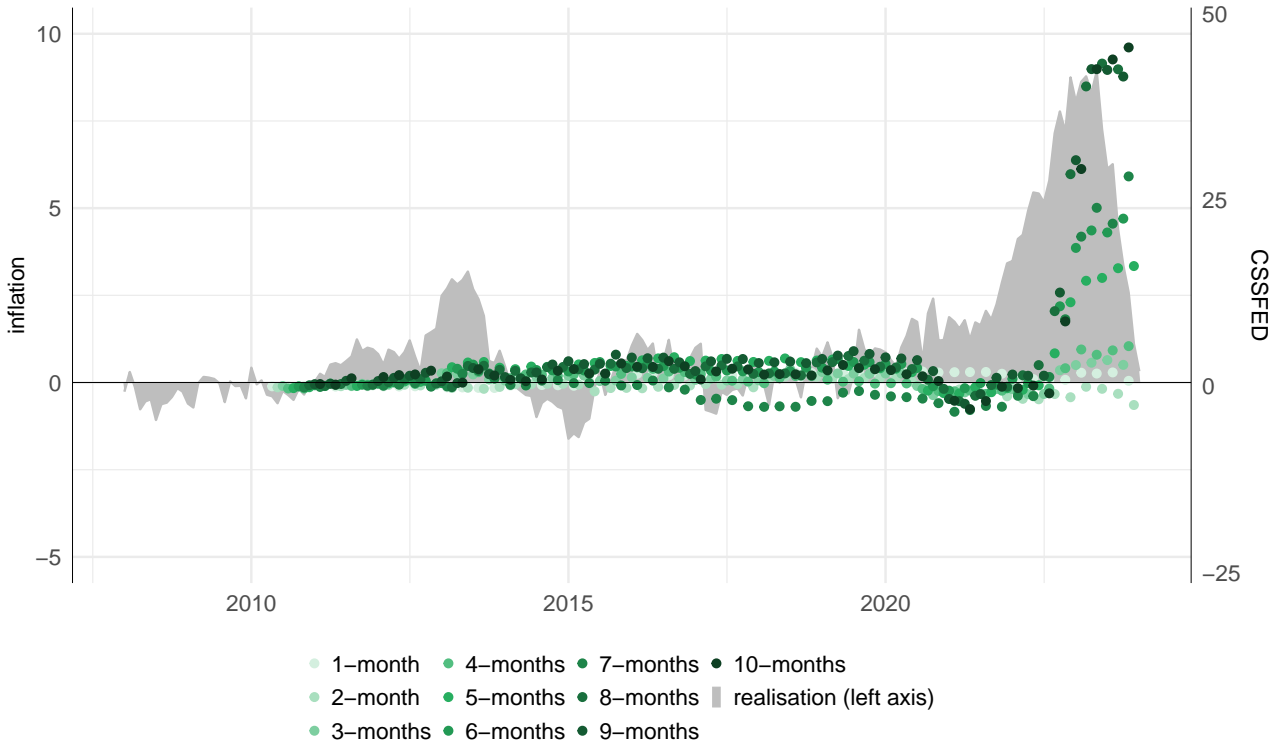Figure 18: Heatmap relative RMSFE, full sample (2010M1-2023M1) †

† Heatmap colors range from a gain in forecasting accuracy against the benchmark model of 5%(light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the benchmark model is *more* than 5% lower *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano-test.

Figure 19: Heatmap relative RMSFE, pre-pandemicsample (2010M1-2019M12) [†]

[†] Heatmap colors range from a gain in forecasting accuracy against the benchmark model of 5%(light blue) to 50% (dark blue). The relative gains are only shown if the forecast gain against the benchmark model is *more* than 5% lower *and* the gain is statistically significant at the 10% level according to the Diebold and Mariano-test.

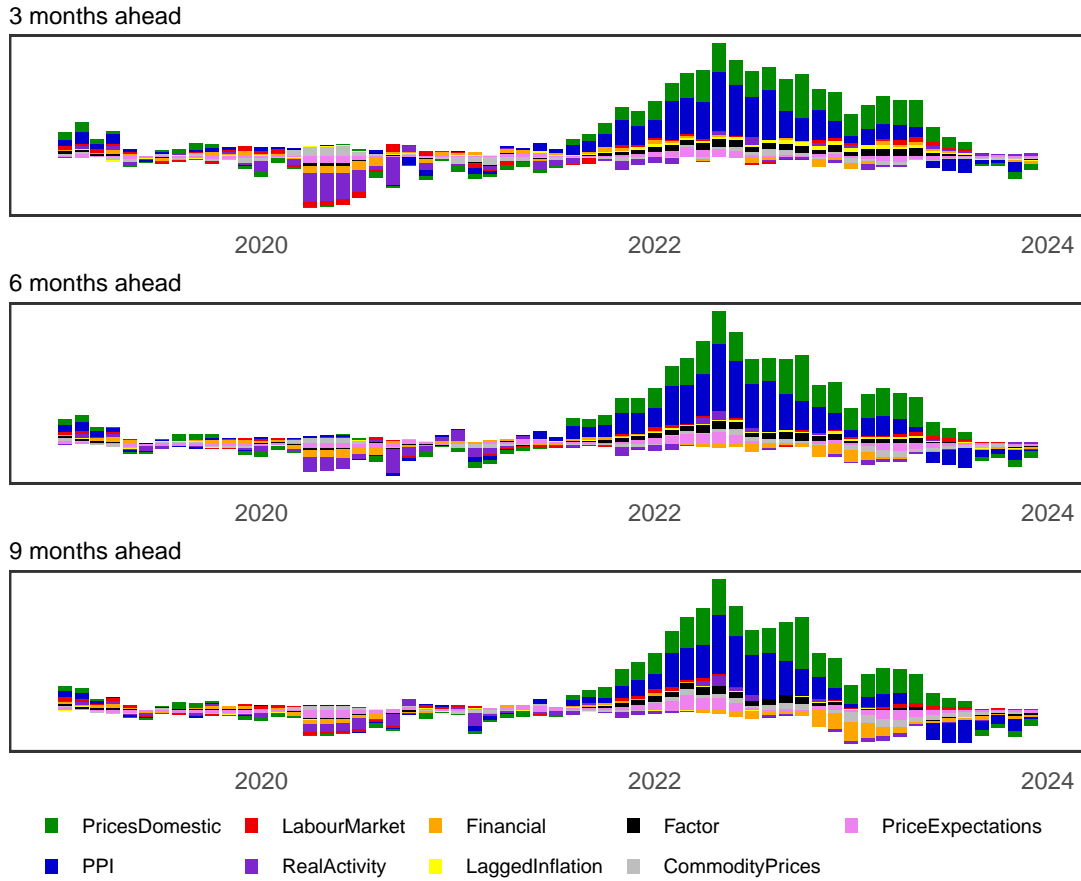Figure 20: NEIG inflation: CSSFED `M.RR.F.BIC` versus NIPE benchmark



