

# The Macroeconomy as a Random Forest

Philippe Goulet Coulombe\*  
University of Pennsylvania

First Draft: November 15, 2019

This Draft: November 4, 2020

[Latest Draft Here](#)

## Abstract

I develop *Macroeconomic Random Forest* (MRF), an algorithm adapting the canonical Machine Learning (ML) tool to flexibly model evolving parameters in a linear macro equation. Its main output, *Generalized Time-Varying Parameters* (GTVPs), is a versatile device nesting many popular nonlinearities (threshold/switching, smooth transition, structural breaks/change) and allowing for sophisticated new ones. The approach delivers clear forecasting gains over numerous alternatives, predicts the 2008 drastic rise in unemployment, and performs well for inflation. Unlike most ML-based methods, MRF is directly interpretable — via its GTVPs. For instance, the successful unemployment forecast is due to the influence of forward-looking variables (e.g., term spreads, housing starts) nearly doubling before every recession. Interestingly, the Phillips curve has indeed flattened, *and* its might is highly cyclical.

---

\*Department of Economics, [gouletc@sas.upenn.edu](mailto:gouletc@sas.upenn.edu). For many helpful discussions, I would like to thankgraphis Karun Adusumilli, Edvard Bakhitov, Francesco Corsello, Frank Diebold, Maximilian Göbel, Frank Schorfheide, Dalibor Stevanovic, David Wigglesworth and Boyuan Zhang. For excellent research assistance during different eras of this project, I am grateful to Tony Liu and JiWhan Moon. MRF is now available as an [R package](#).

# 1 Introduction

The rise of Machine Learning (ML) led to great excitement in the econometrics community. In applied macroeconomics, a first wave of papers took ML algorithms off the shelf and went hunting for forecasting gains. With the emerging consensus that some ML offerings can appreciably increase predictive accuracy, a question emerges: what is the place of economics in all that?

The conditional mean is the most basic input to any empirical macroeconomic analysis. Anything else that follows (e.g., structural analysis) depends on it. Thus, getting it right is not merely useful, it is *necessary*. Clearly, in that regard, ML can help. However, while the latter gladly delivers prediction accuracy gains (and ergo a conditional mean closer to the truth), it is much more reluctant to disclose its inherent model. Consequently, ML is currently of great use to macroeconomic forecasting, but of little help to macroeconomics. I propose a simple remedy: shifting the focus of the algorithmic arsenal away from predicting  $y_t$  into modeling  $\beta_t$ , which are economically meaningful coefficients in a time-varying macroeconomic equation. The newly proposed algorithm, *Macroeconomic Random Forest* (MRF) kills two coveted birds with one stone. First, in most instances, MRF forecasts better than off-the-shelf ML algorithms and traditional econometric approaches. Second, its main output, *Generalized Time-Varying Parameters* (GTVPs), can be interpreted. Their versatility comes from nesting many popular specifications (structural breaks/change, threshold effects, regime-switching, etc.) and letting the data decide whichever combination of them is most suitable. Ultimately, we get a new methodology leveraging the power of ML and big data to provide a modern take on the decades-old challenge of estimating latent states driving linear macroeconomic equations.

**THE STATE OF EMPIRICAL MACRO AFFAIRS.** Answering positively two questions guarantees a viable conditional mean: "are all the relevant variables included in the model?" and at a higher level of sophistication, "is linearity a valid approximation of reality?". The first one led to the successful development of factor models and large Bayesian Vector Autoregressions (VARs) over the last two decades. To address the second, applied macroeconomic researchers have proposed many non-linear time series models based on reasonable economic intuition. Most of them amount to have regression coefficients  $\beta_t$  in

$$y_t = X_t\beta_t + \varepsilon_t$$

evolving through time. The  $\beta_t$  process can take many forms, and a choice must be made *a priori* out of many equally plausible alternatives. Notable members of the vast time-variations catalog are threshold/switching regressions (Hansen, 2011), smooth transition (Teräsvirta, 1994), structural breaks (Perron et al., 2006; Stock, 1994), and random walk time-varying parameters (Sims, 1993; Cogley and Sargent, 2001; Primiceri, 2005). While it is uncontroversial that factor

models and large Bayesian VARs have gone a long way in meeting their original goals, less victorious statements are available for the various time-variation proposals. Why?

More often than not, nonlinear time series models use little data and/or restrict stringently the shape of  $\beta_t$ 's path. While the consequences for forecasting are direct and obvious, those for analysis of macroeconomic relationships are equally problematic. Is the evolving Taylor rule characterized by switching regimes (Sims and Zha, 2006), a Volker structural break (Clarida et al., 2000), or gradually evolving parameters (Boivin, 2005; Primiceri, 2005)? This discordance interferes with our understanding of the past while impacting our expectations for tomorrow's  $\beta_t$ . I now divide popular time-variation approaches into two strands, discuss their shortcomings, and complete by explaining how MRF addresses them.

**OBSERVABLE TIME-VARIATION VIA INTERACTION TERMS.** Using interaction terms and related refinements is a parsimonious way to create time variation in a linear equation. For instance, switching regimes based on an observed regressor can be obtained by interacting the linear equation with the indicator function  $I(q_t > c)$ , where  $c$  is some value, and  $q_t$  is a threshold variable chosen by the researcher. However, using the FRED-QD US macro data set (McCracken and Ng, 2016) reveals an overwhelmingly large number of candidates for  $q_t$ . Additionally, there may be multiple regimes interacting together. Or the "true"  $q_t$  could be an unknown function of available regressors. And structural breaks or slow exogenous variation could get in the way. The list goes on. This renders a credible exploration of the threshold structures' space impossible and the enterprise of manually specifying the model very much compromised.

Here is an empirical example. Auerbach and Gorodnichenko (2012b) and Ramey and Zubairy (2018) use a GDP/unemployment indicator to let the effects of fiscal stimulus (potentially) vary with the state of the economy. Batini et al. (2012) allow for additional dependence on the origin of the impulse (revenue or spending). Such honorable explorations could go on endlessly. MRF provides a hammer solution to the problem. First, the near-universe of threshold structures can be characterized by regression trees — see section 2.1. Second, MRF embeds, among other things, a powerful greedy algorithm designed to explore such "structure" spaces.

**LATENT TIME-VARIATION.** Some methods with an aura of greater flexibility are labeled as "latent change". In this line of work,  $\beta_t$  either follows a law of motion (random walk, Markov process) or could be subject to discrete breaks.<sup>1</sup> At first glance, this appears to solve many of the problems of interaction terms approaches. By treating  $\beta_t$  as a state to be filtered/estimated within the model, the complexity of characterizing its path correctly out of abundant data seems to vanish. Alas, estimating  $\beta_t$ 's path implies a great number of parameters (in fact, often greater than the number of observations, Goulet Coulombe 2020a) which inevitably necessitates strong

---

<sup>1</sup>Simpler derivatives are often used in applied work. In forecasting, rolling-window estimation drops early observations. In empirical macro, pre-defined subsamples are popular (Clarida et al., 2000; Del Negro et al., 2020).

regularization. That regularization is the law of motion itself, a choice far from innocuous – and akin to that of  $q_t$  in "observable" change models. Accordingly, whether it is latent regime-switching, exogenous breaks, or slow change, none can easily accommodate for the additional presence of the other. Yet, these models are routinely fitted *separately* on the *same data*. Consequently, methods often detect what they are designed to detect, in near-complete abstraction of imaginable interference from other nonlinearities.

Additionally, while "latent" approaches may sometimes rationalize the data well in-sample, many of them will struggle to outperform a simple benchmark *out-of-sample*. Often, the very nature of  $\beta_t$ 's law of motion creates forecasting headaches. Classical TVPs imply a two-sided vs one-sided filtering problem. Analogously, detecting a structural break is much harder without a great amount of data on both sides of it. Moreover, there is the obvious problem of statistical efficiency. If the Phillips curve flattened because an economy became increasingly open, including an interaction term with imports/exports is wildly more efficient than obtaining the whole  $\beta_t$  path non-parametrically. Thus, exogenous structural change should be, in some sense, a time variation of last resort. The advantage of MRF is that it algorithmically search for "observable" low-hanging fruits, and turn to split the sample with  $t$  only if necessary. Further, it implicitly creates a forecasting function for  $\beta_t$  which is an RF in its own right. This is, almost in any case, much more powerful than existing alternatives – like random walks.

**MECHANICS.** The key difference when adding the M to MRF is the inclusion of a linear part within each of the tree leaves, rather than just an intercept. Motivated in cross-sectional applications to improve the efficiency of nonparametric estimation (in the spirit of local linear regression), trees with linear parts have been considered (among others) in [Alexander and Grimshaw \(1996\)](#) and [Wang and Witten \(1996\)](#). [Friedberg et al. \(2018\)](#) expand on this by considering an ensemble of them (i.e., a forest) and focusing on the problem of treatment effect heterogeneity. Of course, the difference here is that a linear part is much more meaningful when one can look at  $\beta_t$  as a process of its own – and as a synthesis of nonlinear time series models. Finally, it is noteworthy that the approach may come in semiparametric partially linear clothing, yet it makes no compromise on the range of nonlinearities it captures. This is a virtue of time-varying coefficients models being able to approximate any nonlinear function ([Granger, 2008](#)).

The paper also introduces new devices enhancing MRF's predictive and interpretability potential. First, I propose Moving Average Factors (MAFs) as a simple way to compress ex-ante the information contained in the lags of a regressor entering the RF part of MRF. They boost the meaningfulness of tree splits and helps avoid running out of them quickly. The transformation is motivated by the literature on constraining/regularizing lag polynomials ([Shiller, 1973](#)). Precisely, MAFs' contribution is to induce similar shrinkage when there are no explicit coefficients to shrink. When it comes to GTVPs themselves, I provide a regularization scheme better suited

for time series which procures a desirably smoother path with respect to time. It is inspired by the random walk shrinkage of the classical TVP literature and is implemented within the tree procedure by weighted least-squares. Finally, a variant of the Bayesian Bootstrap provides credible regions that are instrumental for the interpretation of GTVPs.

**RESULTS.** In simulations, the tool does comparably well to traditional nonlinear time series models when the data generating process (DGP) matches what the latter is designed for. When the time-variation structure becomes out of reach for classical approaches, MRF wins. Additionally, it supplants plain RF whenever persistence is pervasive. In a forecasting application, the MRFs gains are present for almost all variables and horizons under study, a rarity for nonlinear forecasting approaches. For instance, the Autoregressive Random Forest (ARRF) almost always supplant its resilient OLS counterpart. Also, an MRF where the linear part is a compact factor-augmented autoregression generates very accurate forecasts of the 2008 downturn for both GDP and the unemployment rate (UR). Inspection of resulting GTVPs reveals they behave differently from random walk TVPs. For instance, in the UR equation, the contribution of forward-looking variables nearly doubles before every recession — including 2008 where the associated  $\beta_t$  is forecasted to do so out-of-sample. This reinforces the view that financial indicators and other market-based expectations proxies can rapidly capture downside risks around business cycle turning points (Adrian et al., 2019). MRF learned and applied it to great success.

Inflation is subject to a variety of time-variations, detection of which would be compromised by approaches lacking the generality of MRF. The long-run mean and the persistence evolved slowly and in an exogenous fashion — this has been repeatedly found in the literature (e.g., Cogley and Sargent 2001). More novel is the finding that the real activity factor's effect on the price level depends positively on the strength of well-known leading indicators, especially housing-related. Following this lead, I complete the analysis by looking at a traditional Phillips' curve specification. I report that the inflation/unemployment trade-off coefficient decreased significantly since the 1980s and also varies strongly along the business cycle. Among other things, it is extremely weak following every recession. This nuances current evidence on the flattening Phillips curve, which, by design, focused almost entirely on long-run exogenous change (Blanchard et al., 2015; Galí and Gambetti, 2019; Del Negro et al., 2020). Overall, MRF suggests inflation can rise from a positive unemployment gap, but it goes down much more timidly from economic slack. These findings are made possible by combining different tools within the new framework, such as credible intervals for the GTVPs, new variable importance measures specifically designed for MRF, and surrogate trees as interpretative devices for  $\beta_t$ .

**OUTLINE.** Section 2 introduces MRF, motivates its use, considers practical aspects, and discusses relationships with available alternatives. Sections 3 and 4 report simulations and forecasting results, respectively. Section 5 analyzes various GTVPs of interest. Section 6 concludes.

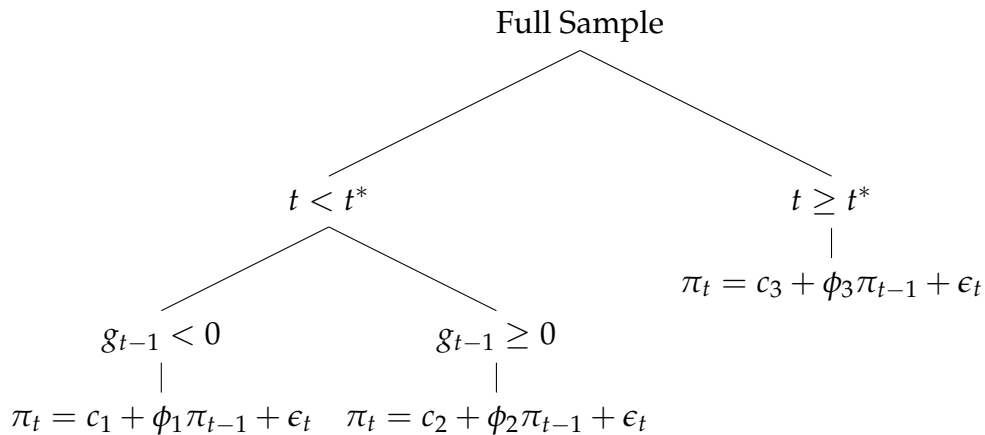
## 2 Macroeconomic Random Forests

This section introduces MRF. I first motivate the use of trees as basis functions by casting standard switching structures for autoregressions as special cases. Second, I detail the MRF mechanics and how it yields GTVPs. Third, I discuss how the approach relates to both standard RF and traditional random walk TVPs. Fourth, I discuss interpretability potential and provide a way to assess parameter uncertainty.

### 2.1 Traditional Macro Non-Linearities as Trees

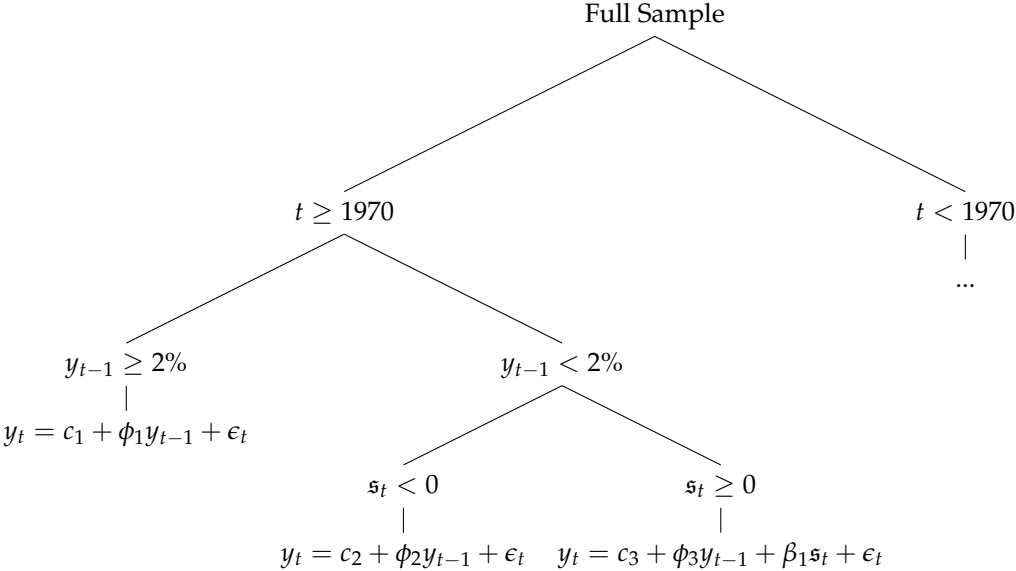
Within the modern ML canon, Random Forest (RF) is an extremely popular algorithm because it allows for complex nonlinearities, handles high-dimensional data, bypasses overfitting, and requires little to no tuning. This is in sharp contrast with, for example, Neural Networks, whose ability to fail upon a bad choice of hyperparameters is largely unmatched. Thus, RF is a reasonable device to look into for constructing GTVPs. But there is more: many common time series nonlinearities fit within a tree structure. Hence, it will be all the more natural to think of MRF as a generalization of previous nonlinear offerings. Overall, it eliminates the arbitrary search for a specification. By creating a unified view, the myriad of time-variations suggested separately can now be tackled jointly.

I now present two examples displaying how common time series nonlinearities imply a tree structure for an AR process. Let us consider the inflation process in a country where inflation targeting (IT) was implemented at a publicly known date (like in Canada). Let  $\pi_t$  be inflation at time  $t$  and  $t^*$  is the onset date of IT. Additionally,  $g_t$  is some measure of output gap. A plausible model is reported in the tree graph below. The story is straightforward. Inflation behaved differently before vs after IT. After IT, it is a simple AR process. Before IT, it was a switching AR process which dynamics and mean depended on the sign of the output gap.<sup>2</sup>



<sup>2</sup>Note that a standard regression tree would set all  $\phi$ 's to 0.

Here is a second (more intricate) example loosely inspired by [Auerbach and Gorodnichenko \(2012a\)](#), [Ramey and Zubairy \(2018\)](#) and others. Let  $y_t$  be GDP growth at time  $t$  and  $s_t$  be some measure of government spending shock. The tree below tells us that only data post-1970 is of "current" interest — the high-growth environment of pre-1970 being characterized by a different process. The effect of spending  $s_t$  on growth  $y_t$  depends on two variables: previous growth  $y_{t-1}$  (the state of the economy) and whether government spending  $s_t$  is expanding or contracting. Hence, this tree allows for different mean/dynamics of growth and state-dependent effects of spending conditional on three variables:  $t$ ,  $y_{t-1}$  and  $s_t$ .



These are two stories out of many that trees can characterize. In practice, none of the above is known. The structure, the splitting variables, and the splitting points could be different. This is both good and bad news. It highlights the flexibility of trees. It also suggests that designing the "true" one from economic deduction is a daunting task — equally plausible alternatives are easily imaginable. Fortunately, algorithms can point out which trees in better agreement with the data.

A global grid search is computationally unfeasible if either  $S_t$  is large or if we want to consider more than a few splits (examples above included 2 and 3, respectively). A natural way forward is recursive partitioning of the data set via a *greedy* algorithm ([Breiman et al., 1984](#)).<sup>3</sup> A greedy algorithm optimizes functions by iteratively doing the best local update, rather than directly solving for a global optimum. As a result, it is prone to high variance ([Friedman et al., 2001](#)). Hence, considering a diversified portfolio of trees appears as the most sensible route. To achieve that, it is highly effective to use Bootstrap Aggregation (*Bagging*, [Breiman 1996](#)) of many de-correlated trees. This is the famous Random Forest proposition of [Breiman \(2001\)](#).

<sup>3</sup>A single autoregressive tree was proposed in [Meek et al. \(2002\)](#).

## 2.2 Generalized Time-Varying Parameters

The general model is

$$\begin{aligned} y_t &= X_t \beta_t + \epsilon_t \\ \beta_t &= \mathcal{F}(S_t) \end{aligned}$$

where  $S_t$  are the state variables governing time variation and  $\mathcal{F}$  a forest.  $S_t$  is observed macroeconomic data which composition is motivated in section 2.6 and laid out explicitly in section 4.  $X$  determines the *linear* model that we want to be time-varying. For instance, an autoregressive random forests (ARRF) – which generalizes the cases of the previous section – uses lags of  $y_t$  for  $X_t$ . The tree fitting procedure underlying *plain* RF is not adequate, as it sets  $X_t = \mathbf{1}$  by default. Thus, analogously to [Friedberg et al. \(2018\)](#), it is modified to

$$\begin{aligned} \min_{j \in \mathcal{J}^-, c \in \mathbb{R}} \left[ \min_{\beta_1} \sum_{\{t \in l | S_{j,t} \leq c\}} (y_t - X_t \beta_1)^2 + \lambda \|\beta_1\|_2 \right. \\ \left. + \min_{\beta_2} \sum_{\{t \in l | S_{j,t} > c\}} (y_t - X_t \beta_2)^2 + \lambda \|\beta_2\|_2 \right]. \end{aligned} \quad (1)$$

The purpose of this problem is to find the optimal variable  $S_j$  (so, finding the best  $j$  out of the random subset of predictors indexes  $\mathcal{J}^-$ ) to split the sample with, and at which value  $c$  of that variable should we split.<sup>4</sup> It outputs  $j^*$  and  $c^*$  which are used to split  $l$  (the parent node) into two children nodes,  $l_1$  and  $l_2$ . We start with the leaf  $l$  being the full sample. Then, we perform a split according to the minimization problem, which procures us with 2 subsamples. Within each of these two newly created subsamples, we run (1) again. Repeating this process recursively constructs an ever-growing set of  $l$ 's which are of ever-shrinking size. Doing so until a stopping criteria is met generates a tree.

**LET THE TREES RUN DEEP.** Recursively splitting  $\beta_0$  into  $\beta_1$  and  $\beta_2$  eventually leads to  $\beta_t$ . However,  $\beta_t$ , by construction, has very little company within its terminal node/leaf. As result, a single tree has low bias, but also very high variance for  $\beta_t$ . When fitting a single tree, the (early) stopping point must be tuned to avoid overfitting. However, this is not necessary when a sufficiently diversified ensemble of trees is considered. Originally, [Breiman \(2001\)](#) himself provided a bound on the generalization error that grows with the correlation between trees.<sup>5</sup> In

<sup>4</sup>Note that, unlike [Friedberg et al. \(2018\)](#),  $S_t$  and  $X_t$  will differ, which is natural when motivated from a TVP perspective (but not so much from local linear regression one). Forcing their equivalence is not feasible nor desirable in a macro environment.

<sup>5</sup>Also, [Duroux and Scornet \(2016\)](#) derive a formula (for a "median" forest) linking tuning parameters related to the depth of the trees and that of diversification.



Goulet Coulombe (2020b), I go further by showing that RF's out-of-sample prediction is equivalent to the optimally "stopped" or "pruned" one, provided sufficiently diversified trees. The desirable property is attributed to the peculiar behavior of "randomized greedy algorithms", which are often overlooked as mere computational necessities. Those insights are of even greater use when it comes to time series since dependence and structural change pose challenges to hyperparameter tuning. Given a large enough  $B$ , a reasonable `mtry` and standard subsampling rate, we can be confident that the out-of-bag prediction and  $\beta_t$ 's exclude fitted noise. In our specific context, it means the sample will not be over-split, and we are not going to see time variation when it is not there. Naturally, the credible regions proposed in section 2.7 will also help in that regard. The property will be illustrated in section 3.2.

(M)RF prediction is the *simple average* from those of its single trees. Same goes for  $\beta_t$ . RF is a clever diversification scheme which generates sufficient randomization for that average to inherit the above properties. To achieve that, it mixes elements of re-sampling and model averaging: Bagging and de-correlated trees.<sup>6</sup>

**BAGGING.** Each tree is "grown" on a bootstrapped sample (or a random subsample) (Breiman, 1996). When the base learner is highly nonlinear in observation and/or unstable, gains from Bagging can be large (Breiman, 1996; Grandvalet, 2004). Nonparametric (or "pairs" MacKinnon 2006) bootstrap is being used — i.e., we are *not* shuffling residuals.<sup>7</sup> Rather, we are randomly selecting many observations triples  $[y_t \ X_t \ S_t]$  (or pairs  $[y_t \ S_t]$  for Plain RF), and then fit a tree on them. For reasons to be detailed in section 2.7, a slightly more sophisticated bootstrapping/-subsampling procedure will be used for MRF.

**DE-CORRELATION.** The second ingredient, proposed in Breiman (2001), is to consider "de-correlated" trees. RF is an average of many trees, and any averaging scheme reduces variance at a much faster rate if its components are uncorrelated. In our context, this is obtained by growing trees semi-stochastically. In equation (1), this is made operational by using  $\mathcal{J}^- \subset \mathcal{J}$  rather than  $\mathcal{J}$ . In words, this means that at each step of the recursion, a different subsample of regressors is drawn to constitute candidates for the split. This prevents the greedy algorithm (which, as we know, only "thinks" locally) to always embark on the same optimization route. As a result, trees are further diversified and computing time, reduced. The fraction of randomly selected predictors is a tuning parameter typically referred to as `mtry` in the literature (and all software), with a default value of  $\frac{1}{3}$  for regression settings. This, other algorithmic parameter settings, and some practical aspects are discussed in appendix A.4.