tion of the CSSFED of `M.RR.F.BIC` compared to the NIPE model for forecasting horizons of 1 to 10 months ahead (right-hand axis).

Between 2010 and 2019, `M.RR.F.BIC` shows some marginal forecasting gains over the NIPE forecasts for most forecasting horizons. From 2020 onwards, as NEIG inflation edged up sligthly, there was an initial slight deterioration in the CSSFED. However, as NEIG inflation surged later on, `M.RR.F.BIC` clearly outperformed the NIPE forecasts, especially at longer forecast horizons.

Figure 21 highlights the predictors that most significantly contributed to the enhanced performance of `M.RR.F.BIC` in forecasting NEIG inflation during the pandemic and post-pandemic periods. The figure displays the contributions of various predictor sets. The contribution of each individual predictor is calculated as the product of its regression coefficient and its value (in deviation from the mean). These contributions are closely related to the concept of Shapley values, which are commonly used to quantify variable importance in ML models. In linear regression models, such as the ridge regression model, predictor contributions are equivalent to Shapley values if the predictors are orthogonal (which is not the case in our dataset).

During the pandemic, real activity predictors, expectations, and financial predictors were the

35

Figure 21: Decomposition of the `M.RR.F.BIC` forecast



3 months ahead

6 months ahead

9 months ahead

■ PricesDomestic  ■ LabourMarket  ■ Financial  ■ Factor  ■ PriceExpectations
■ PPI  ■ RealActivity  ■ LaggedInflation  ■ CommodityPrices

primary contributors to the `M.RR.F.BIC` forecasts. From the end of 2021 onwards, producer prices (PPI) and domestic price pressures, including headline HICP inflation and its main sub-components (excluding NEIG inflation), became the dominant driving forces. This shift reflects the pass-through of energy prices, ongoing supply chain disruptions, and some overheating of the Dutch economy.

# 5 Concluding remarks

In this paper, we investigate whether a range of machine learning methods can produce accurate short- and long-term inflation forecasts for the Netherlands. Our findings suggest that compared to forecasts from simple benchmark models the gains in terms of RMSFE range from negligible at short horizons to more than 40% at longer horizons, depending on the inflation measure. For headline HICP inflation, factor models and ridge regression models are often among the best-

performing models. Tree-based models and ensemble models are underrepresented among the best-performing models for almost all forecasting horizons and inflation measures. In most cases, path average forecasts are not inferior to direct forecasts, and using path average forecasts is a 'no-regret' policy. Adding factors and targeting the predictors can be considered a no-regret policy as well, although the gains in terms of forecasting accuracy are marginal.

Some machine learning models outperform DNB's official inflation forecast, particularly for NEIG inflation. Therefore, incorporating machine learning models into the central bank's model toolkit could enhance forecasting accuracy.

# References

Araujo, G. S. and W. P. Gaglianone (2023). Machine learning methods for inflation forecasting in Brazil: new contenders versus classical models. *Latin American Journal of Central Banking 4*(2), 1–29. link.

Atkeson, A. and L. E. Ohanian (2001). Are Phillips curves useful for forecasting inflation? Quarterly Review 2511, Federal Reserve Bank of Minneapolis. link.

Bai, J. and S. Ng (2002). Determining the number of factors in approximate factor models. *Econometrica 70*(1), 191–221. link.

Bai, J. and S. Ng (2008). Forecasting economic time series using targeted predictors. *Journal of Econometrics 146*(2), 304–317. link.

Bai, J. and S. Ng (2009). Boosting diffusion indices. *Journal of Applied Econometrics 24*, 607–629. link.

Barkan, O., J. Benchimol, I. Caspi, E. Cohen, A. Hammer, and N. Koenigstein (2023). Forecasting cpi inflation components with hierarchical recurrent neural networks. *International Journal of Forecasting 39*(3), 1145–1162. link.

Bernanke, B. S. and J. Boivin (2003). Monetary policy in a data-rich environment. *Journal of Monetary Economics 50*(3), 525–546. link.

Bobeica, E. and B. Hartwig (2023). The COVID-19 shock and challenges for inflation modelling. *International Journal of Forecasting 39*(1), 519–539. link.

Boot, T. and D. Nibbering (2019). Forecasting using random subspace methods. *Journal of Econometrics 209*(2), 391–406. link.

Borup, D., B. J. Christensen, N. S. Mühlbach, and M. Slot Nielsen (2023). Targeting predictors in random forest regression. *International Journal of Forecasting 39*(2), 841–868. link.

Breiman, L. (2001). Random forests. *Machine Learning 45*, 5–32. link.

Chipman, H. A., E. I. George, and R. E. McCulloch (2010). BART: Bayesian additive regression trees. *Annals of Applied Statistics 4*(1), 266–298. link.

Conrad, C. and Z. Enders (2024). The limits to the ECB's inflation projections. Policy Brief 945, SUERF. link.

Darracq Pariès, M., A. Notarpietro, J. Kilponen, N. Papadopoulou, S. Zimic, P. Aldama, G. Langenus, L. J. Alvarez, M. Lemoine, and E. Angelini (2021). Review of macroeconomic modelling in the eurosystem: current practices and scope for improvement. Occasional Paper Series 267, European Central Bank. link.

Das, P. K. and P. K. Das (2024). Forecasting and analyzing predictors of inflation rate: Using machine learning approach. *Journal of Quantitative Economics 22*, 493–517. link.

Den Reijer, A. and P. Vlaar (2006). Forecasting inflation: An art as well as a science! *De Economist 154*, 19–40. link.

Diebold, F. X. and R. S. Mariano (1995). Comparing predictive accuracy. *Journal of Business & Economic Statistics 13*(3), 134–144. link.

Elliott, G., A. Gargano, and A. Timmermann (2013). Complete subset regressions. *Journal of Econometrics 177*(2), 357–373. link.

Faust, J. and J. H. Wright (2013). Forecasting inflation. *Handbook of Forecasting 2*(A), 2–56. link.

Goulet Coulombe, P., M. Leroux, D. Stevanovic, and S. Surprenant (2021). Macroeconomic data transformations matter. *International Journal of Forecasting 37*(4), 1338–1354. link.

Goulet Coulombe, P., M. Leroux, D. Stevanovic, and S. Surprenant (2022). How is machine learning useful for macroeconomic forecasting? *Journal of Applied Econometrics 37*(5), 920–964. link.

Hansen, P. R., A. Lunde, and J. M. Nason (2011). The model confidence set. *Econometrica 69*(2), 453–497. link.

Hoerl, A. E. and R. W. Kennard (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics 12*(1), 55–67. here.

Huang, N., Y. Qi, and J. Xia (2024). China's inflation forecasting in a data-rich environment: based on machine learning algorithms. *Applied Economics forthcoming*, 1–26. link.

Joseph, A., G. Potjagailo, C. Chakraborty, and G. Kapetanios (2024). Forecasting UK inflation bottom up. *International Journal of Forecasting 40*(4), 1521–1538. link.

Kock, A. B. and T. Teräsvirta (2016). Forecasting macroeconomic variables using neural network models and three automated model selection techniques. *Econometric Reviews 35*(8), 1753–1779. link.

Kohlscheen, E. (2022). What does machine learning say about the drivers of inflation? Working Papers 980, BIS. link.

Kotchoni, R., M. Leroux, and D. Stevanovic (2019). Macroeconomic forecast accuracy in a data-rich environment. *Journal of Applied Econometrics 34*, 1050–1072. link.

Lenza, M., I. Moutachaker, and J. Paredes (2023). Density forecasts of inflation: a quantile regression forest approach. Working Paper 2830, European Central Bank. link.

Lenza, M. and G. E. Primiceri (2022). How to estimate a vector autoregression after March 2020. *Journal of Applied Econometrics 37*(4), 688–699. link.

Maehashi, K. and M. Shintani (2020). Macroeconomic forecasting using factor models and machine learning: an application to Japan. *Journal of The Japanese and International Economies 58*, 101104. link.

Marcellino, M., J. H. Stock, and M. W. Watson (2006). A comparison of direct and iterated multi-step AR methods for forecasting macroeconomic time series. *Journal of Econometrics 135*(1), 499–526. link.

McCracken, M. W. and S. Ng (2016). FRED-MD: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics 34*(4), 574–589. link.

Medeiros, M. C., G. F. R. Vasconcelos, A. Veiga, and E. Zilberman (2021). Forecasting inflation in a data-rich environment: The benefits of machine learning methods. *Journal of Business & Economic Statistics 39*(1), 98–119. link.

Naghi, A. A., E. O. E., and M. D. Zaharieva (2024). The benefits of forecasting inflation with machine learning: New evidence. *Journal of Applied Econometrics forthcoming*, 1321—1331. link.

Prüser, J. (2019). Forecasting with many predictors using Bayesian additive regression trees. *Journal of Forecasting 38*(1), 40–53. link.

Quaedvlieg, R. (2021). Multi-horizon forecast comparison. *Journal of Business & Economic Statistics 39*, 621–631. link.

Stock, J. H. and M. W. Watson (1999). Forecasting inflation. *Journal of Monetary Economics 44*, 293–335. link.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B 58*, 267–288. link.