---

<sup>6</sup>See Goulet Coulombe (2020b) for a discussion on how RF compares and contrast with the forecast combination-/averaging literature.

<sup>7</sup>Nonetheless, Bagging in itself is not estranged to macro forecasting (Inoue and Kilian, 2008; Hillebrand and Medeiros, 2010; Hillebrand et al., 2020). However, nearly all studies consider the more common problem of variable selection via hard-thresholding rules – like t-tests (Lee et al., 2020).

Plain RF has many qualities readily transferable to MRF. It is easy to implement and to tune. That is, it has few tuning parameters that are usually of little importance to the overall performance – robustness. It is relatively immune to the adverse effects of including many irrelevant features (Friedman et al., 2001). Given the standard ratio of regressors to observations in macro data, this is a non-negligible advantage. Furthermore, with a sufficiently high  $m_{\text{TRY}}$ , it can adapt nicely to sparsity and discard useless predictors (Olson and Wyner, 2018). Finally, its vanilla version already shows good forecasting performance for US inflation (Medeiros et al., 2019) and macro data in general (Chen et al., 2019; Goulet Coulombe et al., 2019).

### 2.3 Random Walk Regularization

Equation (1) uses Ridge shrinkage which implies that each time-varying coefficient is implicitly shrunk to 0 at every point in time.  $\lambda$  and the prior it entails can exert a significant influence. For instance, if a process is highly persistent (AR coefficient lower than 1 but nevertheless quite high) as it is the case for SPREAD (see section 4), shrinking the first lag heavily to 0 could incur serious bias. Fortunately, this can easily be refined to a Minnesota-style prior if  $X_t$  corresponds to a Bayesian VAR equation. If  $X_t$  is low-dimensional (as it will often be), a simpler alternative consists in using OLS coefficients as prior means. Nonetheless, the specification of previous sections implies that if  $\lambda$  grows large,  $\forall t \beta_t = 0$  (or whatever the prior mean is).  $\beta_i = 0$  is a natural stochastic constraint in a cross-sectional setting, but its time series translation  $\beta_t = 0$  can easily be suboptimal. The traditional regularization employed in macro is rather the random walk

$$\beta_t = \beta_{t-1} + u_t.$$

Thus, it is desirable to transform (1) so that it implements the prior that coefficients evolve smoothly, which is just shrinking  $\beta_t$  to be in the neighborhood of  $\beta_{t-1}$  and  $\beta_{t+1}$  rather than 0. The random walk regularization ensure that the parameter's path will be smooth to some extent. This is in line with the view that economic states (as expressed by  $\beta_t$  here) last for at least a few consecutive periods. Moreover, such shrinkage will greatly facilitate interpretation of resulting GTVPs.

I implement the desired regularization by taking the "rolling-window view" of time-varying parameters. That is, the tree, instead of solving a plethora of small ridge problems, will rather solve many weighted least squares problems (WLS) which includes close-by observations. The latter are in the neighborhood (in time) of observations within current leaf. They are included in estimation, but are allocated a smaller weight.

For simplicity and to keep computational demand low, the kernel used by WLS is rather rudimentary: it is a symmetric 5-step Olympic podium. Informally, the kernel puts a weight of

1 on observation  $t$ , a weight of  $\zeta < 1$  for observations  $t - 1$  and  $t + 1$  and a weight of  $\zeta^2$  for observations  $t - 2$  and  $t + 2$ . Since some specific  $t$ 's will come up many times (for instance, if both observations  $t$  and  $t + 1$  are within the same leaf, podiums overlap), I take the maximal weight allocated to  $t$  as the final weight  $w(t; \zeta)$ .

Formally, define  $l_{-1}$  as the "lagged" version of leaf  $l$ . In other words,  $l_{-1}$  is a set containing each observation from  $l$ , with all of them lagged one step.  $l_{+1}$  is the "forwarded" version.  $l_{-2}$  and  $l_{+2}$  are two-steps equivalents. For a given candidate subsample  $l$ , the podium is

$$w(t; \zeta) = \begin{cases} 1, & \text{if } t \in l \\ \zeta, & \text{if } t \in (l_{+1} \cup l_{-1})/l \\ \zeta^2, & \text{if } t \in (l_{+2} \cup l_{-2}) / (l \cup (l_{+1} \cup l_{-1})) \\ 0, & \text{otherwise} \end{cases}$$

where  $\zeta < 1$ , a tuning parameter guiding the level of time-smoothing. Then, it is only a matter of how to include those additional (but down weighted) observations in the tree search procedure. The usual candidate splitting sets

$$l_1(j, c) \equiv \{t \in l | S_{j,t} \leq c\} \quad \text{and} \quad l_2(j, c) \equiv \{t \in l | S_{j,t} > c\}$$

are expanded to include all observations of relevance to the podium

$$\text{for } i = 1, 2: \quad l_i^{RW}(j, c) \equiv l_i(j, c) \cup l_i(j, c)_{-1} \cup l_i(j, c)_{+1} \cup l_i(j, c)_{-2} \cup l_i(j, c)_{+2}.$$

The splitting rule becomes

$$\min_{j \in \mathcal{J}^-, c \in \mathbb{R}} \left[ \min_{\beta_1} \sum_{t \in l_1^{RW}(j, c)} w(t; \zeta) (y_t - X_t \beta_1)^2 + \lambda \|\beta_1\|_2 \right. \\ \left. + \min_{\beta_2} \sum_{t \in l_2^{RW}(j, c)} w(t; \zeta) (y_t - X_t \beta_2)^2 + \lambda \|\beta_2\|_2 \right]. \quad (2)$$

Note that the Ridge penalty is kept in anyway, so the final model has in fact two sources of regularization. With  $\zeta \rightarrow 0$ , we are heading back to pure Ridge.

Although not considered in the main applications of this paper, models with a larger linear part  $X_t$  are possible. For instance, one could estimate, equation by equation, a high-dimensional VAR. In practice, this simply requires harsher regularization via higher values of  $\lambda$ ,  $\zeta$  and a larger minimum leaf size. Nevertheless, the forecasting benefits from this strategy could prove limited: MRF is "high-dimensional" whenever  $S_t$  is large. The time-varying constant in MRF is

a RF in its own right. It can be seen as a complex misspecification function (in the deep learning jargon, it is effectively called the bias) that adaptively controls for omitted variables in a way that is both non-linear and strongly regularized via randomization. Consequently, the cost from omitting a regressor of minor importance in  $X_t$  is low since it can be picked up by the time-varying intercept.

Of course, the small  $X_t$  strategy treats the extra regressors as exogenous, which could be at odds with some researchers' will to investigate a large web of impulse response functions. Anyhow, both approaches are possible. It turns out large VAR specifications also deliver good results within MRF. In appendix A.1, the high-dimensional VAR MRF (HD-VARRF) provides the best 1-year ahead forecasts for both unemployment and GDP – signaling a (albeit smaller than realized) recession up to a year ahead.

## 2.4 Relationship to Random Walk Time-Varying Parameters

GTVPs have many advantages over classical TVPs. While it is known that any nonlinear model can be approximated by a linear one with TVPs (Granger, 2008), nothing is said about how efficient that estimation is going to be. As it turns out, efficiency crucially matters in a macro context, and random-walk TVPs can be quite inefficient (Aruoba et al., 2017). For example, if the true  $\beta_t$  follows a recurrent switching mechanism, random walk parameters already have two strikes against them. Some dimensionality reduction techniques – like reduced-rank restrictions (de Wind and Gambetti, 2014; Stevanovic, 2016; Chan et al., 2018; Goulet Coulombe, 2020a) – can help, but nothing in that paradigm can come close to the parsimony of simply interacting  $X_t$  with relevant variables. In contrast, MRF considers all time-variations options, and choose the "obvious thing", which may or may not be splitting on  $t$ . Also, it is absolutely possible that the resulting  $\mathcal{F}$  pools *both* latent and observable time variation.

Even though MRF is remarkably flexible, its variance remains low thanks to the diversified portfolio of trees. The variance of classical TVPs can be controlled by cross-validation (Goulet Coulombe, 2020a) or via an elaborate hierarchical prior (Amir-Ahmadi et al., 2018). A number of applications opt for a "manual" approach (D'Agostino et al., 2013). However, it is understood that no tuning, however careful it may be, can overcome the hardship of fitting random-walks when the true  $\beta_t$ 's look nothing like it.

Econometrically, one way to more formally connect this paradigm to recent work on TVPs is to adopt the view that RF are adaptive kernel estimators (Meinshausen, 2006; Athey et al., 2019; Friedberg et al., 2018). That is, the tree ensemble is a machine generating kernel weights. Once those are obtained, estimation amounts to weighted least squares (WLS) problem with a Ridge penalty. By running (1) recursively, one obtains terminal nodes/leaves  $L_b()$  to construct kernel

weights

$$\alpha_t(x_0) = \frac{1}{B} \sum_{b=1}^B \frac{1 \{X_t \in L_b(x_0)\}}{|L_b(x_0)|}$$

to use in

$$\forall t : \operatorname{argmin}_{\beta_t} \left\{ \sum_{\tau=1}^T \alpha_t(\mathbf{s}_\tau) (Y_\tau - X_\tau \beta_\tau)^2 + \lambda \|\beta_t\|_2 \right\}. \quad (3)$$

As shown in [Goulet Coulombe \(2020a\)](#), standard random walk TVPs are in fact a smoothing splines problem, and for those, a reproducing kernel exists ([Dagum and Bianconcini, 2009](#)). [Giraitis et al. \(2014\)](#) drop the random walk altogether and proposed to use kernels directly. Anyhow, in both cases, the only variable entering the kernel is  $t$ . In other words, only proximity in time is considered for the clustering of observations. This makes the seemingly flexible estimator in fact quite restrictive – and dependent on its inherent smoothness prior. Moreover, standard kernel methods are known to break down even in medium dimensions (say <10 variables) ([Friedberg et al., 2018](#)). Therefore, augmenting  $t$  with additional regressors is not an option. No such constraints bind on the RF approach.

## 2.5 Relationship to Standard Random Forest

The standard RF is a restricted version of MRF where  $X_t = \iota$ ,  $\lambda = 0$  and  $\zeta = 0$ . In words, the only regressor is a constant and there is no within-leaf shrinkage. Previous sections motivated MRF as a natural generalization of non-linear time-series models. At this point, a reasonable question emerges from a ML standpoint. Why should we prefer the partially linear MRF to the fully nonparametric RF? One reason is statistical efficiency. The other is potential for interpretation.

### 2.5.1 Smooth Relationships are Hard Relationships (to estimate)

In finite samples, plain RF can have a hard time learning smooth relationships – like a AR(1) process. This is bad news for time series applications. For prediction purposes, estimating

$$y_t = \phi y_{t-1} + \varepsilon_t$$

by OLS implies a single parameter. However, approximating the same relationship with a tree (or an ensemble of them) is far more consuming in terms of degrees of freedom. To get close to the straight line once parsimoniously parametrized by  $\phi$ , we now need a succession of many step functions.<sup>8</sup> With short time series, modeling smooth/linear relationships in such a way is a

---

<sup>8</sup>In a standard regression setup, nobody would model a continuous variable as an ordinal one unless some wild nonlinearities are suspected.

luxury one rarely can afford. The mechanical consequence is that RF will waste many splits on capturing the linear part, and may run out of them before it gets to focus more subtle nonlinear phenomena.<sup>9</sup> In a language more familiar to economists, this is simply running out (quickly) of degrees of freedom. MRF provides a workaround. Modeling the linear part concisely leaves more room to estimate the nonlinear one. By its more strategic budgeting of degrees of freedom, the resulting (estimated) partially linear model could be, in fact, more non-linear than the fully nonparametric one.

This paper is not the first to recognize the potential need for a linear part in tree-based models. For instance, both [Alexander and Grimshaw \(1996\)](#) and [Wang and Witten \(1996\)](#) proposed linear regressions within a leaf of a tree, respectively denominated "Treed Regression" and "Model Trees". More focused on real activity forecasting, [Woloszko \(2020\)](#) and [Wochner \(2020\)](#) blend insights from macroeconomics to build better-performing tree-based models.<sup>10</sup> On a different end of the econometrics spectrum, [Friedberg et al. \(2018\)](#) proposed to improve the nonparametric estimation of treatment effect heterogeneity by combining those ideas developed for trees into a forest.<sup>11</sup> To my knowledge, this paper is the first to exploit the link between this strand of work and the sempiternal search for the "true" state-dependence in empirical macroeconomic models.

## 2.5.2 A Note on Interpretability

The interpretation of ML outputs is now a field of its own ([Molnar, 2019](#)). RF is widely regarded as a black box model which needs to be interpreted using an external device. Indeed, it usually averages over 100 trees of substantial depth, which makes individual inspection impossible. MRFs partially circumvent the problem by providing time series  $\beta_t$  which can be examined, and have a meaning as time-varying parameters for the linear model. Thus, whatever one may do with TVPs, it can be done with GTVPs. There are also some new avenues. For instance, Variable Importance (VI) measures usually deployed to dissect RF's prediction can be used to inspect what is driving  $\beta_t$ 's. Those will be used in section 5.3.

A popular approach to dissect a standard RF is to use interpretable surrogate tree models to partially replicate the black box model's fit. The idea can be transferred to MRF ([Molnar, 2019](#)). In fact, partial linearity facilitates such an exercise. The linear part in MRF splits the nonparametric atom into different pieces ( $X_{t,k}\beta_{t,k}$ ) which can be analyzed separately. Each time series  $\beta_{t,k}$  can be dissected with its own surrogate model, and meaningful combination/transformations

---

<sup>9</sup>One necessary (but not sufficient) symptom is AR terms being flagged as really important by typical RF variable importance measures (one example is [Borup et al. \(2020b\)](#)).

<sup>10</sup>Specifically, [Wochner \(2020\)](#) also note that using trees in conjunction with factor models can improve GDP forecasting. An analogous finding will be reported in section 4.

<sup>11</sup>More broadly, this is extending to trees and ensemble of trees the "classical" non-parametrics literature's knowledge that local linear regression usually has much better properties (especially at the sample boundaries) than the Nadaraya-Watson estimator.

of coefficients can be considered.

## 2.6 Engineering $S_t$

This section discusses principles guiding the composition of  $S_t$ , which is the raw material for  $\mathcal{F}$  in both MRF and plain RF. Macroeconomic data sets (e.g. FRED, [McCracken and Ng 2020](#)) typically contains many regressors and few observations. After incorporating lags for each variable, it can easily be the case that predictors outnumber observations. The curse of dimensionality has both computational and statistical ramifications. The former is mostly avoided in RF since it does not rely on inverting a matrix. However, the statistical curse of dimensionality, a feature of the regressors/observations ratio, remains a difficulty to overcome.

**NO NEED TO CHOOSE.** There are two extreme ways of reducing dimensionality: sparse or dense. The former selects a small number of features out of the large pool in a supervised way (e.g. LASSO), the latter compresses the data in a set of latent factors that should span most of the original regressors space. This is often seen as a necessity to choose *one* of them.<sup>12</sup> However, in a regularized model, both can be included, and we can let the algorithm select an optimal combination of original features and factors. This is useful — it is not hard to imagine a situation where opting for one or the other would prove suboptimal to a more nuanced solution.

To appreciate this point, let us put RF aside for a moment, and look at a high-dimensional linear regression problem. Suppose we define  $S_t = [X_t \ F_t]$  and by construction the factors are some linear combination of original features ( $F_t = X_t R$ ).<sup>13</sup> We can estimate

$$y_{t+1} = X_t \beta + X_t R \gamma + u_t \tag{4}$$

using LASSO. Of course, this would not run with OLS because of perfect collinearity, which is the standard motivation for not mixing dense and sparse approaches. By Frisch-Waugh-Lowell theorem and the factor model

$$X_t = \Lambda F_t + e_t,$$

(4) above is equivalent to

$$y_{t+1} = e_t \beta + F_t \gamma + u_t.$$

At first sight, this has more parameters than either the dense or sparse approach. However, with some adequate penalization of  $\beta$  and  $\gamma$ , the model can balance a proper mix of dense and sparse. For instance, activating some  $\beta$ 's "corrects" the overall prediction when the factor model repre-

---

<sup>12</sup>In macro forecasting work using RF, [Goulet Coulombe et al. \(2019\)](#) follow a dense approach by only including factors while [Borup et al. \(2020a\)](#) opt for sparsity by proposing a Lasso pre-selection step.

<sup>13</sup>Note that in this section only,  $X_t$  denotes generic raw regressors rather than MRF's linear part. This switch allows for the use of familiar-looking notation.

sensation is too restrictive for the effect of a specific regressor  $X_k$  on  $y_{t+1}$ .<sup>14</sup> This representation has been studied in [Hahn et al. \(2013\)](#) and [Hansen and Liao \(2019\)](#) to enhance hard-thresholding methods' performance (like LASSO) in the presence of highly correlated regressors. Coming back to RF, this means its strong regularization/selection allows for both the original data and its rotation to be included in  $S_t$ . This also suggests it is relatively costless to explore alternative rotations of  $X_t$ .

**LAG POLYNOMIALS.** From a predictive standpoint, residuals autocorrelation implies there is forecasting power left on the table. To get rid of it, many lags might be necessary. In multivariate contexts (like that of a VAR), doing so quickly pushes the model to overfit. A standard solution is Bayesian estimation and the use of priors in the line of [Doan et al. \(1984\)](#), which are specially designed for blocks of lags structures. Outside of the VAR paradigm, there is an older literature estimating restricted/regularized lag polynomials in Autoregressive Distributed Lags (ARDL) models ([Almon, 1965](#); [Shiller, 1973](#)). More recently, these methods have found new applications in mixed-frequency models ([Ghysels et al., 2007](#)) where the design of the model leads to an explosion of lag parameters.

(M)RF experiences an analogous situation. A tree may waste many splits trying to efficiently extract information out of a lag polynomial: for instance, splitting on the first lag, then the 7th one, then the 3rd one. In linear parametric models, the above methods can extract the relevant information out of a lag polynomial without sacrificing many degrees of freedom. A significant roadblock to this enterprise in the RF paradigm is that there are no explicit lag polynomials to penalize. An alternative route is to exploit the insight that RF can choose for itself relevant restrictions. We just have to construct regressors that embodies those, and include them in  $S_t$ .

**MOVING AVERAGE FACTORS.** To extract the essential information out of the lag polynomial of a specific variable, a linear transformation can do the job. Consider forming a panel of  $P$  lags of variable  $j$ :

$$X_{t,j}^{1:P} \equiv [X_{t-1,j} \dots X_{t-P,j}] .$$

We want to form weighted averages of the  $P$  lags so that it summarizes most efficiently the temporal information of the feature indexed by  $j$ .<sup>15</sup> The weighted averages with that property will be the first few factors (extracted by PCA) of  $X_{t,j}^{1:P}$ .<sup>16</sup> This can be seen as the time-dimension analog to the traditional cross-sectional factors. The latter are defined such as to maximize their capacity to replicate the cross-sectional distribution of  $X_{t,j}$  fixing  $t$  while the Moving Average Factors (MAFs) proposed here seek to represent the temporal distribution of  $X_{t,j}$  for a fixed  $j$

<sup>14</sup>That problem has been documented in [Bai and Ng \(2008\)](#) and others.

<sup>15</sup> $P$  is a tuning parameter the same way the set of included variables in a standard factor model is one.

<sup>16</sup>While I work directly with the latent factors, a related decomposition called singular spectrum analysis works with the estimate of the *summed* common components. Since this decomposition naturally yields a recursive formula, it has been used to forecast macroeconomic and financial variables ([Hassani et al., 2009, 2013](#)).



in a lower-dimensional space.<sup>17</sup> By doing so, our goal to summarize the information of  $X_{t,j}^{1:P}$  without modifying the RF algorithm (or any other) is achieved: rather than using the numerous lags as regressors, we can use the MAFs which compress information ex-ante. As it is the case for standard factors, MAF are designed to maximize the explained variance in  $X_{t,j}^{1:P}$ , not the fit of the final target. It is the RF part’s job to select the relevant linear combinations among  $S_t$  so to maximize the fit. Finally, it is noteworthy that MAFs facilitate interpretation. As these are moderately sophisticated averages of a single time series, they can be viewed as a smooth index for a specific (but tangible) economic indicator. This is arguably much easier to interpret than a plethora of lags coefficients.

The take-away message from this subsection can be summarized in three points. First, there is no need to choose ex-ante between sparse and dense when the model performs selection/regularization. We can let the algorithm find the optimal balance. Second, to make the inclusion of many lags useful, we need to regularize the lag polynomial. Third, such compression can be achieved most easily by generating MAFs and using those as regressors in RF – or any algorithm.

## 2.7 Quantifying Uncertainty of $\beta_t$ ’s Estimates

Taddy et al. (2015) and Taddy et al. (2016) interpret RF’s prediction as the posterior mean of a tree functional  $\mathcal{T}$  (the splitting algorithm) obtained by an approximate Bayesian bootstrap.<sup>18</sup> Through those lenses, each tree is a posterior draw. Seeing  $\mathcal{T}$  as a Bayesian nonparametric statistic (independently of the DGP) is of even greater interest in the case of MRF.<sup>19</sup> It provides inference for meaningful time-varying parameters  $\beta_t$  rather than an opaque conditional mean function. Such techniques, originating from Ferguson (1973), have seldomly found applications in econometrics, such as Chamberlain and Imbens (2003) for instrumental variable and quantile regressions.

While the Bayesian Bootstrap desirably does not assume many things about the data, it yet makes the assumption that  $Z_t = [y_t \ X_t \ S_t]$  is an *iid* random variable. Thus, it cannot be used directly as a proper theoretical motivation for using the bag of trees directly to conduct inference. I propose a block extension to make Taddy et al. (2015)’s convenient approach amenable to this paper’s setup.

---

<sup>17</sup>In the spirit of the Minnesota prior, one can assign decaying (in  $p$ ) weights to each lag before running PCA. This has the analogous effect of shrinking more heavily the distant lags and less so the recent ones.

<sup>18</sup>The connection between Breiman (1996)’s bagging and Rubin (1981)’s Bayesian Bootstrap was acknowledged earlier in Clyde and Lee (2001).

<sup>19</sup>An alternative (frequentist) inferential approach is that of Friedberg et al. (2018). However, their asymptotic argument requires estimating the linear coefficients and the kernel weights on two different subsamples. This is hard to reconcile with our goal of modeling time-variation and different regimes throughout the entire sample. Furthermore, when the sample size is small, splitting the sample in such a way carries binding limitations on the complexity of the estimated function.

**BLOCK BAYESIAN BOOTSTRAP.** BBB is a conceptual workaround to reconcile time series data with multinomial sampling. First, I briefly review the *standard* Bayesian Bootstrap. Let all the available data be cast in the matrix  $Z_t = [y_t \ X_t \ S_t]$ .  $Z$  is considered as a discrete *iid* random variable with  $T$  support points. Define  $N_t = \sum_{\tau=1}^T I(Z_\tau = z_t)$ , which is the number of occurrences of  $z_t$  in the sample. The goal is to conduct inference on the data weight vector  $\theta_{1:T}$ , and then obtain credible regions for the posterior functional  $\beta_t = \mathcal{T}(\theta_{1:T})$ . To do so, we need to characterize the posterior distribution of vector  $\theta$  (stripped of its subscript for readability)

$$\pi(\theta|\mathbf{z}) = \frac{f(\mathbf{z}|\theta)\pi(\theta)}{\int f(\mathbf{z}|\theta)\pi(\theta)d\theta}.$$

Conditional on  $\theta$ , the likelihood of the data is multinomial. The prior is Dirichlet. Since Dirichlet is the conjugate prior of the multinomial distribution, the posterior is also Dirichlet. That is, it can be shown that combining the likelihood

$$f(\mathbf{z}|\theta) = \frac{N!}{N_1! \cdots N_T!} \prod_{t=1}^T \theta_t^{N_t} \quad \text{with prior distribution} \quad \pi(\theta) = \frac{1}{B(\alpha_{1:T})} \prod_{t=1}^T \theta_t^{N_t + \alpha_t - 1}$$

gives rise to the posterior distribution

$$\pi(\theta|\mathbf{z}) = \frac{1}{B(\bar{\alpha}_{1:T})} \prod_{t=1}^T \theta_t^{N_t + \alpha_t - 1} .$$

where  $\bar{\alpha}_t = \alpha_t + N_t$  and  $B(\bar{\alpha}_{1:T}) = \frac{\prod_{t=1}^T \Gamma(\bar{\alpha}_t)}{\Gamma(\sum_{t=1}^T \bar{\alpha}_t)}$ . Using the uninformative (and improper) prior  $\alpha_t = 0 \ \forall t$ , we can simulate draws from the (proper) posterior using  $\theta_t \sim \text{Exp}(1)$ . The object of scientific interest is typically not  $\theta$  *per se* but rather a functional of it. In [Taddy et al. \(2015\)](#), the functional of interest is a tree and inference is obtained by computing  $\mathcal{T}(\theta_{1:T})$  for each  $\theta_{1:T}$  draw.

BBB is a simple redefinition of  $Z$  so that it is plausibly *iid*. Hence, in the spirit of traditional frequentist block bootstrap ([MacKinnon, 2006](#)), blocks of a well-chosen size will be exchangeable. Thus, a new variable can be defined  $Z_{\mathfrak{b}} \equiv [y_{\mathfrak{b}:\bar{\mathfrak{b}}} \ X_{\mathfrak{b}:\bar{\mathfrak{b}}} \ S_{\mathfrak{b}:\bar{\mathfrak{b}}}]$ . There will be a total of  $\mathfrak{B} = T/\text{block size}$  fixed and non-overlapping blocks. Under covariance stationarity,  $\tilde{Z}_{\mathfrak{b}} = \text{vec}(Z_{\mathfrak{b}})$  are *iid*, for a properly chosen block length.<sup>20</sup> The derivations above can be carried by replacing  $t$  by  $\mathfrak{b}$  and  $T$  by  $\mathfrak{B}$ . Practically, this implies drawing  $\theta_{\mathfrak{b}} \sim \text{Exp}(1)$  which means observations within the same block ( $\mathfrak{b} : \bar{\mathfrak{b}}$ ) share the same weight. As an alternative to this BBB that would also be valid under dependent data, [Cirillo and Muliere \(2013\)](#) provide a more sophisticated urn-based approach with theoretical guarantees. It turns out their approach contains the well-known non-overlapping block bootstrap as a special case, which the above is only its Bayesian rendition.

<sup>20</sup>In practice, I will use block of two years for both quarterly or monthly data.

Analogously to [Taddy et al. \(2015\)](#), block-subsampling is preferred to BBB in implementations since it is faster and gives nearly identical results.

It is reasonable to wonder how the above procedure deals with the possible presence of heteroscedasticity. Fortunately, the nonparametric bootstrap/subsampling that RF uses is in fact the "pairs" bootstrap of [Freedman et al. \(1981\)](#) which is valid under general forms of heteroscedasticity ([MacKinnon, 2006](#)).<sup>21</sup> From a Bayesian point of view, [Lancaster \(2003\)](#) show that the obtained variance for OLS from using such a bootstrap is asymptotically equivalent to that of White's sandwich formula.<sup>22</sup> Hence, in the spirit of heteroscedasticity-robust estimation, no attempt will be made at directly evolving volatility (which is a GLS approach). Rather, it will be reflected in larger bands for periods of smaller signal-to-noise ratio.

### 3 Simulations

Simulations are divided in two parts. The first shows that Autoregressive Random Forest (ARRF) delivers forecasting gains over standard nonlinear time series model when the true DGP mixes both endogenous and exogenous time-variation. Moreover, the former is very resilient against traditional approaches, even when the DGP matches the latter's restrictive assumptions. Additionally, those simulations will numerically document the superiority of ARRF over RF when the AR part is pervasive (as discussed in section 2.5.1). Overall, this helps rationalizing forecasting results from section 4, where ARRF supplants  $\sim$ TARs for the vast majority of targets.

The second simulations section considers simpler linear parts and look at how the algorithm behaves when  $S_t$  is large. Further, I focus on  $\beta_t$  itself and its credible regions. The main point is to visually show that (i) GTVPs adapts nicely to a wide range of DGPs and (ii) are not prone to crying wolf on time-variation.

#### 3.1 Comparison of ARRF to Traditional Nonlinear Autoregressions

I consider 6 DGPs: Autoregression (AR), two Self-Exciting Threshold ARs (SETAR), SETAR with a structural break, AR with a structural break and finally a SETAR model that collapse to an AR (via a structural break). Those DGPs allow to span the space of time-variations I wish to investigate: endogenous, exogenous and both together. Precisely, the DGPs include two types of

---

<sup>21</sup>From a purely predictive point of view, [Grandvalet \(2004\)](#) also stresses the point that bagging provides important improvements when there is "badness" in the data, that is, the presence of uninformative leverage points. Those improvements are shown to be especially meaningful for unstable algorithms such as regression trees.

<sup>22</sup>[Poirier \(2011\)](#) propose better priors and [Karabatsos \(2016\)](#) incorporate such ideas into a generalized ridge regression.

switching variable:  $y_{t-1}$  and  $t$ .<sup>23</sup> For all DGPs,  $X_t = [1 \ y_{t-1} \ y_{t-2}]$ . The simulated series sample size is either  $T = 150$  or  $T = 300$ .<sup>24</sup> The last 40 observations of each sample consist the hold-out sample for evaluation. I forecast 4 different horizons:  $h = 1, 2, 3, 4$ . Models are estimated once at the last available data point.

**MODELS.** SETAR, Rolling-Window (RW) AR, Random Forest (RF) and Autoregressive Random Forest (ARRF) are included. Iterated SETAR forecasts are obtained via the standard bootstrap method (Clements and Smith, 1997) and all the others are generated via direct forecasting. That is, in the latter case, I fit the model directly on  $y_{t+h}$  rather than iterating forward the one-step ahead forecast. To certify that the observed differences between SETAR and other models is not merely due to the choice of iterated vs direct forecasts – a non-trivial choice in many environments (Chevillon, 2007) –, I also include SETAR-d where "d" means its forecasts were alternatively obtained by direct forecasting.

In all simulations, ARRF's  $S_t$  includes 8 lags of  $y_t$  and a time trend, which match what will be referred to in section 4 as "Tiny ARRF". Thus, unlike  $\sim$ TARs, it is "allowed" to split on what we know (by the DGP choices) to be useless regressors (especially at horizon  $h = 1$ ).

**PERFORMANCE METRIC.** Performance is evaluated using the mean squared prediction error (MSPE). In simulation  $s$ , for the forecasted value at time  $t$  made  $h$  steps ahead, I compute

$$RMSE_{h,m} = \sqrt{\frac{1}{40 \times 100} \sum_{s=1}^{100} \sum_{t \in \text{OOS}} (y_t^s - \hat{y}_{t-h}^{s,h,m})^2}.$$

100 different simulations are considered, which means the total number of squared errors being averaged for a given horizon and model is  $100 \times 40 = 4000$ . To provide a visually useful normalization, bar plots report  $RMSE_{h,m}$ 's relative to that of the oracle. Formally, the metric is

$$\Delta_o RMSE_{h,m} = \frac{RMSE_{h,m}}{RMSE_{h,o}} - 1.$$

It is worth specifying what is meant by "oracle". It knows perfectly the law of motion of time-varying parameters  $\beta_t$ . Precisely, if the model has a break and a switching variable, it knows exactly the break points, thresholds and AR parameter values in each regime. The only things the oracle does not know are the future shocks ( $\epsilon_{t+h}$ ), and the out-of-sample evolution of parameters ( $\beta_{t+h}$ ) – unless the latter is purely deterministic.

<sup>23</sup>Since a structural break is just a threshold effect with respect to variable  $t$ , one can conclude without loss of generality that similar results would be obtained using different additional switching variable.

<sup>24</sup> $\sim$ TAR packages in R provide functions to simulate from nonlinear models.

### 3.1.1 No Time-Variation

Given the incredible resilience of AR models in any macroeconomic forecasting exercise, a time-invariant DGP is the inevitable place to start.

**DGP 1: Plain AR(2).** The first DGP being considered is an autoregressive process of order 2

$$y_t = X_t\beta + \epsilon_t, \quad \epsilon_t \sim N(0, 0.25^2)$$
$$\beta = [0.7 \quad -0.2]$$

which results are reported in Figure 1a. As it should be, AR is the best model for all horizons and both sample sizes. The RW-AR suffers from high variance and it is assumed that tuning the window length in a data-driven way would help, but is not the point here. Plain RF struggles, irrespective of the sample size.<sup>25</sup> For the smaller sample, ARRF performs as well as the tightly parametrized SETARs. Their marginal increases in RMSE with respect to the oracle are typically less than 10%, which is small in contrast to simulations yet to come. More observations generally helps AR, the iterated SETAR, and ARRF especially at longer horizons.

### 3.1.2 Endogenous Time-Variation

I now consider cases where parameters vary according to past values of  $y_t$  itself.

**DGP 2: SETAR.** The DGP represents an endogenous switching process which could plausibly suit well real activity variables: it includes high/low regimes, and mildly different dynamics in each of them. In this first SETAR example

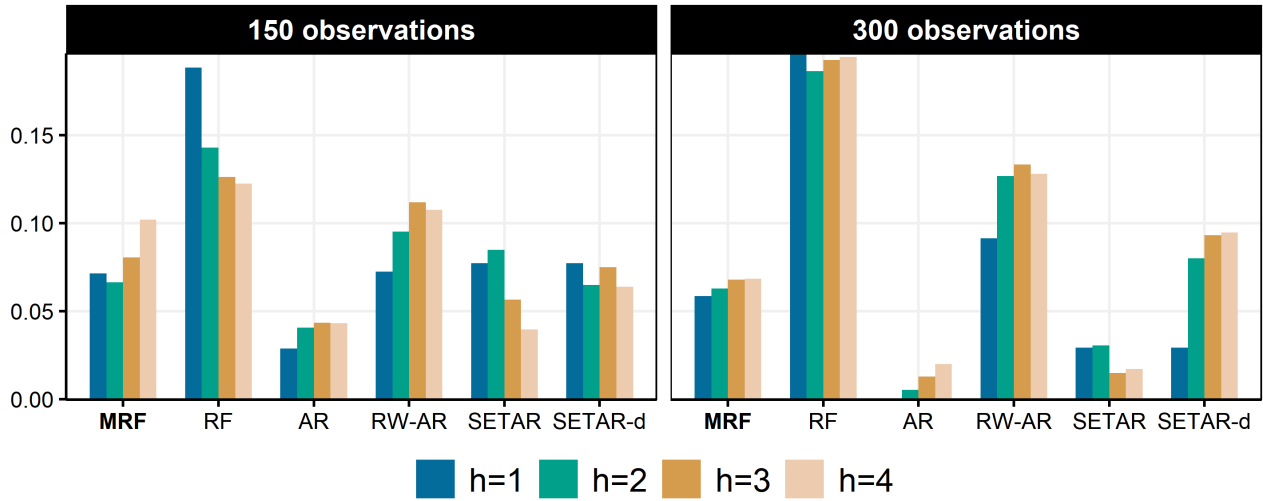
$$y_t = X_t\beta_t + \epsilon_t, \quad \epsilon_t \sim N(0, 0.5^2)$$
$$\beta_t = \begin{cases} [2 \ 0.8 \ -0.2], & \text{if } y_{t-1} \geq 1 \\ [0 \ 0.4 \ -0.2], & \text{otherwise,} \end{cases}$$

AR models are doing badly by not capturing the change in mean and dynamics. It is noteworthy that in this DGP, predictive power quickly vanishes after  $h = 1$ , which is why we observe little performance heterogeneity at longer horizons in Figure 1b: those are dominated by the unshrinkable prediction error.

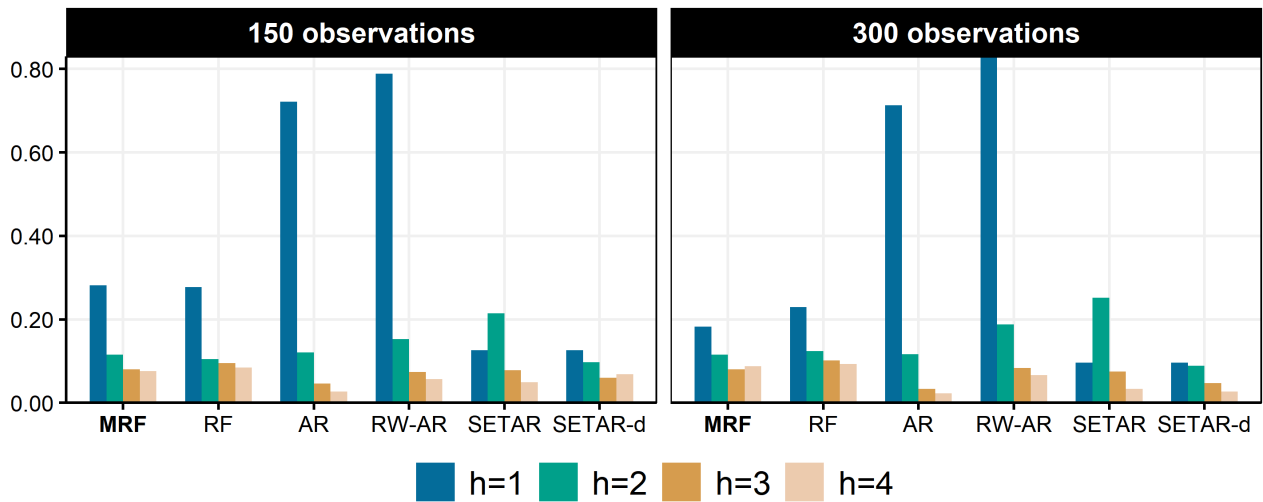
Specifically tailored for this class of DGPs, the two SETARs are offering the best performance. A less trivial observation is that ARRF and RF, while much more general, perform only marginally worse than SETARs. The tie between ARRF and RF is attributable the importance of

---

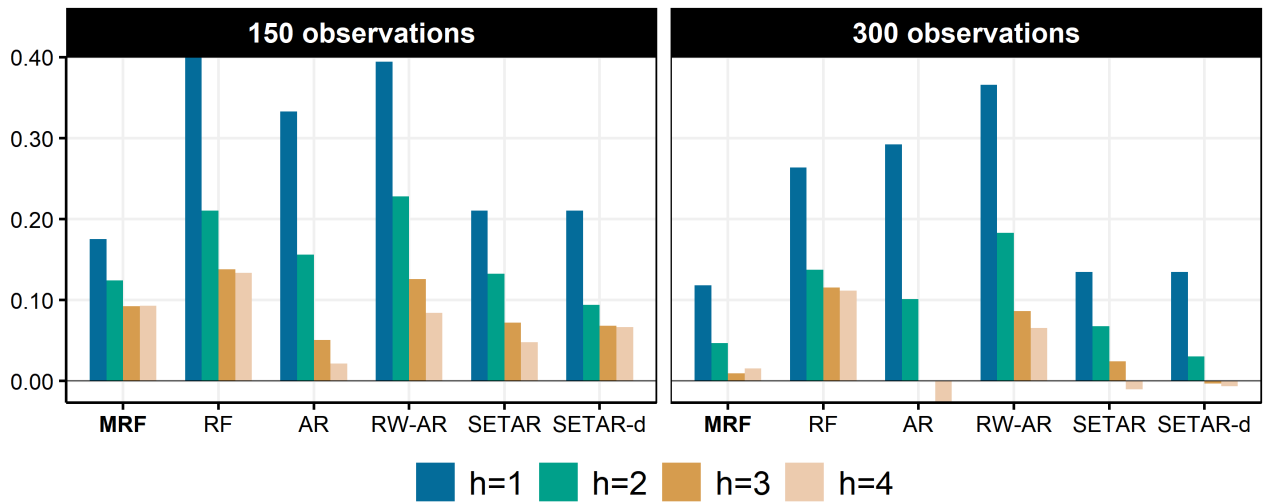
<sup>25</sup>This will be a recurring theme. If the DGP is linear, RF never performs well. The strength of this finding is only magnified when  $X_t$ 's dimension grows, in line with the discussion in section 2.5.1.



(a) DGP is AR(2).



(b) DGP is SETAR.



(c) DGP is Persistent SETAR.

Figure 1: Displayed are increases in relative RMSE with respect to the oracle.

the switching constant in the current DGP, which both models allow for.

**DGP 3: MORE PERSISTENT SETAR.** The increased persistence in

$$y_t = X_t\beta_t + \epsilon_t, \quad \epsilon_t \sim N(0, 0.5^2)$$

$$\beta_t = \begin{cases} [2 \ 0.8 \ -0.2], & \text{if } y_{t-1} \geq 0 \\ [0.25 \ 1.1 \ -0.4], & \text{otherwise} \end{cases}$$

makes results at higher horizons of greater interest in Figure 1c. In the previously considered SETAR, the forecasting ability of the oracle was practically null beyond  $h = 2$ . For all horizons and sample sizes considered, ARRF is practically as good as SETAR, the optimal model in this context. With the increased importance of changing dynamics relative to that of a changing mean, RF is now trailing behind with RW-AR. Nevertheless, the former improves substantially at shorter horizons when the sample size increase. Finally, AR is resilient at longer horizons but is much worse than ARRF and SETAR at shorter ones.

### 3.1.3 Exogenous Time-Variation

I report results for a simple case where  $\beta_t$  varies exogenously – that is, according to time  $t$ .

**DGP 4: AR(2) WITH A BREAK.** Results for

$$y_t = X_t\beta_t + \epsilon_t, \quad \epsilon_t \sim N(0, 0.3^2)$$

$$\beta_t = \begin{cases} [0 \ 0.7 \ -0.35], & \text{if } t < T/2 \\ [0.15 \ 0.6 \ 0], & \text{otherwise} \end{cases}$$

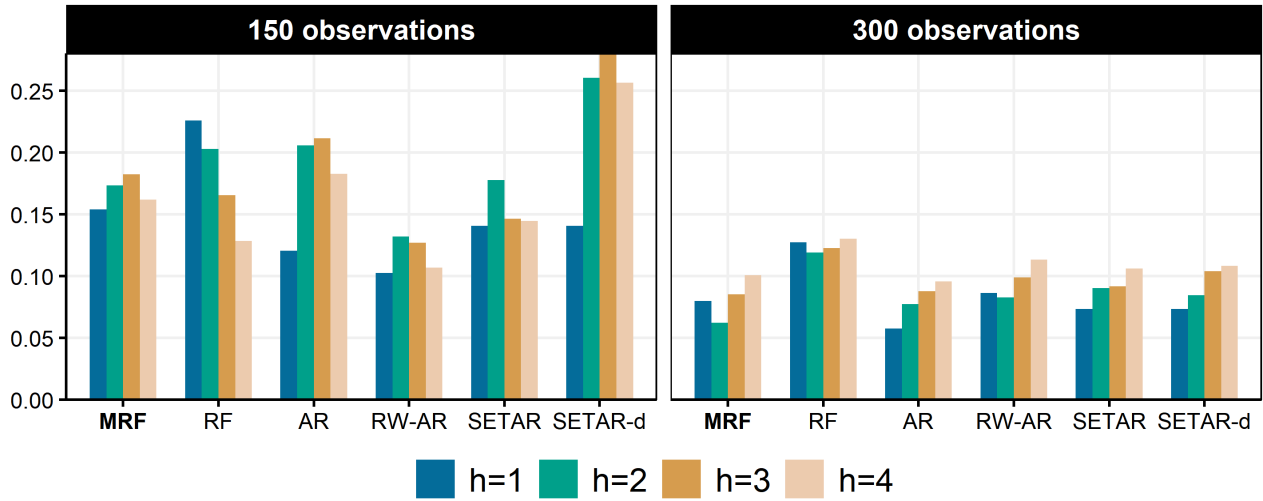
are reported in Figure 1e. In this setup, RW-AR is expected to have an edge, with the estimation window excluding pre-break data. At horizon 1, both RW-AR and ARRF are the best model, beating the robust AR by a thin margin. For  $h > 1$ , ARRF emerges as the best model at both 150 and 300 sample sizes. Naturally, RW-AR is always close behind.<sup>26</sup> As expected, the two models are better than the remaining alternatives by allowing for exogenous structural change (which SETARs and AR do not) and explicitly modeling the autoregressive part (which RF does not).

### 3.1.4 Exogenous & Endogenous Time-Variation

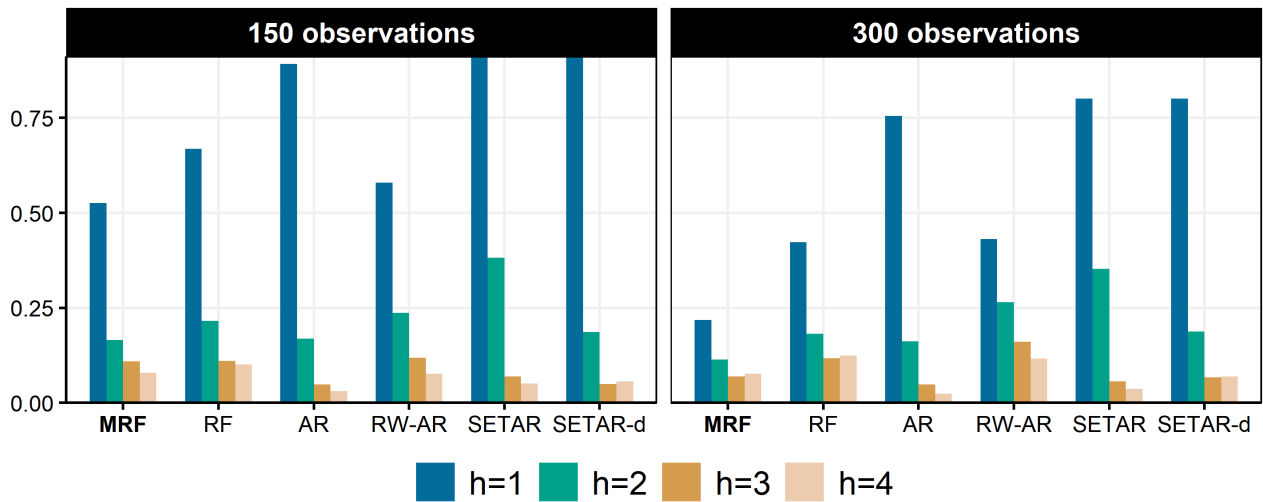
The two sources of time-variation are now combined to display ARRF's edge in these not so implausible situations.

---

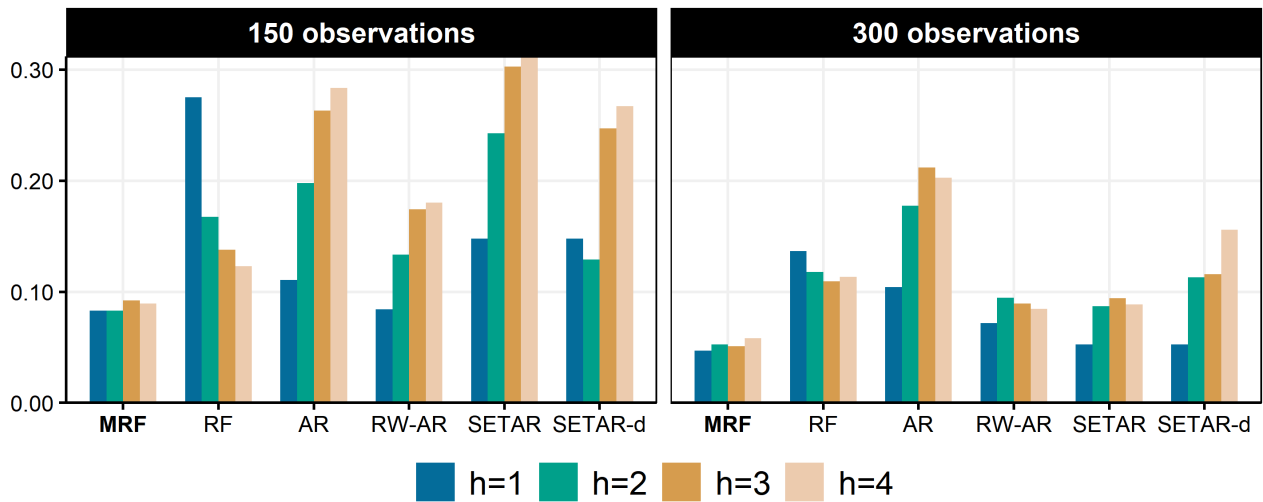
<sup>26</sup>Although not reported here, I considered a simple linear model where I search for a single break (in time) and use the data after the break for forecasting. This option does as well as ARRF for this particular DGP.