Vedder, C. and E. van de Winkel (2024). Is machine learning beneficial for macroeconomic forecasting with limited observations? Discussion paper, CPB. link.

Yoon, J. (2021). Forecasting of real GDP growth using machine learning models: Gradient boosting and random forest approach. *Computational Economics 57*, 247–265. link.

Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association 101*(476), 1418–1429. link.

# A   Data set

Table A.1 provides the 129 monthly series that have been used for the estimation of the models in the main text. As mentioned in the main text, the data set can be split-up into eleven groups: 1) output & income: 27 series, 2) labor market: 10 series, 3) consumption: 14 series, 4) orders & inventories: 6 series, 5) money & credit: 9 series, 6) interest & exchange rates: 14 series, 7) commodity prices: 7 series, 8) producer prices: 20 series, 9) domestic prices: 13 series, 10) price expectations: 3 series, and 11) stock market: 6 series. A complete list of the series is provided in Table A.1.

We constructed two databases: one for computing path average forecasts and another for calculating direct forecasts. First, we collected all monthly series in the economics and finance sections from data warehouses (see column 'Source') that started in January 1990 or earlier, and hand-picked series within each of the eleven groups we defined. Next, for the path average database, we performed seasonal adjustment on all series that were not seasonally adjusted at the source using the US Census X13 ARIMA-SEATS method. Then, we transformed all variables to stationarity. For the path average forecast database we mostly applied the first difference or log first difference operator to the series, whereas for the direct forecast database we employed annual differences or annual log differences[6]. Details on the transformations applied in both databases are given in the 'Transf.' column in Table A.1. We use the monthly current and one year ahead consumer price inflation forecast for the Netherlands from Consensus Forecasts to calculate a monthly series of consumer price inflation expectations 12 months ahead.

Table  A.1:  Description monthly database

| Nr. | Description | Source | Transf. | Last |
|---|---|---|---|---|
| **1. Output & Income** | | | | |
| 1. | Manufacturing (total) | EUR | 2,5 | Dec-23 |
| 2. | Manufacturing of food products, beverages & tobacco | EUR | 2,5 | Dec-23 |
| 3. | Manufacturing of textiles | EUR | 2,5 | Dec-23 |
| 4. | Manufacturing of wearing apparel | EUR | 2,5 | Dec-23 |
| 5. | Manufacturing of leather & related products | EUR | 2,5 | Dec-23 |
| 6. | Manufacturing of wood & of products of wood & cork | EUR | 2,5 | Dec-23 |
| 7. | Manufacturing of paper & paper products | EUR | 2,5 | Dec-23 |
| 8. | Printing & reproduction of recorded media | EUR | 2,5 | Dec-23 |
| 9. | Manufacturing of coke & refined petroleum products | EUR | 2,5 | Dec-23 |

---

[6]There are two exceptions to this rule: the economic sentiment indicator and business climate indicator are included in levels as proxies for the output gap.

| Nr. | Description | Source | Transf. | Last |
|---|---|---|---|---|
| 10. | Manufacturing of chemicals & chemical products | EUR | 2,5 | Dec-23 |
| 11. | Manufacturing of rubber & plastic products | EUR | 2,5 | Dec-23 |
| 12. | Manufacturing of other non-metallic mineral products | EUR | 2,5 | Dec-23 |
| 13. | Manufacturing of basic metals | EUR | 2,5 | Dec-23 |
| 14. | Manufacturing of fabricated metal products | EUR | 2,5 | Dec-23 |
| 15. | Manufacturing of computer, electronic & optical products | EUR | 2,5 | Dec-23 |
| 16. | Manufacturing of machinery & equipment | EUR | 2,5 | Dec-23 |
| 17. | Manufacturing of motor vehicles & trailers | EUR | 2,5 | Dec-23 |
| 18. | Manufacturing of furniture | EUR | 2,5 | Dec-23 |
| 19. | Other manufacturing | EUR | 2,5 | Dec-23 |
| 20. | Supply of natural gas | CBS | 2,5 | Jan-24 |
| 21. | Use of natural gas | CBS | 2,5 | Jan-24 |
| 22. | Hotels & similar accommodation, arrivals, foreigners | EUR | 2,5 | Dec-23 |
| 23. | Hotels & similar accommodation, arrivals, total | EUR | 2,5 | Dec-23 |
| 24. | Hotels & similar accommodation, nights spend, foreigners | EUR | 2,5 | Dec-23 |
| 25. | Hotels & similar accommodation, nights spend, total | EUR | 2,5 | Dec-23 |
| 26. | Economic sentiment indicator | EUR | 0,0 | Feb-24 |
| 27. | Business climate indicator | EUR | 0,0 | Feb-24 |
| **2. Labor market** | | | | |
| 28. | Unemployment rate, total | ECB | 1,4 | Jan-24 |
| 29. | Unemployment rate, female | ECB | 1,4 | Jan-24 |
| 30. | Unemployment rate, total, < 25 years | ECB | 1,4 | Jan-24 |
| 31. | Unemployment rate, female, < 25 years | ECB | 1,4 | Jan-24 |
| 32. | Hourly wages | CBS | 2,5 | Oct-23 |
| 33. | Consumer confidence, unemployment > 12 months | EUR | 1,4 | Feb-24 |
| 34. | Construction confidence, employment > next 3 months | EUR | 1,4 | Feb-24 |
| 35. | Construction confidence, limiting factors, shortage of labour | EUR | 1,4 | Feb-24 |
| 36. | Industrial confidence, employment > 3 months | EUR | 1,4 | Feb-24 |
| 37. | Retail confidence, employment expectations > 3 months | EUR | 1,4 | Feb-24 |
| **3. Consumption** | | | | |
| 38. | Consumer confidence, headline | EUR | 1,4 | Feb-24 |
| 39. | Consumer confidence, financial situation < 12 months | EUR | 1,4 | Feb-24 |
| 40. | Consumer confidence, financial situation > 12 months | EUR | 1,4 | Feb-24 |
| 41. | Consumer confidence, general economic situation < 12 months | EUR | 1,4 | Feb-24 |
| 42. | Consumer confidence, general economic situation > 12 months | EUR | 1,4 | Feb-24 |
| 43. | Consumer confidence, major purchase > 12 months | EUR | 1,4 | Feb-24 |
| 44. | Consumer confidence, major purchases, current | EUR | 1,4 | Feb-24 |
| 45. | Consumer confidence, statement on financial situation of household | EUR | 1,4 | Feb-24 |
| 46. | Consumer confidence, savings > 12 months | EUR | 1,4 | Feb-24 |
| 47. | Consumer confidence, current ec. situation is adequate for savings | EUR | 1,4 | Feb-24 |
| 48. | Retail confidence, headline | EUR | 1,4 | Feb-24 |
| 49. | Retail confidence, volume of stocks currently held | EUR | 1,4 | Feb-24 |
| 50. | New passenger car | ECB | 2,5 | Jan-24 |
| 51. | Retail trade | ECB | 2,5 | Jan-24 |
| **4. Orders & Inventories** | | | | |
| 52. | Industrial confidence, assessment of the current level of stocks | EUR | 1,4 | Feb-24 |
| 53. | Construction confidence, limiting factors, none | EUR | 1,4 | Feb-24 |
| 54. | Construction confidence, limiting factors, insufficient dem& | EUR | 1,4 | Feb-24 |
| 55. | Construction confidence, limiting factors, weather conditions | EUR | 1,4 | Feb-24 |
| 56. | Construction confidence, limiting factors, material/equipment | EUR | 1,4 | Feb-24 |
| 57. | Construction confidence, limiting factors, other | EUR | 1,4 | Feb-24 |
| **5. Money & Credit** | | | | |
| 58. | Loans, excl. government (EUR) | ECB | 3,6 | Jan-24 |

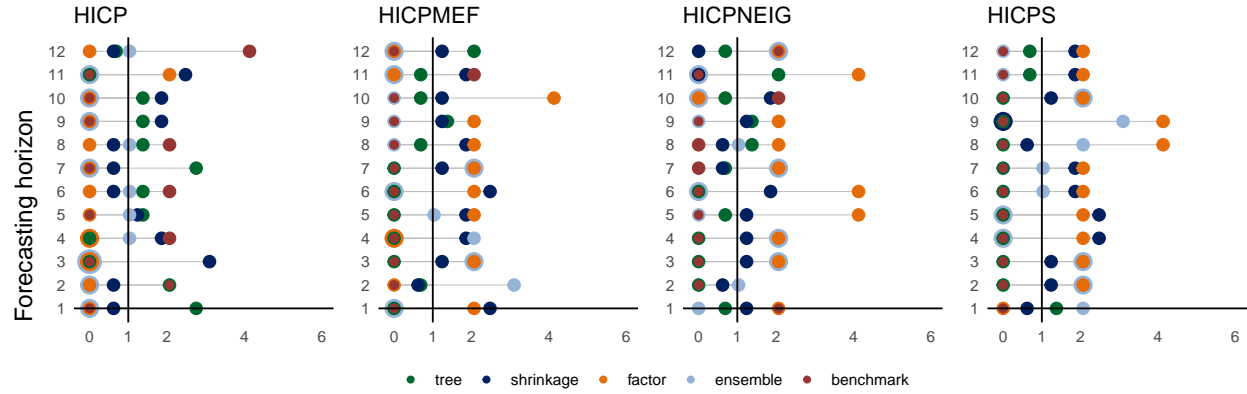| Nr. | Description | Source | Transf. | Last |
|---|---|---|---|---|
| 59. | External assets (EUR) | ECB | 3,6 | Jan-24 |
| 60. | External liabilities (EUR) | ECB | 3,6 | Jan-24 |
| 61. | Overnight deposits (EUR) | ECB | 3,6 | Jan-24 |
| 62. | Deposits <2 years, redeemable at notice < 3 months (EUR) | ECB | 3,6 | Jan-24 |
| 63. | Repo's, debt securities, shares < 2 years (EUR) | ECB | 3,6 | Jan-24 |
| 64. | M1 (EUR) | ECB | 3,6 | Jan-24 |
| 65. | M2 (EUR) | ECB | 3,6 | Jan-24 |
| 66. | M3 (EUR) | ECB | 3,6 | Jan-24 |
| **6. Interest & Exchange Rates** | | | | |
| 67. | Real effective exchange rate, deflator consumer price index | ECB | 2,5 | Feb-24 |
| 68. | Real effective exchange rate, deflator producer price index | ECB | 2,5 | Feb-24 |
| 69. | Nominal effective exchange rate euro | ECB | 2,5 | Feb-24 |
| 70. | 10–year government bond interest rate (%) | ECB | 1,4 | Feb-24 |
| 71. | 3–month interbank interest rate (%) | ECB | 1,4 | Feb-24 |
| 72. | 3–month deposits interest rate (%) | ECB | 1,4 | Feb-24 |
| 73. | Loans non-financial corporations, new, total (%) | ECB | 1,4 | Jan-24 |
| 74. | Loans non-financial corporations, new, $\leq$ EUR 1 million (%) | ECB | 1,4 | Jan-24 |
| 75. | Loans non-financial corporations, new, > EUR 1 million (%) | ECB | 1,4 | Jan-24 |
| 76. | Loans consumption, new, total (%) | ECB | 1,4 | Jan-24 |
| 77. | Loans house purchases, new, total (%) | ECB | 1,4 | Jan-24 |
| 78. | UK pound sterling/EUR exchange rate (%) | EUR | 2,5 | Feb-24 |
| 79. | Japenese yen/EUR exchange rate | EUR | 2,5 | Feb-24 |
| 80. | US dollar/EUR exchange rate | EUR | 2,5 | Feb-24 |
| **7. Commodity Prices** | | | | |
| 81. | Harmonized index of consumer prices, energy | ECB | 2,7 | Jan-24 |
| 82. | Europe brent spot price, USD (barrel) | ECB | 2,5 | Feb-24 |
| 83. | Food commodities | ECB | 2,5 | Feb-24 |
| 84. | Non-energy commodities | ECB | 2,5 | Feb-24 |
| 85. | Terms of trade | CBS | 2,5 | Dec-23 |
| 86. | Import prices | CBS | 2,5 | Dec-23 |
| 87. | Export prices | CBS | 2,5 | Dec-23 |
| **8. Producer Prices** | | | | |
| 88. | Producer price index, headline | EUR | 2,5 | Dec-23 |
| 89. | Producer price index, mining & quarrying | EUR | 2,5 | Dec-23 |
| 90. | Producer price index, food products, beverages & tobacco products | EUR | 2,5 | Dec-23 |
| 91. | Producer price index, textiles | EUR | 2,5 | Dec-23 |
| 92. | Producer price index, wearing apparel | EUR | 2,5 | Dec-23 |
| 93. | Producer price index, leather & related products | EUR | 2,5 | Dec-23 |
| 94. | Producer price index, wood & of products of wood & cork | EUR | 2,5 | Dec-23 |
| 95. | Producer price index, paper & paper products | EUR | 2,5 | Dec-23 |
| 96. | Producer price index, printing & reproduction of recorded media | EUR | 2,5 | Dec-23 |
| 97. | Producer price index, coke & refined petroleum products | EUR | 2,5 | Dec-23 |
| 98. | Producer price index, chemicals & chemical products | EUR | 2,5 | Dec-23 |
| 99. | Producer price index, rubber & plastic products | EUR | 2,5 | Dec-23 |
| 100. | Producer price index, other non-metallic mineral products | EUR | 2,5 | Dec-23 |
| 101. | Producer price index, basic metals | EUR | 2,5 | Dec-23 |
| 102. | Producer price index, fabricated metal products | EUR | 2,5 | Dec-23 |
| 103. | Producer price index, computer, electronic & optical products | EUR | 2,5 | Dec-23 |
| 104. | Producer price index, machinery & equipment | EUR | 2,5 | Dec-23 |
| 105. | Producer price index, motor vehicles & trailers | EUR | 2,5 | Dec-23 |
| 106. | Producer price index, electricity, gas, steam & air conditioning supply | EUR | 2,5 | Dec-23 |
| 107. | Producer price index, water collection, treatment & supply | EUR | 2,5 | Dec-23 |
| **9. Domestic Prices** | | | | |