(d) DGP is AR(2) with structural break.



(e) DGP is SETAR with a break.



(f) DGP is SETAR morphing into AR(2).

Figure 1: (Continued) Displayed are increases in relative RMSE with respect to the oracle.



**DGP 5: SETAR WITH A STRUCTURAL BREAK.** Until now, I have focused on dynamics that can be captured successfully by currently available time series models. The point of this paper is that most of these models may suffer from serious misspecification issues when estimated on real data. Hence, I introduce here what is possibly the simplest example where structural breaks and switching interact.<sup>27</sup> SETARs are expected to fail because they are not designed to catch breaks. RW-AR is also expected to fail because it does not model switching. RF is general enough, but is anticipated to be inefficient. All these heuristics for

$$\text{DGP 4} = \begin{cases} \text{DGP 2,} & \text{if } t < T/2 \\ \text{DGP 3,} & \text{otherwise} \end{cases}$$

are verified in Figure 1d: ARRF is the better model followed closely by RW-AR and RF for short horizons. With 300 observations, the lead of ARRF, as well as the second position of RF, are both strengthened. At longer horizons, all models perform poorly (including the oracle) due to the fundamental unpredictability of the law of motion for  $\beta_t$ . For these horizons, misspecification only plays a minor role in total forecast error variance, explaining the small and homogeneous decrease in performance with respect to the oracle.

**DGP 6: SETAR MORPHING INSTANTLY INTO AR(2).** The goal of this last DGP is to consider a mixture of endogenous and exogenous time-variation with more interesting results at longer horizons. Further,

$$\text{DGP 6} = \begin{cases} \text{DGP 2,} & \text{if } t < T/2 \\ \text{DGP 1,} & \text{otherwise} \end{cases}$$

can rightfully be hypothesized for some economic time series: complex dynamics up until the mid-1980's followed by a very simple autoregressive structure during the Great Moderation. ARRF comes out as the best model for all horizons in the smaller sample. For horizon 1, RW-AR does equally well, which is expected in this DGP. With respect to the plain exogenous time-variation scenario in Figure 1e, both SETAR and RF's performances have further deteriorated in the smaller sample size.

**ABOUT MISSPECIFICATION IN ARRF.** Most of the reported gains from using ARRF come from avoiding misspecification when a more complex DGP arises. What happens if the arbitrary linear part in ARRF,  $X_t$ , is itself misspecified? Figure 15 in the appendix report corresponding results. For all DGPs under consideration, a "Bad" ARRF, where  $X_t$  is composed of two white

---

<sup>27</sup>Without loss of generality, the threshold according to  $t$  could be replaced by a threshold according to any other variable.

noise series (instead of the first two lags of  $y_t$ ), performs similarly well (or bad) as plain RF.<sup>28</sup>

**SUMMARY.** First, when the true DGP is that of the tightly parametrized classical nonlinear time series model, those perform better than ARRF – but marginally. Second, when it is not the case, the more flexible ARRF does better. Third, when there are pervasive linear autoregressive relationships, plain RF struggles. Fourth, ARRF and RF relative performance both increase with the number of observations but ARRF’s one increases faster if the linear part is well-chosen.

### 3.2 A Look at GTVPs when $S_t$ is Large

A notable difference between the simulations presented up to now and the applied work being carried in later sections is the size of  $S_t$ . In many macro applications, there is no shortage of variables to include in MRF’s  $\mathcal{F}$ . For instance, the FRED-QD data base (McCracken and Ng, 2016) contains over 200 potential predictors that can join lags of  $y$  and a time trend within  $S_t$ . In addition to lessened misspecification concerns, RFs will also benefit from more data through increased randomization (Breiman, 2001). Precisely, it prevents fully-grown trees (or any greedy algorithm) from overfitting (Goulet Coulombe, 2020b).

The additional simulations go as follow. First, I simplify the analysis by looking at a static model with mutually orthogonal but autocorrelated regressors  $X_1$  and  $X_2$ , both driving  $y_t$  according to some process. I simulate each of them for 1000 periods and estimate the models with the first 400 observations. The remaining 600 are used to evaluate the out-of-sample performance. The signal to noise ratio is calibrated to  $2/3$  which is about what is found (out-of-sample) for most models in the empirical section.

The only remaining questions are that of the constitution of  $S_t$  and the generation of  $\beta_t$ ’s. I create two autocorrelated (but not cross-correlated) factors. Out of each of them, I create 50 series with a varying amount of additional white noise.<sup>29</sup> Joining those 100 series with lags of  $y_t$  and a time trend, the final size of  $S_t$  is slightly above 100. Finally,  $\beta_t$ ’s are functions of the underlying *first* factor which (like the second) is not directly included in the data set. In certain DGPs, some  $\beta_t$ ’s will also be a pure function of  $t$  (like random walks, structural breaks).<sup>30</sup> Table 1 summarizes the six DGPs in words. More illustratively, Figure 16 plots one example of each DGPs as well as the estimated GTVPs and their credible region (as discussed in section 2.7). It is visually obvious

---

<sup>28</sup>This result may not hold, however, when the law of motion for the intercept is highly complex and requires a great number of split (unlike what is considered here). This is due to the linear part restricting the depth of trees (with to what plain RF could allow for), especially if observations are scarce. In that regard, increasing the ridge penalty (via  $\lambda$ ) will help. Nevertheless, in practice, it is a safer bet to use a small linear part if uncertainty around its composition is high. More on this and the effect of hyperparameters can found in appendix A.4.

<sup>29</sup>To be precise, their standard deviation is  $U[0.5, 3]\%$  that of the original factor standard deviation.

<sup>30</sup>To clarify, the second factor and underlying series are completely useless to the true DGP – arguably mimicking the inevitable when using a data base of the size of FRED-QD.

that GTVPs are adaptive in the sense that it can discover which kind of time-variation is present in the data while estimating it.

Table 1: Summary of Data-Rich Simulations DGPs

| DGP # | Intercept | $\beta_t^{X_1}$        | $\beta_t^{X_2}$                | Residuals Variance    |
|-------|-----------|------------------------|--------------------------------|-----------------------|
| 1     | Switching | Switching              | Switching                      | Flat                  |
| 2     | Flat      | Switching              | Slow Change (function of $t$ ) | Flat                  |
| 3     | Flat      | Switching              | Structural Break               | Flat                  |
| 4     | Flat      | Latent factor directly | Slow Change (function of $t$ ) | Flat                  |
| 5     | Flat      | Random Walk            | Random Walk                    | Flat                  |
| 6     | Flat      | Flat                   | Flat                           | Stochastic Volatility |

Figure 17 reports distributions of RMSE differentials with respect the oracle (the forecast that knows the  $\beta_t$ 's law of motion). MRF performance is compared to OLS, Rolling-Window OLS (RW-OLS) and plain RF. As expected, MRF outperforms all alternatives by wide margins for most DGPs. By construction, RW-OLS and OLS also perform well for DGP 5 (random walks) and DGP 6 (constant parameters). Nonetheless, it is reassuring to see that MRF either performs much better than OLS or worse by a thin margin (in cases with no time-variation).

## 4 Macroeconomic Forecasting

In this section, I present results for a pseudo-out-of-sample forecasting experiment at the quarterly frequency using the dataset FRED-QD (McCracken and Ng, 2020). The latter is publicly available at the Federal Reserve of St-Louis's web site and contains 248 US macroeconomic and financial aggregates observed from 1960Q1. The forecasting targets are real GDP, Unemployment Rate (UR), CPI Inflation (INF), 1-Year Treasury Constant Maturity Rate (IR) and the difference between 10-year Treasury Constant Maturity rate and Federal funds rate (SPREAD). These series are representative macroeconomic indicators of the US economy which is based on Goulet Coulombe et al. (2019) exercise for many ML models, itself based on Kotchoni et al. (2019) and a whole literature of extensive horse races in the spirit of Stock and Watson (1998b). The series transformations to induce stationarity for predictors are indicated in McCracken and Ng (2020). For forecasting targets, GDP, UR, CPI and IR are considered  $I(1)$  and are first-differenced. For the first two, the natural logarithm is applied before differencing. SPREAD is kept in "levels". Forecasting horizons are 1, 2, 4, 6 and 8 quarters.

The pseudo-out-of-sample period starts in 2003Q1 and ends 2014Q4. I use expanding window estimation from 1961Q3. Models are estimated (and tuned, when applicable) every two years. For all models except SETAR and STAR, I use direct forecasts, meaning that  $\hat{y}_{t+h}$  is obtained by fitting the model directly to  $y_{t+h}$  rather than iterating one-step ahead forecasts.  $\sim$ TAR

Table 2: Composition of  $S_t$ 

| What                                | Why  | How                  |
|-------------------------------------|--|----------------------|
| 8 lags of $y_t$                     | Endogenous SETAR-like dynamics               | –                    |
| $t$                                 | Exogenous "structural" change/breaks         | –                    |
| 2 lags of FRED                      | Fast switching behavior                      | –                    |
| 8 lags of 5 traditional factors $F$ | Compress cross-sectional information ex-ante | Usual PCA            |
| 2 MAFs for each variable $j$        | Compress lag polynomial information ex-ante  | PCA on 8 lags of $j$ |

iterated forecasts are calculated using the block-bootstrap method which is standard in the literature (Clements and Smith, 1997).

Following standard practice, the quality of point forecasts is evaluated using the root Mean Square Prediction Error (MSPE). For the out-of-sample (OOS) forecasted values at time  $t$  of variable  $v$  made  $h$  steps ahead, I compute

$$RMSE_{E_{v,h,m}} = \sqrt{\frac{1}{\#\text{OOS}} \sum_{t \in \text{OOS}} (y_t^v - \hat{y}_{t-h}^{v,h,m})^2}.$$

The standard Diebold and Mariano (2002) (DM) test procedure is used to compare the predictive accuracy of each model against the reference AR(4) model.  $RMSE$  is the most natural loss function given that all models are trained to minimize the squared loss in-sample.

It has been argued in section 2.6 that feature engineering matters crucially when the number of regressors exceeds the sample size.  $S_t$ , the set of variables from which RF can select, is motivated by such concerns. Its exact composition is detailed in Table 2. Among other things, it includes both cross-sectional and moving average factors, which are compressing information along their respective dimensions.

**MODELS.** To better understand where the gains from MRF are coming from, I include models that use different subsets of ideas developed in earlier sections. Those are summarized in Table 3. The competitive data-rich models are in the benchmarks group. Non-linear time series models are also included as they share an obvious familiarity with ARRF. "Tiny" versions of both ARRF and RF are considered to gauge the effect from only having access only to a small  $S_t$  — this could be the case for many non-US applications. Conversely, this helps quantify how a data-rich environment contributes to the success of ARRF versus its plain flexibility. Indeed, Tiny ARRF corresponds to what was shown in the "data-poor" simulations (section 3) to be a generalization of  $\sim$ TARs and related models.

Here are some remarks motivating some inclusions and specifications choices. To assess the marginal effects of MAFs alone, Lasso, Ridge and RF are considered using  $S_t$  — those are known to handle high-dimensional feature space. When it comes to FA-ARRF, I opt for a parsimonious

linear specification including one lag of the first two factors. First, concise models make interpretation easier. Second, considering compact linear specifications within MRF is usually the better strategy. Parameters (including the intercept) are all RFs in their own right and can palliate to the omission of marginally important features, if need be. Consequently, it is desirable to fix a humble linear part and let  $\beta_t$ 's take care of the rest.<sup>31</sup> Finally, as discussed in [McCracken and Ng \(2020\)](#), the first factor mostly loads on real activity variables while the second is a composite of forward-looking indicators like term spreads, permits and inventories. They are baptized and interpreted accordingly.

Table 3: Forecasting Models

| Name                                     | Acronym          | Linear Part ( $X_t^m$ )                                | RF part                |
|--|------------------|--|------------------------|
| Autoregression                           | <b>AR</b>        | $[1, y_{t-\{1:4\}}]$                                   | $\emptyset$            |
| Factor-Augmented Autoregression          | <b>FA-AR</b>     | $[1, y_{t-\{1:4\}}, F_{1,t-\{1:2\}}, F_{2,t-\{1:2\}}]$ | $\emptyset$            |
| Plain Random Forest                      | <b>RF</b>        | $\emptyset$  | Raw data <sup>32</sup> |
| Low-Dimensional Plain RF                 | <b>Tiny RF</b>   | $\emptyset$  | $[y_{t-\{1:8\}}, t]$   |
| Plain RF but using $S_t$                 | <b>RF-MAF</b>    | $\emptyset$  | $S_t$                  |
| RF-MAF on de-correlated $y_t$            | <b>AR+RF</b>     | Filter $y_t$ first with an AR(2), then RF              | $S_t$                  |
| Autoregressive Random Forest             | <b>ARRF</b>      | $[1, y_{t-\{1:2\}}]$                                   | $S_t$                  |
| Low-Dimensional Autoregressive RF        | <b>Tiny ARRF</b> | $[1, y_{t-\{1:2\}}]$                                   | $[y_{t-\{1:8\}}, t]$   |
| Factor-Augmented Autoregressive RF       | <b>FAARRF</b>    | $[1, y_{t-\{1:2\}}, F_{1,t-1}, F_{2,t-1}]$             | $S_t$                  |
| Vector Autoregressive RF <sup>33</sup>   | <b>VARRF</b>     | $[1, y_{t-\{1:2\}}, GDP_{t-1}, IR_{t-1}, INF_{t-1}]$   | $S_t$                  |
| Self-Exciting Threshold AR               | <b>SETAR</b>     | $[1, y_{t-\{1:2\}}]$                                   | $\emptyset$            |
| Smooth Transition AR <sup>34</sup>       | <b>STAR</b>      | $[1, y_{t-\{1:2\}}]$                                   | $\emptyset$            |
| 10 years Rolling-Window AR               | <b>RW-AR</b>     | $[1, y_{t-\{1:2\}}]$                                   | $\emptyset$            |
| Time-Varying Parameters AR <sup>35</sup> | <b>TV-AR</b>     | $[1, y_{t-\{1:2\}}]$                                   | $\emptyset$            |
| LASSO using $S_t$                        | <b>LASSO-MAF</b> | $S_t$  | $\emptyset$            |
| Ridge using $S_t$                        | <b>Ridge-MAF</b> | $S_t$  | $\emptyset$            |

Notes: models are classified in 3 categories: benchmarks, MRFs (and related prototypes), and misc (non-linear time series models, other reasonable additions). The main analysis in section 4.1 omits the 3<sup>rd</sup> club for parsimony.

<sup>31</sup>Further backing a parsimonious choice (with MRF), [McCracken and Ng \(2020\)](#) report that the first two factors account for 30% of the variation in the data while adding two more only bumps it up to 41%, making the last two presumably more disposable in our context.

<sup>32</sup>Precisely, this means 8 lags of FRED-QD, after usual transformations for stationarity have been applied.

<sup>33</sup>Note that the VAR appellation refers to the linear equation consisting of typical "small monetary VAR". The model remains univariate and direct forecasts are used.

<sup>34</sup>The state variable is  $y_{t-1}$ , as in SETAR.

<sup>35</sup>Estimated and tuned via the Ridge approach proposed in [Goulet Coulombe \(2020a\)](#).

## 4.1 Main Quarterly Frequency Results

Violin plots are used throughout to summarize dense RMSEs tables (like Table 4). I report the distribution of  $RMSE_{v,h,m}/RMSE_{v,h,AR}$ . This is informative about the overall ranking and versatility of considered models. Of course, being ranked first does not imply being the best model for every  $h$  and  $v$ . Rather, it means that it performs better on average, over all targets.

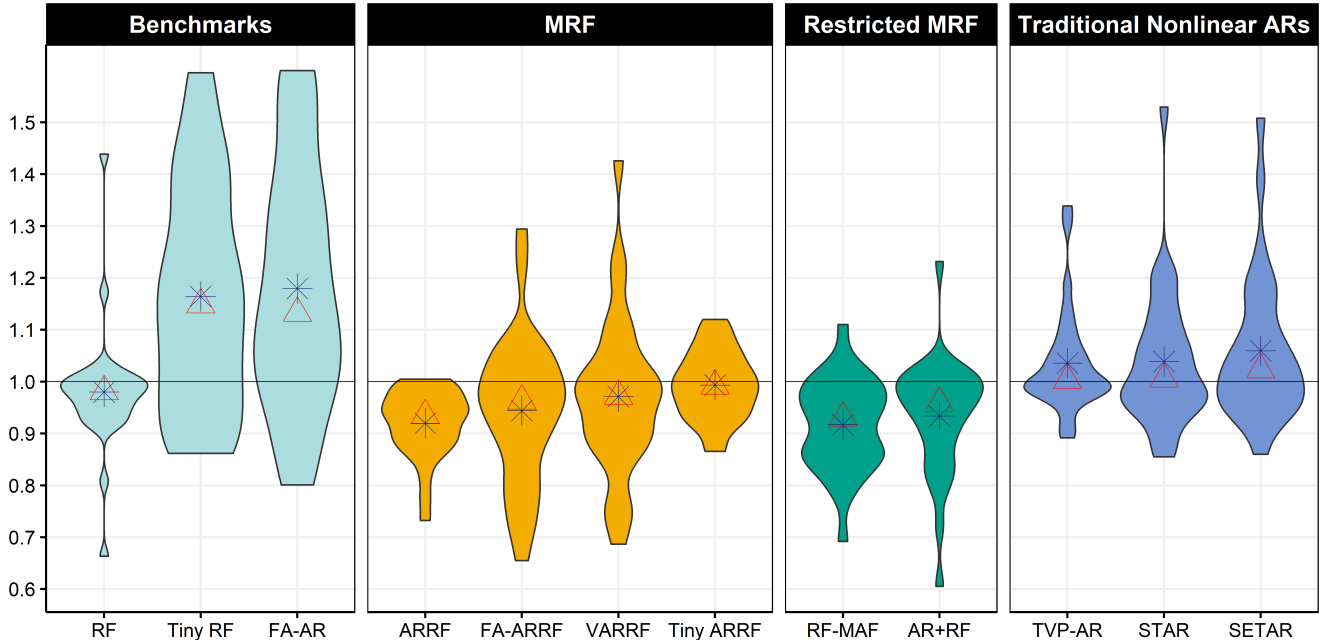


Figure 2: The distribution of  $RMSE_{v,h,m}/RMSE_{v,h,AR}$ . The star is the mean and the triangle is the median.

Here are interesting observations from Figure 2. Clearly, MRFs deliver important gains over both the AR and FAAR benchmarks (the latter is second to last). ARRF has a noticeably small mass above the 1 line. In other words, there are no targets for which ARRF does significantly worse than its OLS counterpart, which makes it atypically adaptable among nonlinear autoregressions. A look at Table 4 confirms this observation also extends to FA-ARRF vs FA-AR. The simplification AR+RF, ranks third with a performance that is much more volatile. This suggests that imposing time-invariant dynamics can sometimes help (see one example in Figure 5), but can also be highly detrimental (as reported for inflation). Of course, that we do not know ex-ante, and it is why AR+RF does not inherit ARRF's "off-the-shelf" quality.

MAFs are useful: RF-MAF does much better than RF which uses the raw data. The latter only exhibits conservative gains over the benchmark. Thus, it is understood that a fraction of MRFs' forecasting gains emanates from considering more sensible transformations of time series data – and which are trivially implementable. The relevance of MAFs is studied more systematically in Appendix A.3 by comparing workhorse high-dimensional models (RF, Lasso, Ridge) with

different information sets.<sup>36</sup>

FAARRF provides very substantial improvements, but can also fail. This is the linear part's doing: FAAR will mostly work well for real activity variables while AR is a jack of all trades. Thus, it is not surprising to see FAARRF inherit some of these uneven properties, albeit to a much milder extent. For instance, in Table 4, FAAR is noticeably worse than AR for all inflation horizons, while FAARRF beats AR for all of them. This phenomenon is well summarized by FAAR being second to last *overall*, well behind FA-ARRF. VARRF has a behavior similar to that of FAARRF, but with less highly noticeable gains.

Does a large  $S_t$  pay off? Most of the time, yes. It is worth re-emphasizing that restricting  $S_t$  restricts the space of time-variations possibilities as well as the potential for trees diversification. Nonetheless, if the restrictions are "true", gains are possible.<sup>37</sup> This is reported to be a rare occurrence, with ARRF  $\succ$  Tiny ARRF (and RF  $\succ$  Tiny RF) for almost any target. Thus, we can safely conclude that a rich  $S_t$  is desirable, with  $\mathcal{F}$  being tasked with selection of relevant items.

As discussed in earlier sections, ARRF connects to the wider family of nonlinear autoregressive models. It clearly does better on average than SETAR and Smooth-Transition TAR. This advantage is attributable to both a more flexible law of motion and a large  $S_t$ . Tiny ARRF is better than the  $\sim$ TAR group, while ARRF is *much* better. Linking this result to those of simulations, this means that no  $\sim$ TAR is likely the true model.

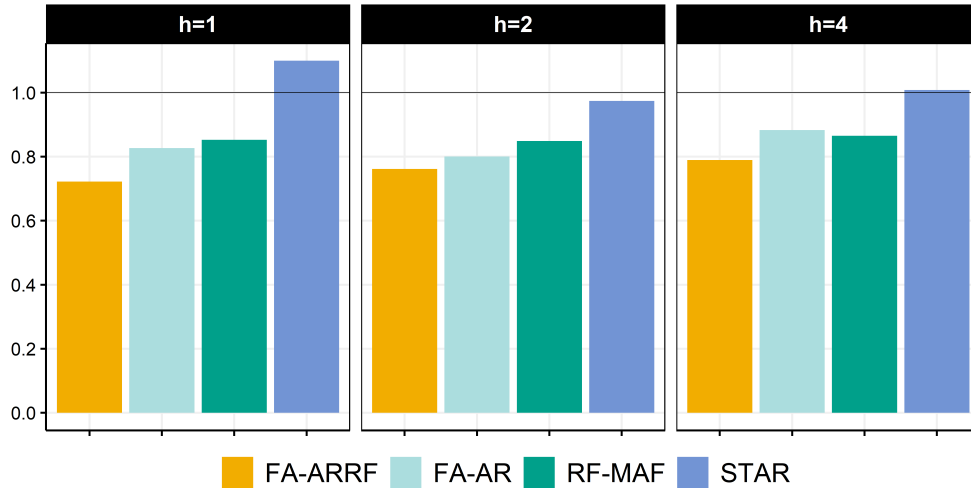
**REAL ACTIVITY TARGETS.** Figure 3 reports results for UR. FAARRF dominates strongly. Table 4 confirms it is the best model for all horizons but the last one (8 quarters ahead, where the encompassed RF-MAF is the best). Clearly, at an horizon of one quarter, the preferred model successfully predicts the drastic rise in unemployment during the Great Recession. Rather than responding with a lag to negative shocks (which is what we observe from AR and ARRF), the model visibly predicts them. As a result, improvements in RMSE are between 25% and 30% over AR for all horizons. Specifically, predicting UR (change) with FAARRF at  $h = 1$  yields an unusually high out-of-sample  $R^2$  of about 80%. The nearly perfect overlap of the yellow and black lines highlight the absence of a one-step ahead shock around 2008. Note that FA-AR and STAR forecasts are omitted from Figure 3b to enhance visibility. STAR forecasts are either similar or worse than the benchmark (as often found for nonlinear time series models). FA-AR forecasts at  $h = 1$  follows a proactive pattern similar to the yellow line, but with a 1 to 2 quarter delay – hence the inferior results.

For  $h = 2$ , the quantitative rise is nowhere near the realized one, but it reveals 6 months ahead the arrival of a significant economic downturn. Additionally, ARRF and FAARRF both flag one year ahead the arrival of a rise in unemployment, which is a quality shared by very few

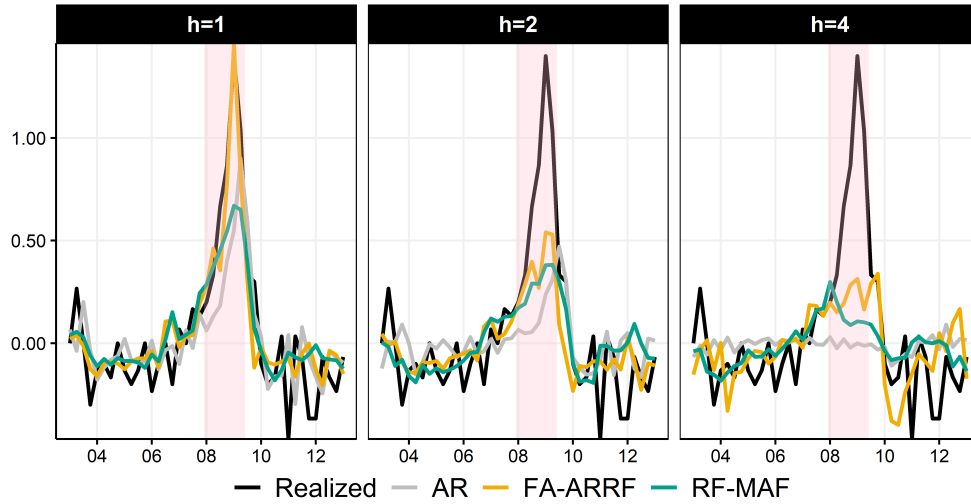
---

<sup>36</sup>Goulet Coulombe et al. (2020a) explores MAFs and more sophisticated derivatives at a larger scale.

<sup>37</sup>An interesting specific case is Tiny ARRF being close behind ARRF for inflation. This is intuitive given that INF has often been associated with exogenous time variation.



(a)  $RMSPE_{UR,h,m} / RMSPE_{UR,h,AR(4)}$



(b) A look at some forecasts

Figure 3: Zooming on best model within each group for UR (change)

models. The barplot in Figure 3 (and Table 4) provides a natural decomposition of FAARRF's gains. Adding the MAFs to an otherwise plain RF procures an improvement of roughly 15% across all horizons (RF-MAF  $\succ$  RF, in Table 4). The linear FAAR part and the rest of algorithmic modifications discussed in section 2 provide an additional reduction of 10% to 15% depending on the forecasted horizon (FAARRF  $\succ$  RF-MAF and FAARRF  $\succ$  FAAR). It is noteworthy that good results for  $h = 1$  are mechanically close to impossible with a plain RF since it cannot extrapolate – i.e., predict values of  $y_t$  that did not occur in-sample. In contrast, this is absolutely feasible within MRF thanks to the linear part.

GDP is known to have a lower signal-to-noise ratio. In Figure 18, FAARRF exhibits a bit less than a 20% drop in RMSE over the AR and nicely grasp the 2008 drop one quarter ahead.<sup>38</sup>

<sup>38</sup>Diebold and Rudebusch (1994) proposed an empirically successful regime-switching factor model. Given that line of work and more recent results in Wochner (2020), the FAARRF's success is not an anomaly.



However, FAARRF performance does not stand apart as much as it did for UR. One reason can be traced visually to predicting higher post-recession growth than its competitors. Finally, RF-MAF closing in on ARRF will be investigated on its own in section 5.2.1. In short, this occurs because once the time-varying intercept is flexibly modeled by RF, there is very little room left for autoregressive behavior (at the quarterly frequency).

**SPREAD AND INFLATION.** VARRF shines for SPREAD (Figure 19) by capturing key movements, even up to a year ahead. The simpler AR+RF also does remarkably well. FAARRF provides successful one-year ahead forecasts. Overall, these results highlight the common importance of the autoregressive part, which is no surprise given SPREAD’s persistence. For INF, Table 4 displays that RF-MAF is the leading model (with ARRF close behind) reducing RMSEs by 12-15% for all horizons. I investigate this with GTVPs in section 5.2.1.

In appendix A.1, I demonstrate that larger VAR linear parts are possible, and sometimes quite helpful. By increasing the strength of MRF’s three main regularizers ( $\lambda, \zeta$ , and the minimal leaf size), impressive results are reported for a VARRF-ALL which linear part includes *all of* FRED (i.e, over 200 regressors). In Figure 13d, UR’s forecast implies the detection of a recession a year ahead.

## 4.2 Monthly Frequency Results

I run a similar exercise as in Goulet Coulombe et al. (2019) which is very close to what has been precedently conducted for quarterly data. FRED-MD is now used. It contains 134 monthly US macroeconomic and financial indicators observed from 1960M01 (McCracken and Ng, 2016). To match the experimental design of Goulet Coulombe et al. (2019) for ML methods, Industrial Production (IP) replaces GDP and IR is dropped. The horizons of interest are  $h = 1, 3, 9, 12, 24$  months. The forecast target is the average growth rate  $\sum_{h'}^h y_{t+h'}^v / h$  which is much less noisy than the monthly growth rate. For example, for inflation 24 months ahead, I target the average inflation rate over the next two years – rather than the monthly inflation rate in 2 years. The OOS period is the same as before.

In Figure 4, VARRF is now doing much better on average, ranking first in terms of mean improvement over AR. ARRF still provides great insurance against doing worse than a plain AR counterpart (here AR(12)).<sup>39</sup> FAARRF remains very competitive. The models that do not have the MAFs (benchmarks) are clearly outperformed by the rest that do. This unsurprisingly indicates that lag polynomial compression can be of even greater use at the monthly frequency.

Table 5 reports specific  $RMSE_{v,h,m} / RMSE_{v,h,AR}$ ’s with Diebold-Mariano tests. Broadly, they show that (i) MAFs are without any doubt the major improvement for the first three variables (IP,

<sup>39</sup>This is also true for the more parsimonious AR, see Table 5.

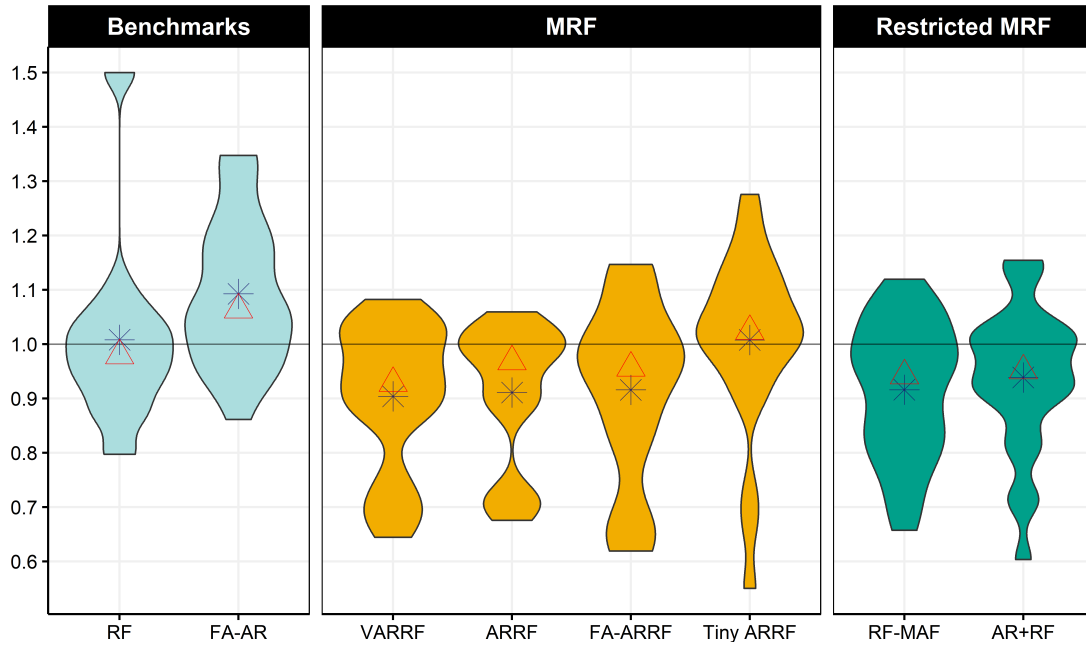


Figure 4: The distribution of  $RMSE_{v,h,m}/RMSE_{v,h,AR}$  for monthly data. The star is the mean and the triangle is the median.

UR, SPREAD), (ii) simpler approaches like RF-MAF and AR+RF do well (*except* for INF) (iii) *all* MRFs do very well for inflation. Particularly, for (iii), ARRF and Tiny ARRF provide significant gains of 33% and 45% over the benchmark at  $h = 12$  and  $h = 24$ , respectively. It is clear from this evidence, and that of the quarterly section, that forcing time-invariant inflation dynamics is costly in terms of RMSPE. GTVPs will confirm that, in accord with classic evidence on the matter (Cogley and Sargent, 2001).

Gains for INF are miles ahead from the usual competition. Table 5 includes forecasts inspired by the contribution of Atkeson et al. (2001): 1,  $h$  and 12 months moving averages are considered (where  $h$  is the targeted horizon). As in the original paper, the "AO-12" forecasts prove remarkably resilient, but are bested with sizable margins at each horizons by ARRF, Tiny ARRF, and FAARRF. For instance, at  $h = 24$ , the next best non-MRF forecast delivers 16% gains over the benchmark AR, whereas the worst MRF provides a gain of 27%. Tiny ARRF supremacy at longer horizons is sensible given that restricting  $S_t$  emphasizes long-run exogenous change, a usual suspect for INF.

Another interesting observation emerges from MRFs successes with monthly inflation. FAARRF is often close to the best model, and that, at all horizons. Naturally, this is intriguing as FAARRF can be thought of as a Phillips' curve forecast, which recurrent failures are well documented (Atkeson et al., 2001; Stock and Watson, 2007). Moreover, it is reported that FAAR, in contrast, does really bad. To sort this out, FAARRF's underlying GTVPs are studied in section 5.3.2.

### 4.3 External Validity

Much attention has been paid to the prediction of US economic aggregates. An even greater challenge is that of forecasting the future state of a small open economy. Such an application is beyond the scope of this paper but is considered in [Goulet Coulombe et al. \(2020b\)](#). The study considers the prediction of more than a dozen key economic variables for Canada and Québec using the large Canadian data base of [Fortin-Gagnon et al. \(2018\)](#). Forecasts from about 50 models and different averages of them are compared, with ARRF and FAARRF among them. MRFs generate substantial improvements especially at the one-quarter horizon for numerous real activity variables (Canadian GDP, Québec GDP, industrial production, real investment). In such cases, ARRF or FAARRF provide reductions (with respect to autoregressive benchmark) that are sizable and statistically significant, going up to 32% in RMSE. That performance is sometimes miles ahead from the next best option (among Complete Subset Regression, Factor models, Neural Networks, Ridge, Lasso, plain RF and different model averaging schemes). [Goulet Coulombe et al. \(2020b\)](#)'s results suggest that MRFs forecasting abilities generalize beyond the traditional exercise of predicting US aggregates.

## 5 Analysis

Based on forecasting results, I concentrate on FAARRF's GTVPs. Additionally, its parameters are easier to interpret (given factors are labeled) since regressors are mechanically orthogonal. First, I look at  $\beta_t$  and analyze their behavior around the Great Recession. Second, I compare GTVPs to random walk TVPs, ex-post vs ex-ante, with a focus on the recessionary episode. Finally, I use a surrogate model approach to explain of the parameters' paths in terms of observed variables.

### 5.1 Forecast Anatomy

$\beta_t$ 's characterize completely MRF's forecasts. Thus, we can investigate GTVPs to understand results from the previous section. The FAARRF forecasting equation is

$$y_{t+h} = \mu_t + \phi_{1,t}y_t + \phi_{2,t}y_{t-1} + \gamma_{1,t}F_{1,t} + \gamma_{2,t}F_{2,t} + u_{t+h}.$$

and naturally  $\beta_t = [\mu_t \ \phi_{1,t} \ \phi_{2,t} \ \gamma_{1,t} \ \gamma_{2,t}]$ . To avoid overfitting,  $\hat{\beta}_t$ 's are (as in section 3.2) the mean over draws that did not include observations  $t - 4$  to  $t + 4$  (a two-year block) in the tree-fitting process. Intuitively, this mimics in-sample the real out-of-sample experiment that starts here in 2007Q2.<sup>40</sup>

---

<sup>40</sup>Note that this is partially different from what gave the results reported in section 4.1, where the model was re-estimated every 2 years. Here, estimation occurs once in 2007Q2.

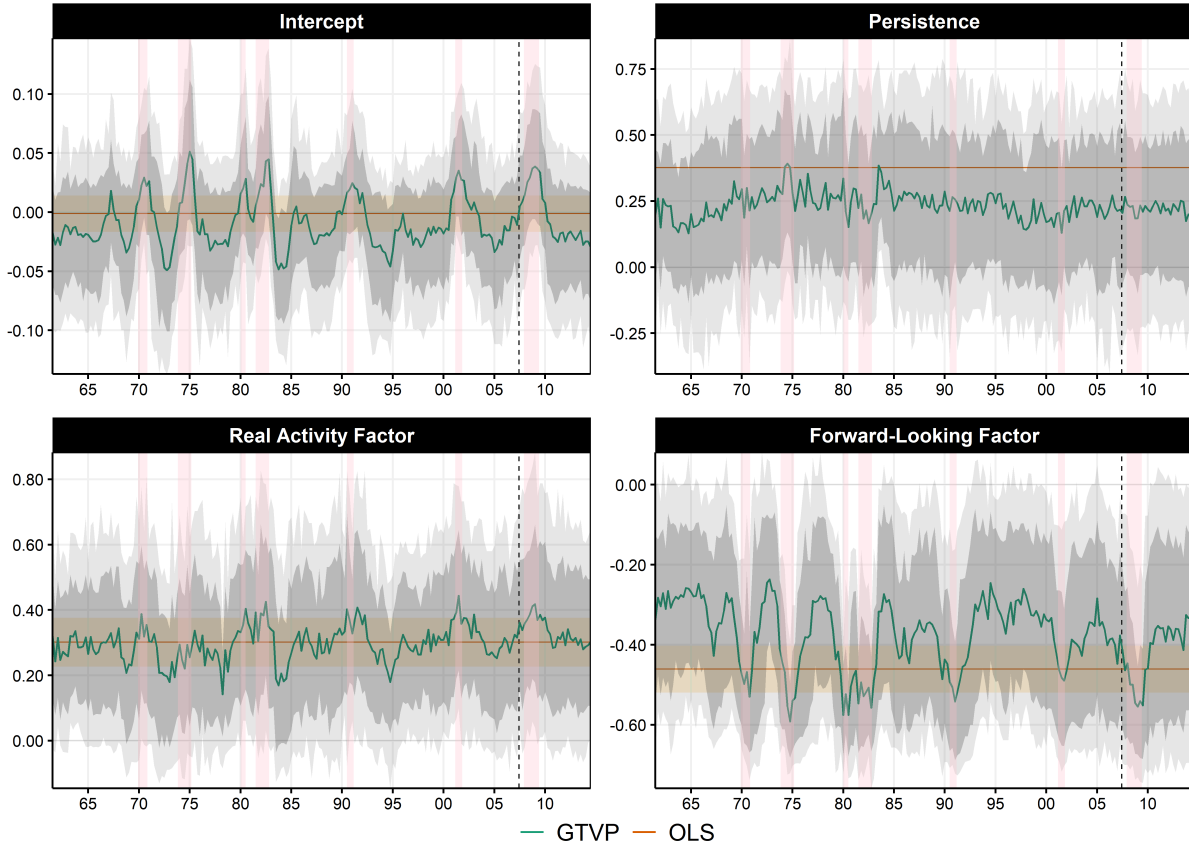


Figure 5: GTVPs of the one quarter ahead UR forecast. Persistence is defined as  $\phi_{1,t} + \phi_{2,t}$ . The gray bands are the 68% and 90% credible region. The pale orange region is the OLS coefficient  $\pm$  one standard error. The vertical dotted blue line is the end of the training sample. Pink shading corresponds to NBER recessions.

Figure 5 displays GTVPs underlying the successful one-step ahead UR *change* forecast. The intercept clearly alternates between at least two regimes and the "increasing UR" one is in effect circa 2008. In levels, this translates to UR alternating between a positive and negative (albeit small) trend. Overall persistence is strikingly time-invariant, and marginally smaller than for OLS estimates. The effect of  $F_1$ , the real activity factor, is generally within OLS confidence intervals, suggesting that while  $\gamma_{1,t}$  almost doubles around recessions, this is subject to great uncertainty.

What is less uncertain, however, is the magnified contribution of the forward-looking factor  $F_2$  during recessionary episodes, which stands out as the key difference with OLS.  $\gamma_{2,t}$  smooth-switching behavior can be best interpreted by remembering that  $F_2$  is highly correlated with capacity utilization, manufacturing sector indicators, building permits and financial indicators (like spreads) (McCracken and Ng, 2020). Many of those variables are considered "leading" indicators and have often been found to increase forecasting performance, mostly before and during recession periods (Stock and Watson, 1989; Estrella and Mishkin, 1998; Leamer, 2007).

Recently, there has been renewed attention on the matter, with financial indicators highlighted as capable of capturing economic activity downside risk (Adrian et al., 2019; Delle Monache et al., 2020). This brand of nonlinearity can translate to a more active  $\gamma_{2,t}$  around business cycle turning points. MRF learns that, while OLS provides a clumsy average of two regimes. In Figure 5, the obvious consequence of OLS' rigidity is being over-responsive to leading indicators during tranquil economic times, and under-responsive when it matters.

Section 5.3 will investigate formally the underlying variables driving this time variation. Figure 21 displays equivalent  $\beta_t$  for GDP one quarter ahead. The pattern  $\gamma_{2,t}$  is also visible for GDP, but it is quantitatively weaker and more uncertain – which is no surprise given GDP being generally noisier than UR. Additionally, slow and relatively mild long-run change is observed. Interestingly,  $\gamma_{1,t}$  has been shrinking since the mid 1980s, and its regime dependence exhibited in the first four recessions is no more.

## 5.2 Comparing Generalized TVPs with Random Walk TVPs

The relationship between random walk TVPs and GTVPs was evoked earlier. I compare them for the small factor model. I estimate standard TVPs using the ridge regression technique developed Goulet Coulombe (2020a). Conveniently, the procedure incorporates a cross-validation step that determines the optimal level of time variation in the random walks.<sup>41</sup>

As Figure 5 suggested for  $\mu_t$  and  $\gamma_{2,t}$ , parameters can be subject to recurrent, rapid and statistically meaningful shifts. Such behavior creates difficulties for random-walk TVPs, which put the accent on smooth and slow structural change. Figure 6 confirms this conjecture. Standard TVPs look for long-run change when regime-switching behavior is the main driving force. As a result, they are flat and within OLS confidence bands, as often reported in the literature (D'Agostino et al., 2013). Of course, more action will mechanically be obtained for TVPs when considering a smaller amount of smoothness than what cross-validation proposed. In appendix A.5, I report the same figures, but using the optimal smoothing parameters (as picked by CV) divided by 1000. This provides much more volatile random walk TVPs that are inclined, at certain specific moments, to follow the GTVPs. However, it is clear in Figure 6 that the end-of-sample/revision problem is worsen by the forced lack of smoothing.

It is known in the traditional TVP literature that there is a balance between flexible (but often erratic)  $\beta_t$  paths and very smooth ones where time variation may simply vanish.<sup>42</sup> Since random-walk TVPs are unfit for many forms of the time-variation present in macroeconomic data, high bias estimates are usually reported as only they can keep variance at a manageable level. This

<sup>41</sup>I show with simulations that this much easier approach performs similarly well (and sometimes better) to traditional Bayesian TVP-VAR, for model sizes that the latter is able to estimate.

<sup>42</sup>In the case of ridge regression-based TVPs, cross-validation is just a data-driven way of backing this necessary empirical choice.



Figure 6: UR equation  $\beta_t$ 's obtained with different techniques. Persistence is defined as  $\phi_{1,t} + \phi_{2,t}$ . TVPs estimated with a ridge regression as in Goulet Coulombe (2020a) and the parameter volatility is tuned with k-fold cross-validation — see Figure 26a for a case where TVP parameter volatility is forced to be higher. Ex Post TVP means using the full sample for estimation and tuning as opposed to only using pre-2002 data as for GTVPs. The pale orange region is the OLS coefficient  $\pm$  one standard error. Pink shading corresponds to NBER recessions.

can have serious implications. Relying too much on time-smoothing can create a mirage of long-run change and/or dissimulate parameters that mostly (but not solely) vary according to expansions/recessions.

Another concern, particularly consequential to the act of forecasting, is the boundary problem. As discussed earlier, random-walk TVP models forecasts can suffer greatly from it because by construction, forecasts are always made at the boundary of the variable on which the kernel is based – i.e., time. One can deploy a 1-sided kernel, but this only alleviate a few pressing symptoms without attacking the heart of the problem. In sharp contrast, GTVPs use a large information set  $S_t$  to create the kernel, which implies that the likelihood of making a forecast at the boundary is rather low, unless the RF part constantly selects  $t$  as splitting variable.

Figures 6 and 22 show, for both random walk and generalized TVPs, their full-sample versions (up to the end of 2014, "ex post") and their version with a training sample ending in 2007Q2 (the dashed blue line). There are two main observations. First, GTVPs are much less prompt

to rewrite recent history than random-walk TVPs. Indeed, the green line and the magenta one closely follows each other all the way up to the end of the training sample. Second, while GTVPs can change many quarters after 2007Q2 (like the GDP constant), they are generally very close to each other at the boundary – especially when the time variation is statistically meaningful (like that of  $\mu_t$  and  $\gamma_{F_2,t}$ ), which is what matters for forecasting. This is much less true of random walk TVPs as there are clear examples where the two version differ for a long period of time (for instance, the intercept and the coefficient on  $F_2$  in the GDP equation), and this often culminates at the boundary.<sup>43</sup>

### 5.2.1 Why and When MRF Can Fail to Deliver Better Forecasts

MRF can sometimes be outperformed by simpler alternatives, like standard RF that incorporate MAFs. When that occurs, it is usually due to the inadequacy of the linear part rather than GTVPs themselves. Unlike traditional TVPs, GTVPs rarely provides a model worse than OLS.

Trivially,  $\beta_t$  helps understanding relative performance. For instance, in the case of forecasting inflation with the *quarterly* data set, ARRF does not supplant RF-MAF. The critical difference between ARRF (reported in Figure 7a) and its restricted analog is that the two autoregressive coefficients of the former are shut to 0.<sup>44</sup> In Figure 7a, the estimates of ARRF broadly agree with the view that inflation persistence has substantially decreased during and following Volker disinflation (Cogley and Sargent, 2001; Cogley et al., 2010).

In terms of anticipated forecasting performance, such decline in persistence suggests a constrained version simply including  $\mu_t$  may do better. The OOS evaluation period corresponds to the region of Figure 7a where  $\phi_{1,t} + \phi_{2,t}$  is the nearest to 0. Given that observation, RF-MAF mildly improving upon ARRF is less surprising. An analogous finding emerges for GDP at many horizons. ARRF does not outperform RF-MAF like FAARRF and larger VARs versions of MRF do. GTVPs showcased in Figure 7b provide a simple explanation. There is only a limited role for persistence when allowing for a forest-driven  $\mu_t$ .  $\phi_{1,t} + \phi_{2,t}$  is below the OLS counterpart most of the time and the credible 68% credible region frequently includes 0. The ensuing forecast is essentially a time-varying constant, which is what RF-MAF does. In sum, unlike many ML offerings, MRF successes and failures can be understood via a time-varying parameter interpretation. The helpfulness of this attribute cannot be overstated when thinking about future model improvements.

---

<sup>43</sup>In (real) practice, all models would be re-estimated each quarter. However, it is worth pointing out that re-estimating every period is much more important for random-walk TVP than it is for GTVPs. For such reasons, the TV-AR in section 4 was the sole model estimated every period rather than every two years.

<sup>44</sup>Of course, lags of INF can still enter the forest part for  $\mu_t$ , so RF-MAF does not suppress entirely the link between current and recent inflation.