Table A.1 – Continued from previous page

| Nr. | Description | Source | Transf. | Last |
|---|---|---|---|---|
| 108. | Harmonized index of consumer prices, headline | ECB | 2,7 | Jan-24 |
| 109. | Harmonized index of consumer prices, headline excl. energy & food | ECB | 2,7 | Jan-24 |
| 110. | Harmonized index of consumer prices, headline excl. energy | ECB | 2,7 | Jan-24 |
| 111. | Harmonized index of consumer prices, unprocessed food | ECB | 2,7 | Jan-24 |
| 112. | Harmonized index of consumer prices, food | ECB | 2,7 | Jan-24 |
| 113. | Harmonized index of consumer prices, services | ECB | 2,7 | Jan-24 |
| 114. | Harmonized index of consumer prices, industrial goods excl. energy | ECB | 2,7 | Jan-24 |
| 115. | Harmonized index of consumer prices, processed food | ECB | 2,7 | Jan-24 |
| 116. | Consumer price index, all items | ECB | 2,5 | Jan-24 |
| 117. | Construction costs, material prices, residential | CBS | 2,5 | Jan-24 |
| 118. | Housing prices | CBS | 2,5 | Dec-23 |
| 119. | Price materials new construction home | CBS | 2,5 | Jan-24 |
| 120. | Consumer confidence, price trend < 12 months | EUR | 1,4 | Feb-24 |
| **10. Price Expectations** | | | | |
| 121. | Consumer confidence, price trend > 12 months | EUR | 1,4 | Feb-24 |
| 122. | Construction confidence, price expectations > 3 months | EUR | 1,4 | Feb-24 |
| 123. | Consensus forecast consumer price index > 12 months | CF | 1,4 | Feb-24 |
| **11. Stock Market** | | | | |
| 124. | Amsterdam exchange index (AEX) | ECB | 2,5 | Feb-24 |
| 125. | Amsterdam midkap index | DS | 2,5 | Feb-24 |
| 126. | Dow Jones euro stoxx 50 index | ECB | 2,5 | Feb-24 |
| 127. | Financial stability index, Germany | ECB | 1,4 | Jan-24 |
| 128. | Financial stability index, United Kingdom | ECB | 1,4 | Jan-24 |
| 129. | Financial stability index, Netherlands | ECB | 1,4 | Jan-24 |