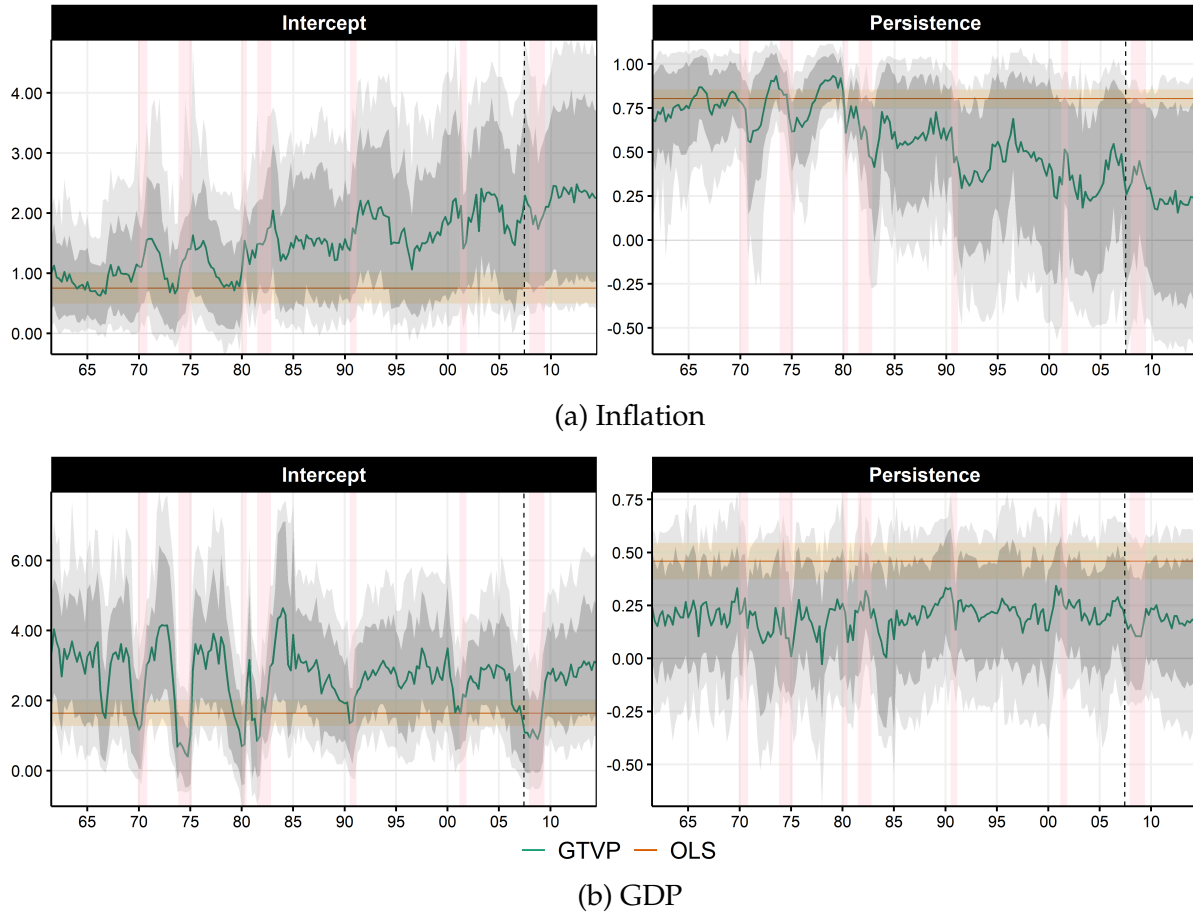


Figure 7: GTVPs of the one-quarter ahead forecasts using ARRF. Persistence is defined as  $\phi_{1,t} + \phi_{2,t}$ . The gray bands are the 68% and 90% credible regions. The pale orange region is the OLS coefficient  $\pm$  one standard error. The vertical dotted line is the end of the training sample. Pink shading corresponds to NBER recessions.

### 5.3 Cutting Down the Forest, One Tree at a Time

Evolving  $\beta_t$  can limit macroeconomists in their ability to use the model for counterfactuals. Complementarily, policy-makers will complain about the limited use for a model in which tomorrow's parameters are unknown (random walks). Fortunately, GTVPs may be the result of an opaque ensemble of trees, but they are made out of observables rather than a multiplicity of latent states. That is, they change, but according to a *fixed* structure. Hence, the reduced-form coefficients could easily change, and yet remain completely predetermined as long as  $\mathcal{F}$  itself is stable. In this paradigm, a changing  $\beta_t$  is not necessarily empirical evidence supporting Lucas (1976)'s critique – rather, a changing  $\mathcal{F}$  could be. Hence, dissecting  $\mathcal{F}$  is inherently interesting. One way to get started on this is to use well-established measures of Variable Importance (VI), originally proposed in Breiman (2001). Those extract features driving the *prediction*. Conveniently, they can be adapted to inquire  $\beta_t$ . Then, one can capitalize on VI's insights to build interpretable small trees parsimoniously approximating  $\beta_{t,k}$ 's path.



The construction of upcoming graphs consists in two steps. I start by computing 3 different VI measures:  $VI_{OOB}$  (out-of-bag predictive performance),  $VI_{OOS}$  (out-of-sample predictive performance) or  $VI_{\beta}$  (for a specific coefficient rather than the whole prediction). Appendix A.2 contains a detailed explanation those and a discussion on how the current approach relates to recent work in the ML interpretability literature. As a potential data set for the construction of a surrogate tree, I consider the union of the 20 most potent predictors as highlighted by any of the three VIs. The tree is pruned with a cost-complexity factor (usually referred to as  $c_p$ ) of 0.075. That tuning parameter is set such as to balance its capacity to mimic the original GTVP and potential for interpretation.

### 5.3.1 Unemployment Equation

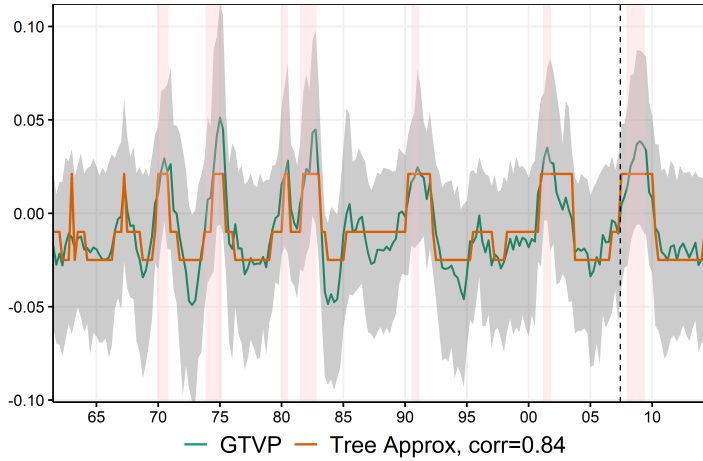
I limit the attention  $\mu_t$  and  $\gamma_{2,t}$  paths, which were argued of greater importance to FAARRF's success in forecasting UR. Also, the nature of their variation is easier to characterize with a single tree (ex-post). Figures 8b and 8d show that paths can sometimes be summarized succinctly using a handful of predictors.

Most of  $\mu_t$  can be captured by two states which are determined by a cut-off on total private sector employees (USPRIV): 0.021 (increasing unemployment) and -0.018 (decreasing). This first layer basically classifies recession vs expansions in a very parsimonious way, which is inevitably crude and imperfect. The additional split on a MAF of non-financial leverage provides a more refined classification: there are more or less three states. The time series plot shows the alternation between two symmetrically opposed states of 0.021 and -0.025 (respectively entering and exiting a recession) and a transitory (and seldomly visited) middle ground around 0.

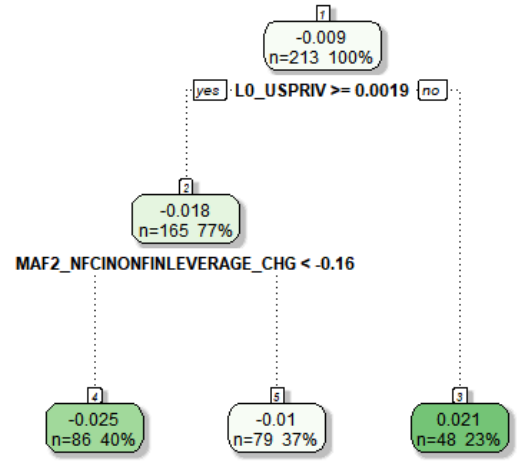
The impact of  $F_2$  on UR switches significantly, and most of the action can be summarized by a private sector employees dummy (USPRIV). The indicator's movement downwards – which usually commence from the *onset* of a recession – can double the effect of  $F_2$  on UR in absolute terms. However, some high (absolute)  $\gamma_{2,t}$  episodes would be left behind when merely using USPRIV. Those are retrieved by an additional split with a MAF of average corporate bonds yield with a BAA rating (lower medium grade).

The GTVP (green line) often plunges earlier than the ex-post surrogate tree's replica (orange). This is important, especially from a forecasting perspective. In Figure 23b, it is clear that leading indicators (especially financial ones) play a prominent role in driving the GTVP  $\gamma_{2,t}$  – well before USPRIV starts showing signs of an imminent downturn. Since  $F_2$  is already composed mostly of forward-looking variables, this hints at a convex effect of market-based expectations proxies.

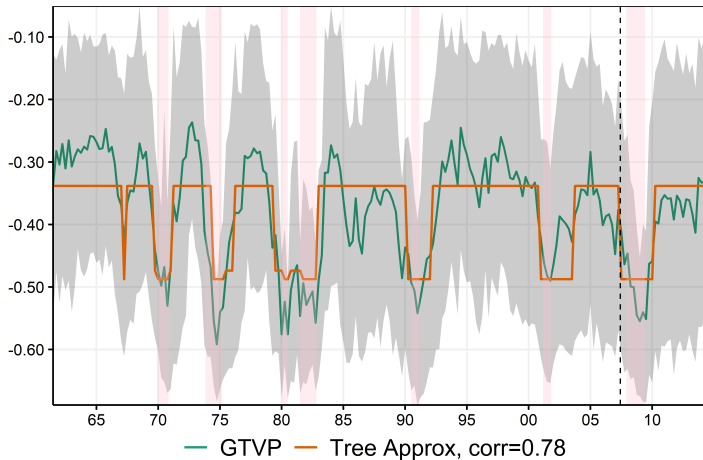
Lastly, a word of caution. Given the points raised earlier in section 2.1, it is more appropriate to see these surrogate trees as suggestive of one potential explanation. It is an open secret that their exact structure is sensitive to small changes in the estimated path. For instance, little



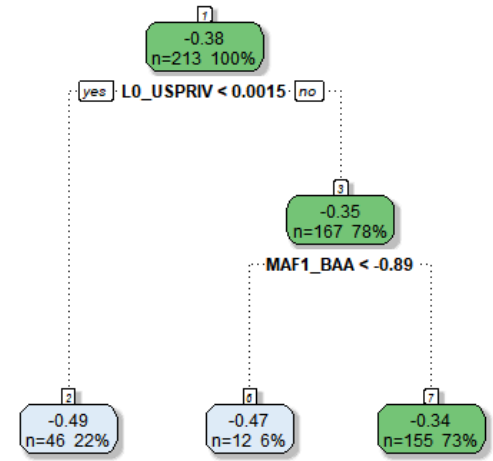
(a)  $\mu_t^{UR,h=1}$ : Surrogate Model Replication



(b)  $\mu_t^{UR,h=1}$ : Corresponding Tree



(c)  $\gamma_{t,F2}^{UR,h=1}$ : Surrogate Model Replication



(d)  $\gamma_{t,F2}^{UR,h=1}$ : Corresponding Tree

Figure 8: Surrogate  $\beta_{t,k}$  Trees. Shade is 68% credible region. Pink shading is NBER recessions.

variation in  $\beta_t$  is needed to observe a change in the exact choice of variables itself. As a result, some of them may rightfully seem exotic when singled out in such a simple tree. GTVPs, as the product of a forest, will more often than not rely on a multitude of indicators from a specific group (which we observe in Figure 23a) rather than a single indicator.

### 5.3.2 Monthly Inflation Equation

As discussed briefly earlier, FAARRF is a very competitive model for *monthly* inflation at all horizons. By its use of  $F_1$ , the real activity factor, it has the familiar flavor of a Phillips' curve (PC).<sup>45</sup> This is of interest given PCs have at best a very uneven forecasting track record (Atkeson et al.,

<sup>45</sup>As noted in Stock and Watson (2008), the plethora of output gap indicators used in literature makes the use of a common statistical factor a credible alternative.

2001; Stock and Watson, 2008; Faust and Wright, 2013). For instance, simple autoregressive/random walk/historical mean benchmarks often do much better.

Given its paramount importance within New Keynesian models, many explanations have been proposed for PC forecasts failures. The curve could be time-varying in a way that annihilates its forecasting potential (Stock and Watson, 2008). Closely related, some have stipulated the PC is nonlinear (Dolado et al., 2005; Doser et al., 2017; Lindé and Trabandt, 2019; Mineyama, 2020). If that were to be true, this should be exploitable. Lastly, an adjacent point of view, which became increasing popular following the Great Recession, is that the PC has irreversibly flattened to the point of predictive desuetude (Blanchard et al., 2015; Blanchard, 2016; Del Negro et al., 2020). Unlike the first two propositions, this one is, by nature, terminal.

Of course, all those explanations amount to hypotheses on the nature of  $\gamma_{1,t}$ 's time variation, of which MRF provides a very flexible account. It is worth emphasizing that MRF is estimated up to 2007Q2, unlike many of the above models explaining the "missing disinflation" after observing that it took place.<sup>46</sup> The variable importance measures reported in Figure 24 showcase a "consensus" subset of variables that matters for inflation time variation. Three popular ones are the trend, MAF of building permits and MAF of housing starts. The leading role for the trend suggests that exogenous time variation is important to explain inflation – to no one's surprise (Cogley and Sargent, 2001). Studying  $\beta_t$ -specific VI's suggest that this is mostly a feature of the intercept and persistence.

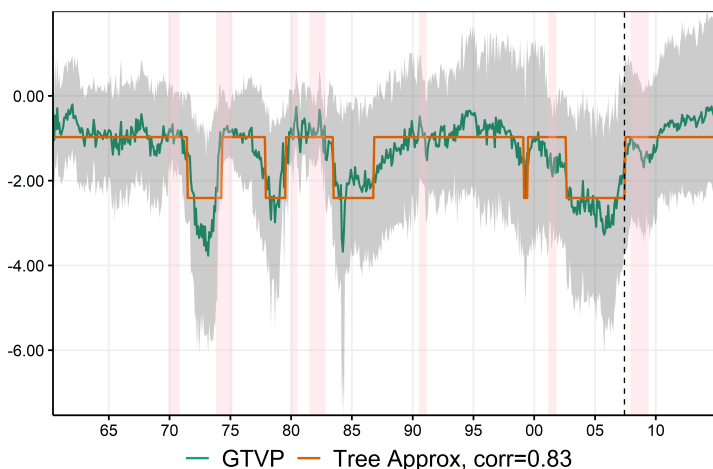
Figures 9, 25a and 25b allow to re-conciliate PC forecasting evidence. For instance, a visible PC death zone spans all of the 90s, which constitutes most of the sample used in Atkeson et al. (2001).<sup>47</sup> It also includes the post-2008 period, which motivated Blanchard et al. (2015)'s inquiry. Most interestingly, for the latter era,  $\gamma_{1,t}$  is predicted to head toward 0 *out-of-sample*. To clarify, the parameter is driven by post-2008 data, but the structure itself ( $\mathcal{F}$ ) is not re-evaluated past the dotted line.

By looking at predictive performance results *ex-post*, Stock and Watson (2008) report that Phillips' curve forecasts usually outperform univariate benchmarks around turning points, but suffer a reversal of fortune when the output/unemployment gap is close to 0. They note that the finding "cannot yet be used to improve forecasts" because their gap relies on a two-sided filter. More recently, Kotchoni et al. (2019) reinforce this view by showing an ARMA(1,1) is triumphant for inflation *except* in recessionary periods, where a data-rich environment can be helpful. But to capitalize on this, one needs a recession/expansion forecast. MRF recognize this potential and relies on leading indicators of the housing market to activate  $\gamma_{1,t}$  in a timely manner. This is particularly evident from looking at  $\gamma_{1,t}$ 's VI measure in Figure 24 and its resulting GTVP

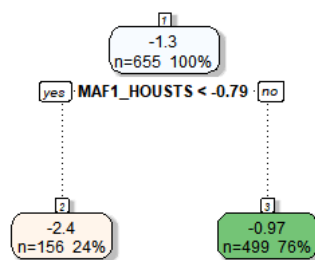
---

<sup>46</sup>Indeed, they do so either by fitting the post-2008 data directly, or by choosing a specification (or building a theoretical model) directly inspired by the experience of the Great Recession.

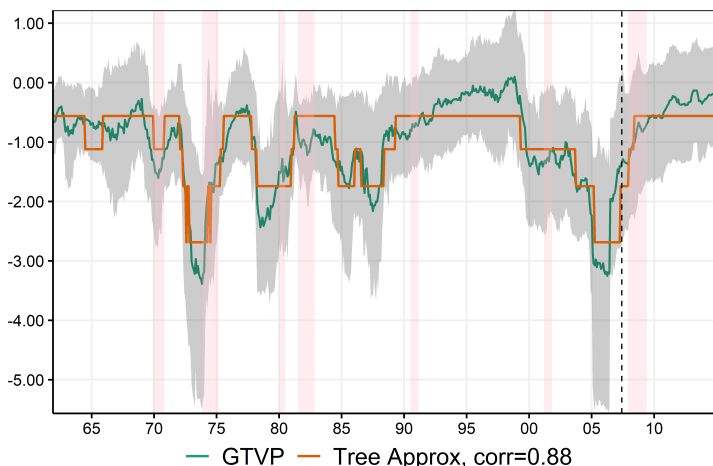
<sup>47</sup>The decade-long wedge between the OLS estimate and GTVP in Figure 25b nicely explains PC failures.



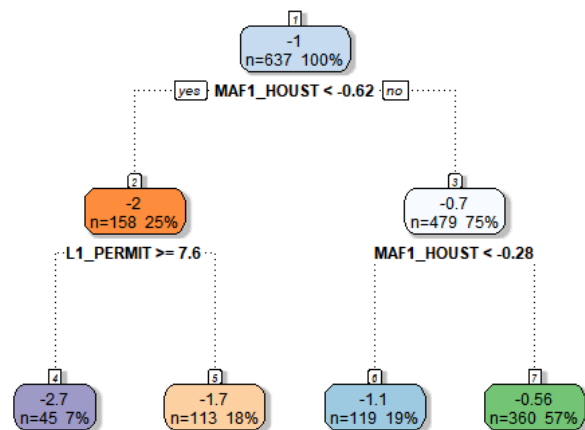
(a)  $\gamma_{1,t}^{INF,h=1}$ : Surrogate Model Replication



(b)  $\gamma_{1,t}^{INF,h=1}$ : Corresponding Tree



(c)  $\gamma_{1,t}^{INF,h=12}$ : Surrogate Model Replication



(d)  $\gamma_{1,t}^{INF,h=12}$ : Corresponding Tree

Figure 9: Surrogate  $\beta_{t,k}$  Trees for Inflation. Shade is 68% credible region. Pink shading is NBER recessions.

in Figure 9. Overall, we see that the relationship between inflation and economic activity is episodic, as conjectured by [Stock and Watson \(2008\)](#), and often prevails before recessions (but not all). Figure 9 proposes a clear-cut answer: inflation responds to the real activity factor when the housing market is booming.

For a long time, housing sector indicators have been known as predictors of future economic activity ([Stock and Watson, 1998a](#); [Leamer, 2007](#)). However, when it comes to forecasting inflation itself, including leading indicators (like permits) does not remedy Phillips' curve forecasts failures ([Stock and Watson, 2007](#)). FAARRF differs by *not* using housing permits/starts as a replacement and/or additional output gap proxy. Rather, its role is to increase the curvature when the time is right. As mentioned above, one explanation is that housing starts and permits

are proxying for future economic activity, resolving the conundrum posed by [Stock and Watson \(2008\)](#). Overall, this implies a PC which would be highly nonlinear in real activity, as further inquired in section 5.4. Another hypothesis is that MRF discovers – through aggregate data – how to leverage [Stock and Watson \(2019\)](#)'s insights that some components of inflation are much more cyclically sensitive than others. [Stock and Watson \(2019\)](#) show that the most cyclical component of inflation is *housing*, followed closely by food components. Accordingly, MRF activating  $\gamma_{1,t}$  with building permits and housing starts is the algorithm's way of predicting when more cyclically sensitive components take the front stage – and by doing so, revive the Philipps' curve. In sum, nonlinearities would be a consequence of aggregation.

Anyhow, the predictive PC studied here differs in many aspects to those studied, for instance, in [Blanchard et al. \(2015\)](#). Most importantly,  $F_1$  summarize indicators that are (for most of them) in first differences. A typical output/unemployment gap measure will be much more persistent. Economically, this means the gap can remain negative for many years following a downturn. In contrast,  $F_1$ , which is strongly correlated with the first difference of UR, will go back up as soon as UR stops growing. To validate current insights and obtain new ones, I complete this section by looking at a prototypical Phillipps' Curve.

## 5.4 A More Traditional Phillipps' Curve

The behavior of inflation since the Great Recession – starting with the missing disinflation and followed by "missing inflation" of recent years – sparked renewed interest in the Phillipps curve. Much attention has been given to its hypothesized flattening ([Blanchard et al., 2015](#); [Galí and Gambetti, 2019](#); [Del Negro et al., 2020](#)). This body of work supports the view that the PC coefficient (either reduced-form or semi-structural) has substantially declined over the last decades. The focus on slow structural change is operationalized by the modeling strategy – either random walk TVPs or sample splitting at a specific date. [Coibion and Gorodnichenko \(2015\)](#) show less worry about PC's health. They rationalize post-2008 inflation with a simple OLS PC where expectations are based on consumer survey data rather than lags or professional forecasters. [Del Negro et al. \(2015\)](#) demonstrate that a standard DSGE (which encompasses a structural New Keynesian PC) is not baffled by post-2008 inflation since it relies on model-based forward-looking expectations of future marginal cost. More recently, [Lindé and Trabandt \(2019\)](#) and [Mineyama \(2020\)](#) articulate theories supporting a nonlinear specification for the reduced-form PC, which could also account for the inflation puzzles punctuating the last 12 years. Given this background and forecasting results reported earlier, a traditional PC must be a fertile ground for MRF-based detective work.

I contribute to the literature by fitting an MRF which linear part corresponds to an expectations-augmented Phillipps' curve.  $X_t$  is inspired by what [Blanchard et al. \(2015\)](#) (henceforth BCS) con-

siders:

$$\pi_t = \theta_t \hat{\pi}_t^{LR} + (1 - \theta_t) \hat{\pi}_t^{SR} + \phi_t u_t^{GAP} + \psi_t \pi_t^{IMP} + \epsilon_t, \quad (5)$$

where  $\pi_t$  stands for CPI inflation,  $\hat{\pi}_t^{LR}$  and  $\hat{\pi}_t^{SR}$  respectively for long-run and short-run inflation expectations.  $u_t^{GAP}$  represents the (negative) unemployment gap and  $\pi_t^{IMP}$  is import prices inflation. I translate this to the MRF framework by making  $\mu_t = \theta_t \hat{\pi}_t^{LR}$  the time-varying intercept, letting  $\beta_{1,t} = 1 - \theta_t$  and by obtaining  $u_t^{GAP}$  by means of Hodrick-Prescott filtering.<sup>48</sup> As in BCS,  $\hat{\pi}_t^{SR}$  is the average inflation over the last four quarters. Hence, the estimated equation

$$\pi_t = \mu_t + \beta_{1,t} \hat{\pi}_t^{SR} + \beta_{2,t} u_t^{GAP} + \beta_{3,t} \pi_t^{IMP} + \epsilon_t \quad (6)$$

does not impose the constraint implied by  $\theta_t$  in equation (5). However, estimation results will desirably have  $\beta_{1,t} \in [0, 1]$  at almost any point in time.  $S_t$  is the same as that considered in the forecasting section. The data set runs up to 2019Q4.

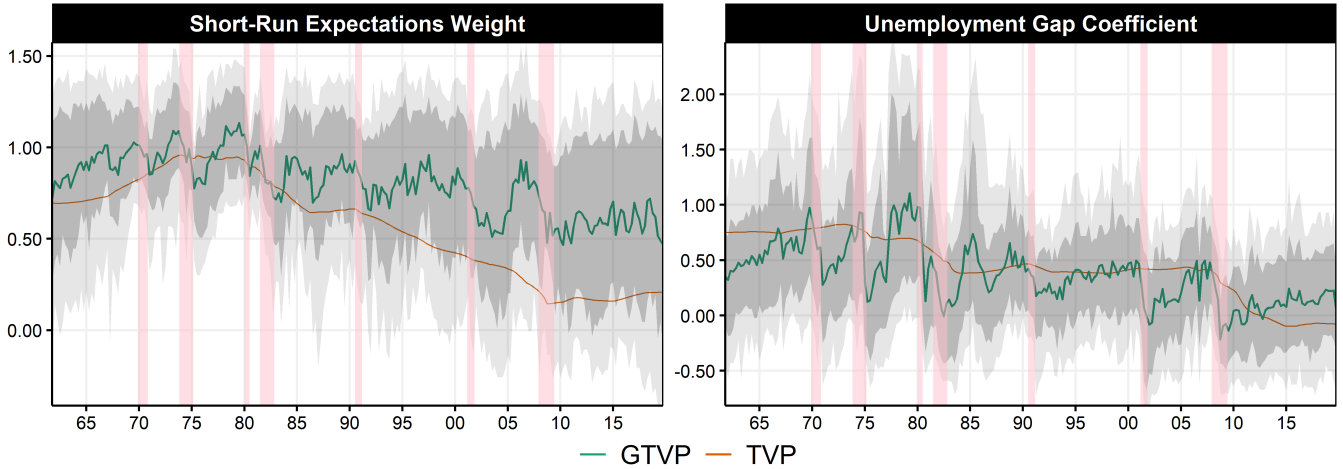


Figure 10: The gray bands are the 68% and 90% credible regions. Pink shading corresponds to NBER recessions.