Notes: Nr.: Number of indicator; Description: Indicator description; Source: CBS: Statistics Netherlands, CF: Consensus Forecasts, DS: Refinitiv datastream, ECB: European Central Bank, EUR: Eurostat; Transf. = x,y: Transformation of variable in path average (x) and direct (y) forecast database, respectively, 0 = level, 1 = first difference, 2 = log difference, 3 = difference of log difference, 4 = annual difference, 5 = annual log difference, 6 = difference of annual log difference, 7 = annual percentage change; Last: Last monthly observation.
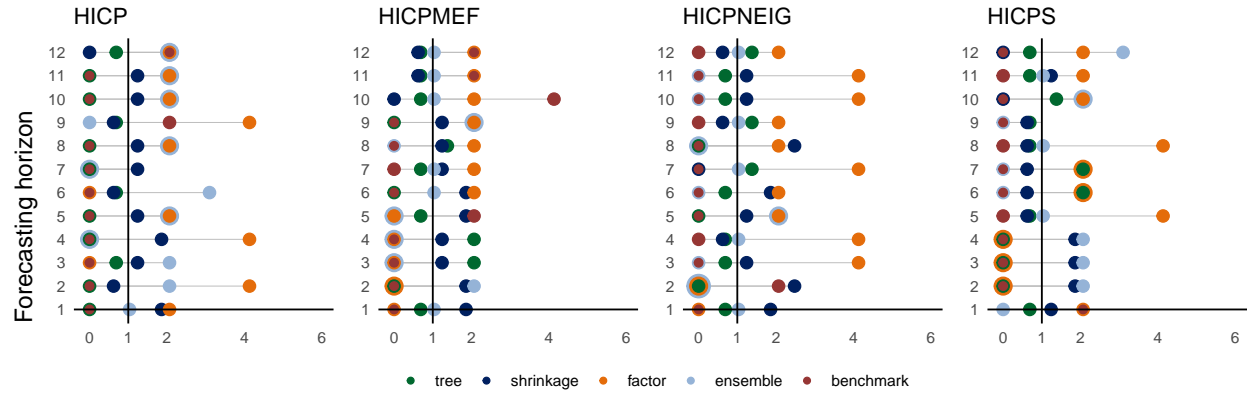
# B  Detailed MCS results

Figure B.1: Over & underrepresentation in top-5 MCS against random draw, full sample (2010M1–2023M12)[†]



[†] (nr. of models of modeltype in top-5 MCS (Hansen et al., 2011) with lowest p-value : [5 × nr. of horizons])/ (nr. of models of modeltype : total number of models); HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

Figure B.2: Over & underrepresentation in top-5 MCS against random draw, pre-pandemic sample (2010M1–2019M12)[†]



[†] (nr. of models of modeltype in top-5 MCS (Hansen et al., 2011) with lowest p-value : [5 × nr. of horizons])/ (nr. of models of modeltype : total number of models); HICP = headline HICP inflation, HICPMEF = HICP inflation excluding food & energy, HICPNEIG = non-energy industrial goods inflation, HICPS = services inflation.

# C   Additional OOS-R2 results

Figure C.1: Test of direct versus path average forecast, pre-pandemic sample (2010M1-2019M12)
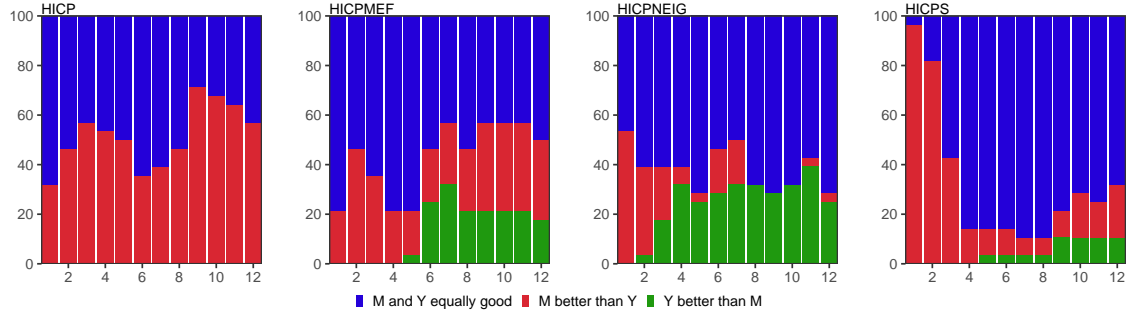