Figure 10 reports GTVPs of interest: the weight on short-run expectations and the output gap coefficient. Additionally, it contains traditional TVP estimates as means of comparison. The latter convey the usual wisdom: inflation expectations slowly start to be more anchored from the mid 1980s. Around the same time, the unemployment/inflation trade-off begins its slow collapse. The updated data shows that the TVP-based Phillips' curve has further flattened to plain 0 in the last decade.

For  $\beta_{1,t}$ , the weight on short-run expectations, both methods agree that it has been decreasing steadily after the 1983 recession. But GTVPs highlight an additional pattern for the importance

<sup>48</sup>Specifically, both this gap and that of BCS get out of negative territory around 2014.

of  $\hat{\pi}_t^{SR}$ : it tends to increase during economic expansions, collapse during recessions then start increasing again until the next downturn. Note that the phenomenon is also observed in Figure 7b for the simpler ARRF on quarterly inflation. The decrease in the coefficient (usually of about 0.25) is observed for *every* recession and usually last for some additional quarters after the end of it. The linear rise in the coefficient occurs for all expansions except those preceding the early 90s and 2000s recessions, where the pattern is punctuated with additional peaks and troughs. The increased importance of short-run expectations with the age of the expansion is also observed for recent expansionary periods. Hence, the phenomenon is not merely a matter of the 70s and 80s recessions being preceded by a sharp acceleration of inflation.

From a more statistical point a view, the sharp decline in  $\beta_{1,t}$  following every recession suggests that in the aftermath of an important downward shock, the long-run inflation expectation is a more reliable predictor as it is minimally affected by recent events. As the expansion slowly progress (and recessionary data points get out of the short-run average),  $\hat{\pi}_t^{SR}$  becomes a more up to date and reliable barometer of future inflation conditions. This narrative is corroborated by variable importance (Figure 27) for  $\beta_{1,t}$ , which highlights the importance of the trend, but also recent lags of inflation.

When it comes to the low-frequency movement of the unemployment gap coefficient, both methods agree about a significant decline starting from the 80s. However, GTVPs uncover additional heterogeneity. **First** and most strikingly,  $\beta_{2,t}$  gets very close to 0 following every recession. This suggests a nonlinear Philipps' curve where inflation responds strongly to a very positive  $u_t^{GAP}$  but not so much to a negative one. **Second**, the 70s and early 80s are characterized as a series of peaks (preceding the first three recessions of the sample) rather than a sustained high coefficient. Traditional TVPs, by excessive time-smoothing, dissimulate the effects of inflationary spirals on  $\beta_{2,t}$ . Such pre-recession accelerations still occur during the Great Moderation but in a much milder way.

**Third**, VI measures (in Figure 27) confirm the importance of activity indicators (like Total Capacity Utilization (TCU)) in driving  $\beta_{2,t}$  itself. The correlation between  $\beta_{2,t}$  and TCU is 0.81, and the correspondence between the two variables is striking in Figure 11. Many notable increases in  $\beta_{2,t}$  are nicely matched (between the two 70s recessions and before 2008). Of course, this simple characterization remains imperfect since it misses some highs (like the end of the 70s) and predicts a higher  $\beta_{2,t}$  in the years following the 2008-2009 recession. Generally, given the strong co-cyclicalities between TCU and  $u_t^{GAP}$ , this is evidence of a *convex* PC.

The collapse of  $\beta_{2,t}$  following recessions is not unique to 2008: it happened following *every* recession since 1960. As a result, inflation will rise when the economy is running well above its potential, but much more timidly will it go down from economic slack. Recently, [Lindé and Trabandt \(2019\)](#) have shown that such a phenomenon can be rationalized by a New Keynesian

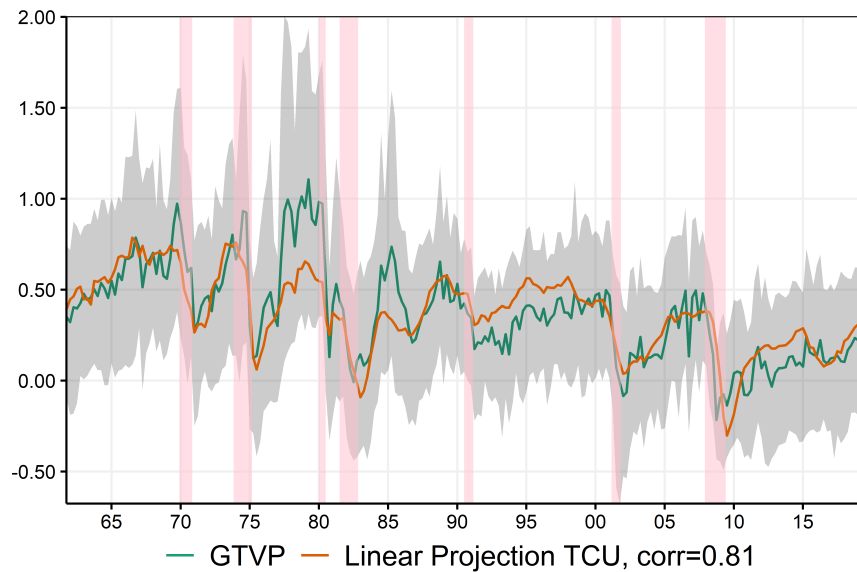


Figure 11: "What Goes Around Comes Around": Capacity Utilization is substantially correlated with the inflation-unemployment trade-off. The gray band is the 68% credible region. Pink shading corresponds to NBER recessions.

DSGE model. Indeed, by allowing for additional strategic complementarity in firms price- and wage-setting behavior and solving the nonlinear model (rather than considering the linear approximation around the steady state), the authors obtain a state-dependent PC which becomes very flat during large downturns. This can explain both the small coefficient during recessions and its subsequent timid increase. Theoretically, convexity can also emerge from downward wage rigidities (Mineyama, 2020), but its empirical plausibility for the post-2008 era has been contested (Coibion and Gorodnichenko, 2015).

This pattern remains when adding controls in the linear part for supply shocks and monetary policy shocks. Those are the usual confounding factors suspected of blurring the relationship by introducing a positive correlation between unemployment and inflation.<sup>49</sup> The economic suspicion particular to this application is that omitting them could create a downward bias in  $\beta_{2,t}$  that only occurs locally, generating the cyclical pattern. As it turns out, controls make cyclicity even more obvious in Figure 28, especially for the later part of the sample.<sup>50</sup> However, the overall strength of the coefficient is smaller (especially for the 70s).

Many hypotheses can be accommodated by a model estimated on two disjoint samples, like in Del Negro et al. (2020). Much fewer of them are compatible with the richer  $\beta_{2,t}$  path extracted by MRF. This is important: learning the type of nonlinearity, rather than partially imposing it, helps in discriminating economic suppositions. Figure 11 and recent theoretical developments

<sup>49</sup>While the time-varying constant can go a long way at controlling for such factors – being a RF in itself, including them in the linear part makes them "stand out" as everything going through the intercept is inevitably heavily regularized.

<sup>50</sup>Results being similar for both curves is reminiscent of Galí and Gambetti (2019) who report little differences between paths of reduced-form and semi-structural wage PCs (although they focus on long-run change).



both suggest that much of the PC's decline is attributable to upward nonlinearities being less solicited in the last 3 decades. This is in accord with the policy hypothesis: since Paul Volker's chairmanship the monetary authority has responded much more aggressively to inflationary pressures, limiting the spirals that gave rise to high  $\beta_{2,t}$ 's in the 70s. Two conclusions emerge from this observation. First, exogenous change cannot so simply be ruled out. Second, knowing what were MRF beliefs about PC nonlinearities at different points in time could be enlightening.

#### 5.4.1 Conditional Coefficient Forecasting

$\beta_{2,t}$ 's lows are getting lower, and longer. Should we have known? Much of the recent work on PC is directly inspired by Great Recession aftermath, and aims at explaining it. Whether it is theoretical or empirical work, much of it could be overfitting: a model can replicate one or two facts it is trained to replicate, but fails to generalize. That is, even if models are tested out-of-sample (which is itself not so often the case in the literature), the choice of nonlinearity itself is often determined in attempt to match the OOS. Beyond the linear part being a PC, MRF does not assume much — and its nonlinearities are certainly not "personalized" to the recent inflation experience. Thus, it is interesting to ask: what was MRF "thinking" about  $\beta_{2,t}$  in 2007? in 1995? Did it know something we did not, or did it learn (as most economists) of PC's collapse from the post-2008 experience? I conduct a  $\beta_{2,t}$  dynamic learning exercise to find out.

To make this operational, MRF is estimated up to 1995, 2007 and 2019, and GTVPs are projected out-of-sample from those dates (when applicable). To be clear,  $\hat{\beta}_{2,t|1995} = \hat{\mathcal{F}}_{1995}(S_t)$  means the coefficient *predictive structure* is last estimated in 1995. Coefficients keep moving out-of-sample because  $S_t$  does.  $\hat{\mathcal{F}}_{1995}(S_t)$  and  $\hat{\mathcal{F}}_{2007}(S_t)$  will differ for two main reasons. The first is estimation error – both in terms of precision and re-evaluating which nonlinearity seems more appropriate.<sup>51</sup> The second is structural change, perhaps completely exogenous or triggered by policy interventions.

Much can be learned from Figure 12. First, GTVPs are all very alike for the pre-1995 period, suggesting little was observed post-1995 that made MRF change its reading of the past. Similarly, the green and the magenta line, which both share the 1995-2007 period within their training sets, are close to one another. Overall, this indicates that OOS difference between paths are very unlikely due to a better re-estimation and/or a completely new choice of  $\mathcal{F}$ .

Second, unlike what we have seen for the unemployment equation (Figure 5), there are important disparities between the ex-ante and the ex-post paths *out-of-sample*. Thus, one can rightfully hypothesize that structural change got in the way, making  $\hat{\mathcal{F}}_{1995}$ 's attempt of replicating the strong nonlinearities of the 70s into the 2000s go wildly off course. An analogous (yet far less noticeable gap) punctuates the post-2007 period. This suggests that while  $\beta_{2,t}$  was expected

<sup>51</sup>The second part has the flavor of model selection "error".

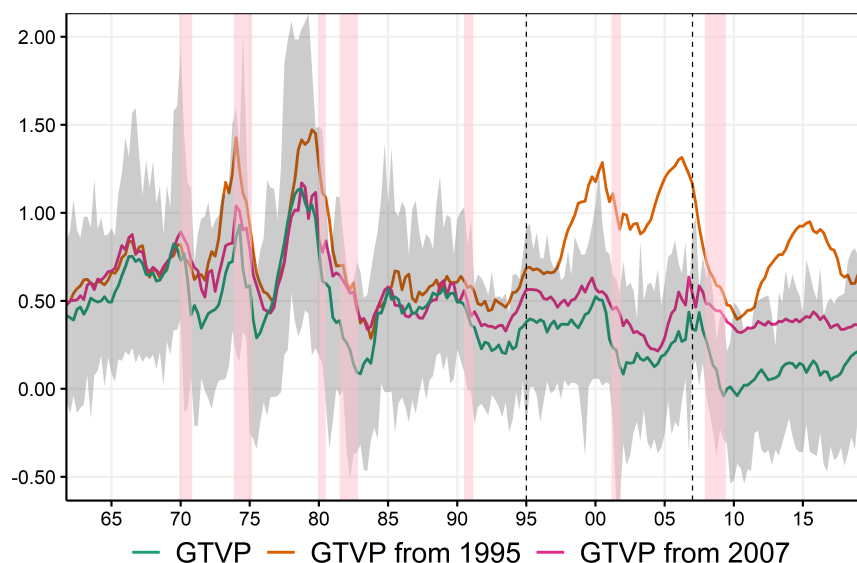


Figure 12: Conditional  $\beta_{2,t}$  Forecasting. The gray band is the 68% credible region for GTVPs estimated up to 2019Q4. Pink shading corresponds to NBER recessions. For enhanced visibility, GTVPs are smoothed with 1-year moving average. The vertical dotted lines are the end of the training samples.

to fall marginally following the crisis and stay low thereafter (according to  $\hat{\mathcal{F}}_{2007}$ ), it was not expected to go *that* low. Indeed, only  $\hat{\mathcal{F}}_{2019}$  hits 0 and stays in its vicinity.

Of course, by design, exogenous structural change cannot be captured out-of-sample – with the results that we know ( $\hat{\mathcal{F}}_{1995}$ ). This dismal predicament does not apply to cyclical behavior: it has been forecastable at least since 1995. Indeed,  $\hat{\mathcal{F}}_{1995}$  propose a  $\beta_{2,t}$  for 2000 and 2008 that is very similar to that of 70s inflation spirals. Moreover,  $\hat{\beta}_{2,t|1995}$ 's collapse following 2008 is of a magnitude only seen during Arthur Burns' days. Hence, a much weaker PC following large downturns is hardly new. However, what  $\hat{\beta}_{2,t|2007}$  and  $\hat{\beta}_{2,t|2019}$  tell us is that the overall amplitude (and level) of those variations has evolved exogenously, forcing MRF to update  $\mathcal{F}$  repeatedly.

This exercise may rightfully seem exotic, with no obvious analog in the literature. The simple explanation is that traditional time variations only give "trivial" parameter forecasts by construction, and there is no clear "leaning" process to analyze. For example, the "forecasted" random walk TVP would be a straight line over the whole OOS. Doing so with a threshold model would only inform us of the increasing precision of estimation as sample size grows – i.e., the model itself cannot be re-evaluated.

An avenue to be explored in future work is to specify, under clear conditions, the meaning of  $\hat{\beta}_{2,t|2007} - \hat{\beta}_{2,t|1995}$ . For instance, if we are willing to assume that any change in *structure* MRF fails to capture dynamically origins from policy shifts (like evolving monetary policy), then the difference is the treatment effect of policy change on the reduced-form coefficients – i.e., the measurable effect of the Lucas critique. In the PC case,  $\hat{\beta}_{2,t|2007} - \hat{\beta}_{2,t|1995}$  indicates that the difference is most salient during periods of economic overheating (within which we know the monetary

authority is now more active). This sort of analysis is possible because, unlike traditional non-linear methods, MRF provides non-trivial "counterfactual"  $\beta_t$  paths out-of-sample. Indeed, it discovers structural change rather than imposing it. Essentially, this line of work could extend some of the conditional forecasting toolbox and insights (Waggoner and Zha, 1999) to conditional coefficients forecasting.

## 6 Conclusion

I proposed a new time series model that **(i)** expands multiple non-linear time series models, **(ii)** adapts Random Forest for Macro forecasting and **(iii)** can be interpreted as Generalized Time-Varying Parameters. On the empirical front, the methodology provides substantial empirical gains over RF and competing non-linear TS models. The resulting Generalized TVPs have a very distinct behavior vis-à-vis standard random walk parameters. For instance, they adapt nicely to regime-switching behavior that seems pervasive for unemployment – while not neglecting potential long-run change. This finding is facilitated by the fact that GTVPs lend themselves much more easily to interpretation than either standard RF or random-walk TVPs. Indeed, rather than trying to open the back-box of an opaque conditional mean function (like one would with plain RF), MRFs can be compartmentalized in different components of the small macro model. Furthermore, GTVPs can be visualized with standard time series plots and credible intervals are provided by a variant of the Bayesian Bootstrap.

When looking at Phillips' curves in general, MRF finds both structural change in the persistence and regime-dependent behavior in the economic activity/inflation trade-off. In particular, a recurrent theme across all specifications is that the slowly decaying curve is also much steeper when the economy is overheating – in line with the convexity/nonlinearity hypothesis. Hence, MRF can be of great help sorting out what is plausible and what is not when it comes to macroeconomic equations with a history of controversy. Since there is no shortage of those, MRF holds many possibilities for future research.

## References

- Adrian, T., Boyarchenko, N., and Giannone, D. (2019). Vulnerable growth. *American Economic Review*, 109(4):1263–89.
- Alexander, W. P. and Grimshaw, S. D. (1996). Treed regression. *Journal of Computational and Graphical Statistics*, 5(2):156–175.
- Almon, S. (1965). The distributed lag between capital appropriations and expenditures. *Econometrica: Journal of the Econometric Society*, pages 178–196.
- Amir-Ahmadi, P., Matthes, C., and Wang, M.-C. (2018). Choosing prior hyperparameters: with applications to time-varying parameter models. *Journal of Business & Economic Statistics*, pages 1–13.
- Aruoba, S. B., Bocola, L., and Schorfheide, F. (2017). Assessing dsge model nonlinearities. *Journal of Economic Dynamics and Control*, 83:34–54.
- Athey, S., Tibshirani, J., Wager, S., et al. (2019). Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178.
- Atkeson, A., Ohanian, L. E., et al. (2001). Are phillips curves useful for forecasting inflation? *Federal Reserve bank of Minneapolis quarterly review*, 25(1):2–11.
- Auerbach, A. J. and Gorodnichenko, Y. (2012a). Fiscal multipliers in recession and expansion. In *Fiscal policy after the financial crisis*, pages 63–98. University of Chicago Press.
- Auerbach, A. J. and Gorodnichenko, Y. (2012b). Measuring the output responses to fiscal policy. *American Economic Journal: Economic Policy*, 4(2):1–27.
- Bai, J. and Ng, S. (2008). Forecasting economic time series using targeted predictors. *Journal of Econometrics*, 146(2):304–317.
- Bañbura, M., Giannone, D., and Reichlin, L. (2010). Large bayesian vector auto regressions. *Journal of Applied Econometrics*, 25(1):71–92.
- Batini, N., Callegari, G., and Melina, G. (2012). Successful austerity in the united states, europe and japan.
- Bergmeir, C., Hyndman, R. J., and Koo, B. (2018). A note on the validity of cross-validation for evaluating autoregressive time series prediction. *Computational Statistics & Data Analysis*, 120:70–83.
- Blanchard, O. (2016). The phillips curve: Back to the’60s? *American Economic Review*, 106(5):31–34.
- Blanchard, O., Cerutti, E., and Summers, L. (2015). Inflation and activity—two explorations and their monetary policy implications. Technical report, National Bureau of Economic Research.
- Boivin, J. (2005). Has us monetary policy changed? evidence from drifting coefficients and real-time data. Technical report, National Bureau of Economic Research.
- Borup, D., Christensen, B. J., Mühlbach, N. N., Nielsen, M. S., et al. (2020a). Targeting predictors in random forest regression. Technical report, Department of Economics and Business Economics, Aarhus University.

- Borup, D., Rapach, D., and Schütte, E. C. M. (2020b). Now-and backcasting initial claims with high-dimensional daily internet search-volume data. *Available at SSRN 3690832*.
- Breiman, L. (1996). Bagging predictors. *Machine learning*, 24(2):123–140.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- Breiman, L., Friedman, J., Stone, C. J., and Olshen, R. A. (1984). *Classification and regression trees*. CRC press.
- Chamberlain, G. and Imbens, G. W. (2003). Nonparametric applications of bayesian inference. *Journal of Business & Economic Statistics*, 21(1):12–18.
- Chan, J. C., Eisenstat, E., and Strachan, R. W. (2018). Reducing dimensions in a large TVP-VAR. CAMA Working Papers 2018-49, Centre for Applied Macroeconomic Analysis, Crawford School of Public Policy, The Australian National University.
- Chen, J. C., Dunn, A., Hood, K. K., Driessen, A., and Batch, A. (2019). Off to the races: A comparison of machine learning and alternative data for predicting economic indicators. In *Big Data for 21st Century Economic Statistics*. University of Chicago Press.
- Chevillon, G. (2007). Direct multi-step estimation and forecasting. *Journal of Economic Surveys*, 21(4):746–785.
- Cirillo, P. and Muliere, P. (2013). An urn-based bayesian block bootstrap. *Metrika*, 76(1):93–106.
- Clarida, R., Gali, J., and Gertler, M. (2000). Monetary policy rules and macroeconomic stability: evidence and some theory. *The Quarterly journal of economics*, 115(1):147–180.
- Clements, M. P. and Smith, J. (1997). The performance of alternative forecasting methods for setar models. *International Journal of Forecasting*, 13(4):463–475.
- Clyde, M. and Lee, H. (2001). Bagging and the bayesian bootstrap. In *AISTATS*.
- Cogley, T., Primiceri, G. E., and Sargent, T. J. (2010). Inflation-gap persistence in the us. *American Economic Journal: Macroeconomics*, 2(1):43–69.
- Cogley, T. and Sargent, T. J. (2001). Evolving post-world war ii us inflation dynamics. *NBER macroeconomics annual*, 16:331–373.
- Coibion, O. and Gorodnichenko, Y. (2015). Is the phillips curve alive and well after all? inflation expectations and the missing disinflation. *American Economic Journal: Macroeconomics*, 7(1):197–232.
- D’Agostino, A., Gambetti, L., and Giannone, D. (2013). Macroeconomic forecasting and structural change. *Journal of Applied Econometrics*, 28(1):82–101.
- Dagum, E. B. and Bianconcini, S. (2009). Equivalent reproducing kernels for smoothing spline predictors. In *Proceedings of the American Statistical Association, Business and Economic Statistics Section*.
- de Wind, J. and Gambetti, L. (2014). Reduced-rank time-varying vector autoregressions. CPB Discussion Paper 270, CPB Netherlands Bureau for Economic Policy Analysis.
- Del Negro, M., Giannoni, M. P., and Schorfheide, F. (2015). Inflation in the great recession and new keynesian models. *American Economic Journal: Macroeconomics*, 7(1):168–96.

- Del Negro, M., Lenza, M., Primiceri, G. E., and Tambalotti, A. (2020). What's up with the phillips curve? Technical report, National Bureau of Economic Research.
- Delle Monache, D., De Polis, A., and Petrella, I. (2020). Modeling and forecasting macroeconomic downside risk.
- Diebold, F. X. and Mariano, R. S. (2002). Comparing predictive accuracy. *Journal of Business & economic statistics*, 20(1):134–144.
- Diebold, F. X. and Rudebusch, G. D. (1994). Measuring business cycles: A modern perspective. Technical report, National Bureau of Economic Research.
- Doan, T., Litterman, R., and Sims, C. (1984). Forecasting and conditional projection using realistic prior distributions. *Econometric reviews*, 3(1):1–100.
- Dolado, J. J., Maria-Dolores, R., and Naveira, M. (2005). Are monetary-policy reaction functions asymmetric?: The role of nonlinearity in the phillips curve. *European Economic Review*, 49(2):485–503.
- Doser, A., Nunes, R. C., Rao, N., and Sheremirov, V. (2017). Inflation expectations and nonlinearities in the phillips curve.
- Duroux, R. and Scornet, E. (2016). Impact of subsampling and pruning on random forests. *arXiv preprint arXiv:1603.04261*.
- Estrella, A. and Mishkin, F. S. (1998). Predicting us recessions: Financial variables as leading indicators. *Review of Economics and Statistics*, 80(1):45–61.
- Faust, J. and Wright, J. H. (2013). Forecasting inflation. In *Handbook of economic forecasting*, volume 2, pages 2–56. Elsevier.
- Ferguson, T. S. (1973). A bayesian analysis of some nonparametric problems. *The annals of statistics*, pages 209–230.
- Fortin-Gagnon, O., Leroux, M., Stevanovic, D., and Surprenant, S. (2018). A large canadian database for macroeconomic analysis. Technical report, Department of Economics, UQAM.
- Freedman, D. A. et al. (1981). Bootstrapping regression models. *The Annals of Statistics*, 9(6):1218–1228.
- Friedberg, R., Tibshirani, J., Athey, S., and Wager, S. (2018). Local linear forests. *arXiv preprint arXiv:1807.11408*.
- Friedman, J., Hastie, T., and Tibshirani, R. (2001). *The elements of statistical learning*, volume 1. Springer series in statistics New York, NY, USA:.
- Galí, J. and Gambetti, L. (2019). Has the us wage phillips curve flattened? a semi-structural exploration. Technical report, National Bureau of Economic Research.
- Ghysels, E., Sinko, A., and Valkanov, R. (2007). Midas regressions: Further results and new directions. *Econometric Reviews*, 26(1):53–90.
- Giraitis, L., Kapetanios, G., and Yates, T. (2014). Inference on stochastic time-varying coefficient models. *Journal of Econometrics*, 179(1):46–65.