Figure C.2: Test of factor versus no factor forecast, pre-pandemic sample (2010M1-2019M12)
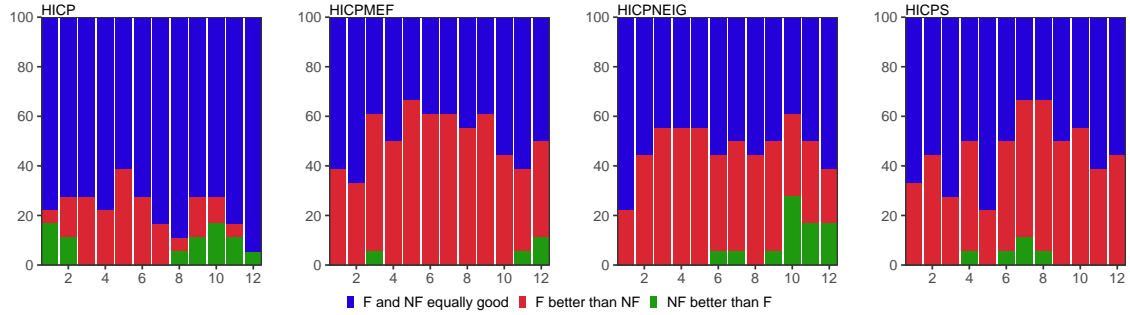


Figure C.3: Test of of targeting versus no targeting of predictors, pre-pandemic sample (2010M1-2019M12)
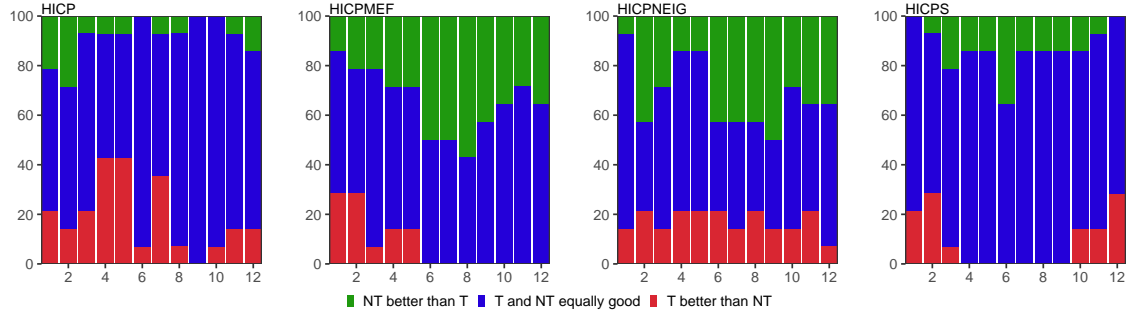
Figure C.4: Gain in OOS-R2 of path average forecast over direct forecast, pre-pandemic sample (2010M1-2019M12)
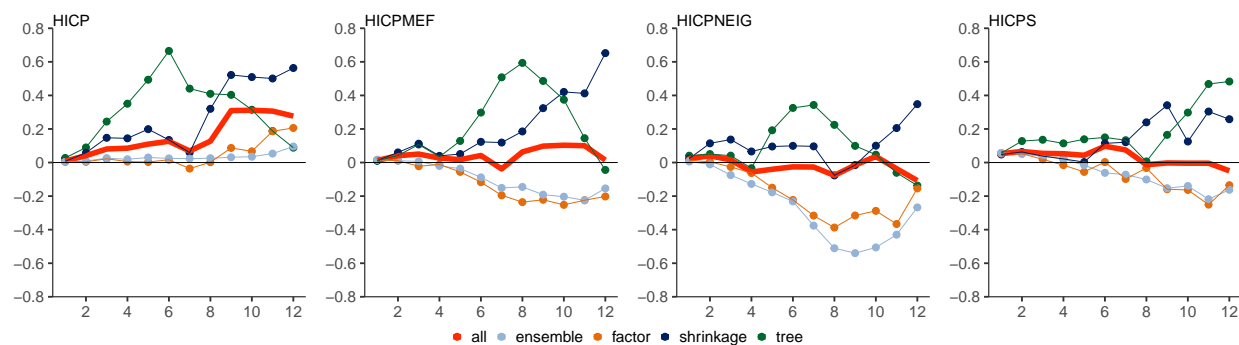


Figure C.5: Gain in OOS-R2 of factor versus no factor forecast, pre-pandemic sample (2010M1-2019M12)
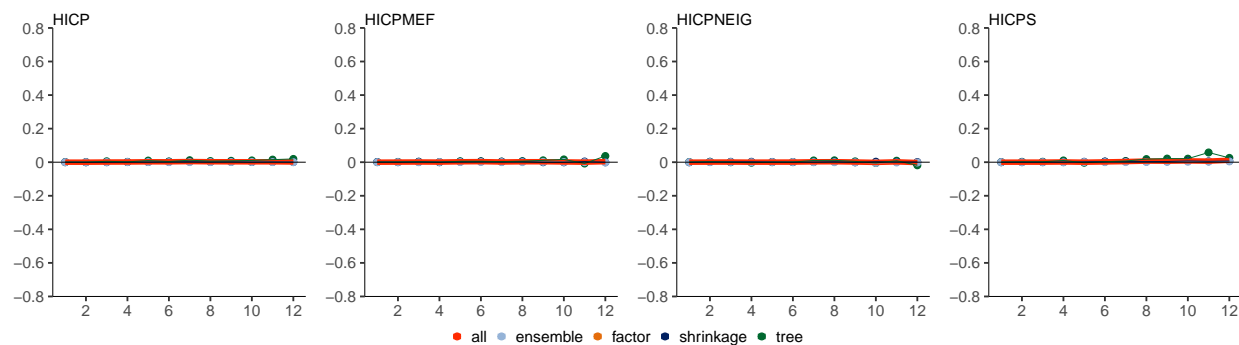


Figure C.6: Gain in OOS-R2 of targeting versus no targeting of predictors, pre-pandemic sample (2010M1-2019M12)