- Goulet Coulombe, P. (2020a). Time-varying parameters as ridge regressions. *arXiv preprint arXiv:2009.00401*.
- Goulet Coulombe, P. (2020b). To bag is to prune. *arXiv preprint arXiv:2008.07063*.
- Goulet Coulombe, P., Leroux, M., Stevanovic, D., Surprenant, S., et al. (2019). How is machine learning useful for macroeconomic forecasting? Technical report, CIRANO.
- Goulet Coulombe, P., Leroux, M., Stevanovic, D., Surprenant, S., et al. (2020a). Macroeconomic data transformations matter. Technical report, CIRANO.
- Goulet Coulombe, P., Leroux, M., Stevanovic, D., Surprenant, S., et al. (2020b). Pr evision de l'activit e  conomique au qu ebec et au canada   l'aide des m ethodes "machine learning". Technical report, Technical report, CIRANO.
- Grandvalet, Y. (2004). Bagging equalizes influence. *Machine Learning*, 55(3):251–270.
- Granger, C. W. (2008). Non-linear models: Where do we go next-time varying parameter models? *Studies in Nonlinear Dynamics & Econometrics*, 12(3).
- Hahn, P. R., Carvalho, C. M., and Mukherjee, S. (2013). Partial factor modeling: predictor-dependent shrinkage for linear regression. *Journal of the American Statistical Association*, 108(503):999–1008.
- Hansen, B. E. (2011). Threshold autoregression in economics. *Statistics and its Interface*, 4(2):123–127.
- Hansen, C. and Liao, Y. (2019). The factor-lasso and k-step bootstrap approach for inference in high-dimensional economic applications. *Econometric Theory*, 35(3):465–509.
- Hassani, H., Heravi, S., and Zhigljavsky, A. (2009). Forecasting european industrial production with singular spectrum analysis. *International journal of forecasting*, 25(1):103–118.
- Hassani, H., Soofi, A. S., and Zhigljavsky, A. (2013). Predicting inflation dynamics with singular spectrum analysis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 176(3):743–760.
- Hillebrand, E., Lukas, M., Wei, W., et al. (2020). Bagging weak predictors. Technical report, Monash University, Department of Econometrics and Business Statistics.
- Hillebrand, E. and Medeiros, M. C. (2010). The benefits of bagging for forecast models of realized volatility. *Econometric Reviews*, 29(5-6):571–593.
- Inoue, A. and Kilian, L. (2008). How useful is bagging in forecasting economic time series? a case study of us consumer price inflation. *Journal of the American Statistical Association*, 103(482):511–522.
- Karabatsos, G. (2016). A dirichlet process functional approach to heteroscedastic-consistent covariance estimation. *International Journal of Approximate Reasoning*, 78:210–222.
- Kotchoni, R., Leroux, M., and Stevanovic, D. (2019). Macroeconomic forecast accuracy in a data-rich environment. *Journal of Applied Econometrics*, 34(7):1050–1072.
- Lancaster, T. (2003). A note on bootstraps and robustness. Available at SSRN 896764.

- Leamer, E. E. (2007). Housing is the business cycle. Technical report, National Bureau of Economic Research.
- Lee, T.-H., Ullah, A., and Wang, R. (2020). Bootstrap aggregating and random forest. In *Macroeconomic Forecasting in the Era of Big Data*, pages 389–429. Springer.
- Lindé, J. and Trabandt, M. (2019). Resolving the missing deflation puzzle.
- Lucas, R. E. (1976). Econometric policy evaluation: A critique. In *Carnegie-Rochester conference series on public policy*, volume 1, pages 19–46.
- MacKinnon, J. G. (2006). Bootstrap methods in econometrics. *Economic Record*, 82:S2–S18.
- McCracken, M. and Ng, S. (2020). Fred-qd: A quarterly database for macroeconomic research. Technical report, National Bureau of Economic Research.
- McCracken, M. W. and Ng, S. (2016). Fred-md: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics*, 34(4):574–589.
- Medeiros, M. C., Vasconcelos, G. F., Veiga, Á., and Zilberman, E. (2019). Forecasting inflation in a data-rich environment: the benefits of machine learning methods. *Journal of Business & Economic Statistics*, (just-accepted):1–45.
- Meek, C., Chickering, D. M., and Heckerman, D. (2002). Autoregressive tree models for time-series analysis. In *Proceedings of the 2002 SIAM International Conference on Data Mining*, pages 229–244. SIAM.
- Meinshausen, N. (2006). Quantile regression forests. *Journal of Machine Learning Research*, 7(Jun):983–999.
- Mineyama, T. (2020). Downward nominal wage rigidity and inflation dynamics during and after the great recession. *Available at SSRN 3157995*.
- Molnar, C. (2019). *Interpretable machine learning*. Lulu.com.
- Olson, M. A. and Wyner, A. J. (2018). Making sense of random forest probabilities: a kernel perspective. *arXiv preprint arXiv:1812.05792*.
- Perron, P. et al. (2006). Dealing with structural breaks. *Palgrave handbook of econometrics*, 1(2):278–352.
- Poirier, D. J. (2011). Bayesian interpretations of heteroskedastic consistent covariance estimators using the informed bayesian bootstrap. *Econometric Reviews*, 30(4):457–468.
- Primiceri, G. E. (2005). Time varying structural vector autoregressions and monetary policy. *The Review of Economic Studies*, 72(3):821–852.
- Ramey, V. A. and Zubairy, S. (2018). Government spending multipliers in good times and in bad: evidence from us historical data. *Journal of Political Economy*, 126(2):850–901.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Rubin, D. B. (1981). The bayesian bootstrap. *The annals of statistics*, pages 130–134.



- Shiller, R. J. (1973). A distributed lag estimator derived from smoothness priors. *Econometrica (pre-1986)*, 41(4):775.
- Sims, C. A. (1993). A nine-variable probabilistic macroeconomic forecasting model. In *Business cycles, indicators and forecasting*, pages 179–212. University of Chicago press.
- Sims, C. A. and Zha, T. (2006). Were there regime switches in us monetary policy? *American Economic Review*, 96(1):54–81.
- Stevanovic, D. (2016). Common time variation of parameters in reduced-form macroeconomic models. *Studies in Nonlinear Dynamics & Econometrics*, 20(2):159–183.
- Stock, J. H. (1994). Unit roots, structural breaks and trends. *Handbook of econometrics*, 4:2739–2841.
- Stock, J. H. and Watson, M. W. (1989). New indexes of coincident and leading economic indicators. *NBER macroeconomics annual*, 4:351–394.
- Stock, J. H. and Watson, M. W. (1998a). Business cycle fluctuations in us macroeconomic time series. Technical report, National Bureau of Economic Research.
- Stock, J. H. and Watson, M. W. (1998b). A comparison of linear and nonlinear univariate models for forecasting macroeconomic time series. Technical report, National Bureau of Economic Research.
- Stock, J. H. and Watson, M. W. (2007). Why has us inflation become harder to forecast? *Journal of Money, Credit and banking*, 39:3–33.
- Stock, J. H. and Watson, M. W. (2008). Phillips curve inflation forecasts. Technical report, National Bureau of Economic Research.
- Stock, J. H. and Watson, M. W. (2019). Slack and cyclically sensitive inflation. Technical report, National Bureau of Economic Research.
- Taddy, M., Chen, C.-S., Yu, J., and Wyle, M. (2015). Bayesian and empirical bayesian forests. *arXiv preprint arXiv:1502.02312*.
- Taddy, M., Gardner, M., Chen, L., and Draper, D. (2016). A nonparametric bayesian analysis of heterogenous treatment effects in digital experimentation. *Journal of Business & Economic Statistics*, 34(4):661–672.
- Teräsvirta, T. (1994). Specification, estimation, and evaluation of smooth transition autoregressive models. *Journal of the american Statistical association*, 89(425):208–218.
- Waggoner, D. and Zha, T. (1999). Conditional forecasts in dynamic multivariate models. *Review of Economics and Statistics*, 81(4):639–651.
- Wang, Y. and Witten, I. H. (1996). Induction of model trees for predicting continuous classes.
- Wei, P., Lu, Z., and Song, J. (2015). Variable importance analysis: a comprehensive review. *Reliability Engineering & System Safety*, 142:399–432.
- Wochner, D. (2020). Dynamic factor trees and forests—a theory-led machine learning framework for non-linear and state-dependent short-term us gdp growth predictions.
- Woloszko, N. (2020). Adaptive trees: a new approach to economic forecasting.

# A Appendix

## A.1 Can Larger Linear Parts Help?

As argued earlier, an advantage of MRF over plain RF is that by taking the TVP view of nonlinearities, we are in a much better position to attempt an interpretation of the successful model. One could rightfully retort that while FAARRF performs nicely, its potential for interpretation is spoiled using factors rather than raw data. While this critique is partially addressable by putting names on factors such as "real activity" and "forward-looking" factors, it is worthwhile to consider alternative dimensionality reductions schemes that keep the data in the original space.

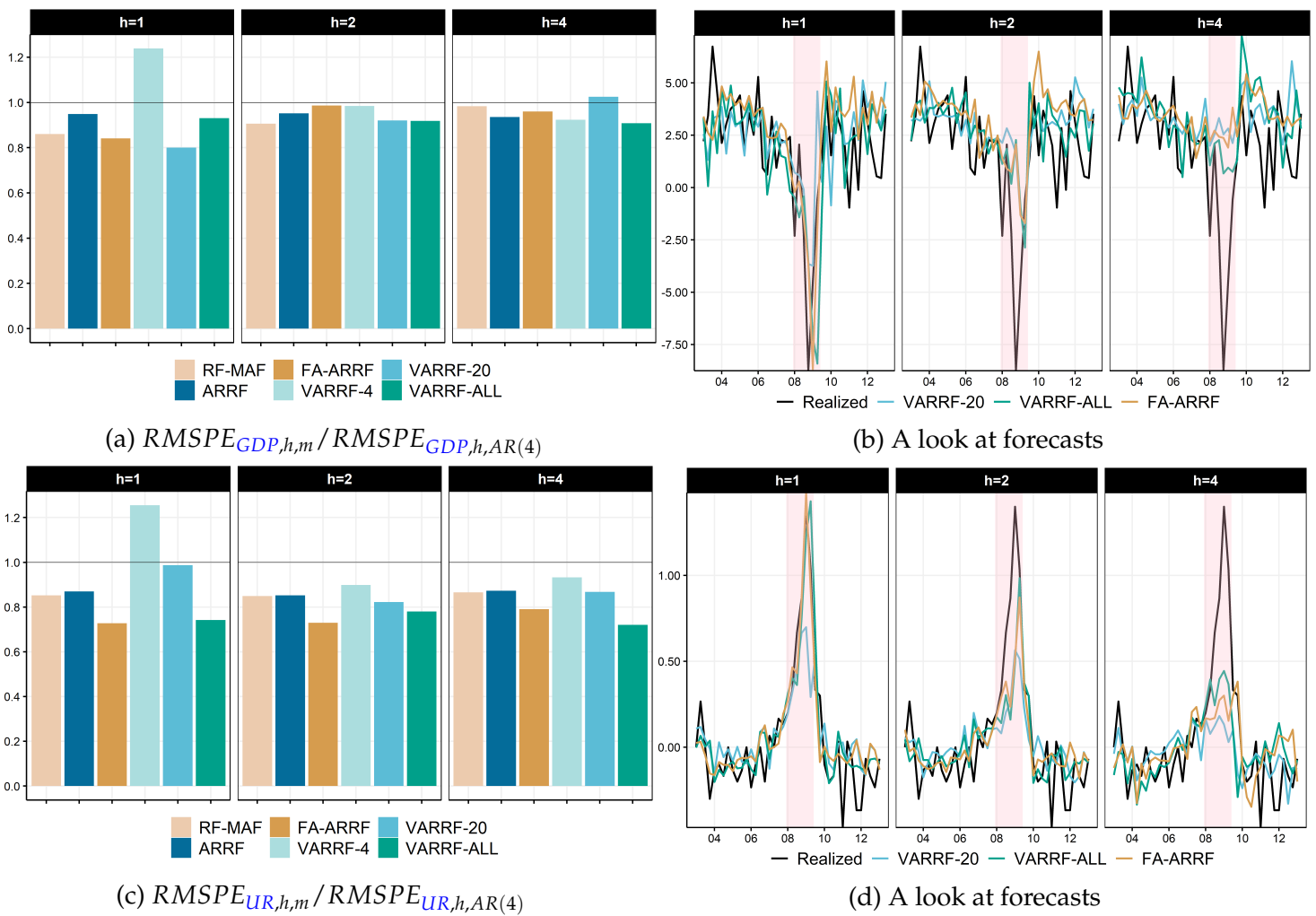


Figure 13: Large Linear Part Results

Since the 4 variables VARRF results are not necessary sterling, the expectations for VARRF of bigger size are rather low. Regardless, Figure 13 show promising results for both a VARRF with 20 variables (in the spirit of Bańbura et al. (2010)'s medium VAR) and a VARRF that includes all FRED-QD 200+ variables. In both cases, especially the later, regularization must be stringent to

keep estimation variance low. Indeed, in VARRF-ALL, the linear part has more regressors than observations even before any split is attempted. This implies a much higher value of  $\lambda$ ,  $\zeta = 1$  and the use of a single lag in the linear part. Quite strikingly, Figure 13d report that VARRF-ALL tracks UR at  $h = 1$  remarkably well, in addition to providing forecasts that clearly hint at a recession for both UR and GDP up to a year ahead. In that latter regard, VARRF-ALL is the best model of the whole lot. Hence, larger models, while not the center of interest for this paper, can be handled by MRF and provide excellent forecasts given proper regularization.

This subsection's results point out that the tool can be easily (and desirably) extended to estimate large GTVP-VARs. The dynamic coefficients can be estimated by either fitting MRF equation by equation. Another possibility (left for future research) is to simply modify the splitting rule in (2) to be multivariate so that each tree is fitted jointly for all equation – pooling time-variation across equations. Finally, elements of the covariance matrix of residuals can be fitted separately with a plain RF, which is very fast.

## A.2 More on Surrogate $\beta_t$ Trees

The approach described in section 5.3 belongs to a family of methods usually referred to as "surrogate models" (Molnar, 2019). Attempting to fit the whole conditional mean obtained from a black-box algorithm using a more transparent model is a global surrogate. An obvious critique of this approach is that if the complicated model justifies its cost in interpretability with its predicting gains, it is hard to believe a simple model can reliably recreate its predictions. Conversely, if the surrogate model is quite successful, this casts some doubts about the relevance of the black box itself. In this line of work, a more promising avenue is a local surrogates model as proposed in Ribeiro et al. (2016), which fits interpretable models *locally*. By following Granger (2008)'s insights, we already have this: by looking at the  $\beta_t$  paths directly, we effectively have a local model – in time. The purpose of surrogate models is to learn about the model, not the data. The former is much easier in MRF than in standard RF since the vector  $\beta_t$  fully characterizes the prediction at a particular point in time.<sup>52</sup> Moreover, the coefficients are attained to predictors that can have themselves a specific economic meaning. Considering this and the earlier discussion of section 2.1, it is natural in a macro time series context to fit surrogate models to time-varying parameters themselves – a blatant divide-and-conquer strategy.

### A.2.1 About $VI_{OOB}$ , $VI_{OOS}$ and $VI_\beta$

I now explain the motivation and mechanics behind the different VI measurements. The first measure,  $VI_{OOB}$ , is the standard out-of-bag (hence OOB) VI permutation measure widely used

---

<sup>52</sup>More generally, any partially linear model in the spirit of MRF has a potential for local surrogate analysis along the linear regression space rather than the observations line.

in RF applications (Wei et al., 2015). It consists of randomly permuting one feature  $S_j$  and comparing predictive accuracy to the full model on observations that were not used to fit the tree.<sup>53</sup> This pseudo evaluation set is convenient because it is a direct byproduct of the construction of the forest. Under a well-specified model that includes enough lags of  $y_t$ , autocorrelation of residuals will not be an issue. This condition is likely to be met here since the analysis focuses on results for  $h = 1$ .<sup>54</sup>  $VI_{OOS}$  considers a different testing set more natural for time series data: the real OOS, which in this section spans from 2007q2 to the end of 2014. By construction, this measure focuses on finding variables which contribution paid off during a specific forecasting experiment, rather than throughout the whole sample. This is not bad *per se* but is a different concept that can be of independent interest. Finally, both  $VI_{OOB}$  and  $VI_{OOS}$  focus on overall fit.  $VI_{\beta}$  implements the same idea as  $VI_{OOB}$  but is calculated using a different loss function. That is,  $VI_{\beta_{k,j}}$  reports a measure of how much the path of  $\beta_k$  is altered (out-of-bag) when variable  $S_j$  is randomly permuted in the forest part. Finally, I use the various VI measurements as devices to narrow down the set of predictors for the construction of intuitive trees.

I restrict the number of considered variables (for the next step) to be 20 for each VI criteria. When VI suggest that a parsimonious set of variables matter, it is very rarely more than 3 or 4 variables. Thus, restricting it to 20 is a constraint that only binds if all variables contribute, but marginally, in the spirit of a Ridge regression (Friedman et al., 2001). When it comes to that, the cut-off is simply the natural reflection of a trade-off between interpretability and fit.

### A.3 Further Investigation of the Importance of $S_t$

Do MAFs matter? For 3 standard ML models (standard RF, LASSO, Ridge) that can handle high-dimensional data sets, I investigate the usefulness of the MAFs advocated in section 2.6. The codes to describe the different information sets are

**CSF** : only standard cross-sectional factors (5 factors, 8 lags of them),

**MAF** :  $S_t$ ,

**ALL** :  $S_t$  + all the raw data (8 lags),

**X** : 8 lags of the raw data.

Figure 14 summarizes results over 18 targets (6 variables and the first 3 horizons). The first striking fact is that the four best models are RF, followed by the LASSO block, Ridge and FAAR.

<sup>53</sup>This is thought as the equivalent for a black-box model to setting a specific coefficient to 0 in a linear regression and then comparing fits. However, VI as implemented here (and in most applications) does not re-estimate the model after dropping  $S_j$ . This differs from a t-test since it is well known that the latter is equivalent to comparing two  $R^2$ 's – the original one and that of a re-estimated model, under the constraint.

<sup>54</sup>Notwithstanding, at longer horizons,  $VI_{OOB}$  could paint a distorted picture in the presence of autocorrelation – the same way K-fold cross validation can be inconsistent for time series data (Bergmeir et al., 2018). This worry can be alleviated by using a block approach like in section 2.7.

This suggests, with an unprecedented level of surprise, that models matter. For RFs, the best model is the one using  $S_t$  followed closely by the one that also adds the raw data to it. However, if we drop the MAFs, we incur a significant loss and obtain the worst of RFs, (so-called RF-X). The RF with cross-sectional factors only performs quite well in an unequal fashion.

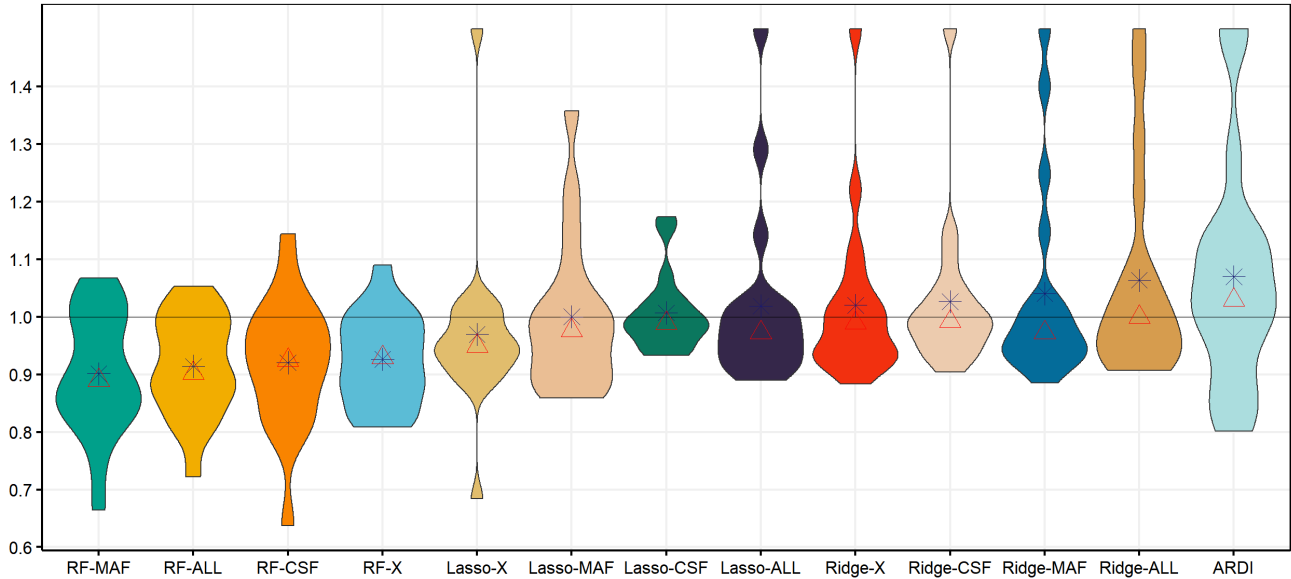


Figure 14: On the usefulness of  $S_t$ .

The best LASSO models must include the raw data. Models with either standard factors or MAFs only do not perform as well. This is not true for Ridge where the best model is the one that uses  $S_t$ . It is however important to note at this point of the ranking that these models are already lagging RFs in a significant way. The usefulness of MAFs is further studied in [Goulet Coulombe et al. \(2020a\)](#) and found to help, mostly with tree-based algorithms. However, it is found supplanted by a more computationally demanding (but more general) transformation of the raw data that [Goulet Coulombe et al. \(2020a\)](#) propose specifically for ML-based macroeconomic forecasting.

## A.4 On Tuning Parameters

The bulk of the discussion on the algorithm's specifics is deferred to the [R package](#). None of the RFs reported in the text were tuned. This is not heresy, as minuscule performance gains from doing so (like optimizing `mtry`) are the norm rather than the exception. Additionally, restraining the terminal nodes size can only alter performance very mildly and it is now clear why ([Goulet Coulombe, 2020b](#)). Nonetheless, reviewing some of those untuned tuning parameters can be insightful about MRFs inner workings.

- `RWR`: stands for Random Walk Regularization strength as discussed in 2.3. It is the  $\zeta$  in equation (2).
- `RL`: stands for Ridge Lambda ( $\lambda$ ) in equation (1). Prior means are OLS estimates.
- `Minimal Node Size`: Minimal parent leaf size to consider a new split. Set to 10 for quarterly data and 15 for monthly.
- `MLF`: stands for Minimum Leaf Fraction. It is the parameter in MRF that has a role complementary to that of minimum node size. The so-called "fraction" is the ratio of parameters in the linear part to that of observations in any node (which includes most importantly the terminal ones). Here is an example. Set  $\text{MLF} = 2$ , the linear part has 3 parameters, and we are trying to split a subset of 15 observations. This setting implies that any split that results in having less than 6 observations in the children node will not be considered. This specific setting ensures that the ratio of parameters to observations never exceeds  $1/2$  in any node. This ensure stability, especially if the two aforementioned HPs are set to 0. However, when `RWR` and `RL` are active, it is possible to consider  $\text{MLF} = 1$  or even lower, like for the large VARs specifications of section A.1. The extra regularization allows in the latter case to have base regressions that have parameters/observations ratio exceeding 1 (high-dimensional setting). This is very desirable in a quarterly macro setting because setting  $\text{MLF} > 2$  or higher seriously restricts the depth of the trees being grown.
- `mtry`: how many  $S_j$ 's do we consider as potential "splitter" at each split? It is easier to think about it as a fraction of the total number of predictors. For regression settings, the suggested value is  $1/3$ . The lower it gets, the more random tree generation gets, and better diversification may ensue. Moreover, `mtry` directly impacts computational burden. It is often found, in a macro context, that lowering `mtry` to 0.2 does not alter performance noticeably, while reducing appreciably computations. In fact, running RF-MAF with  $\text{mtry} \in \{0.1, 0.2, 0.33, 0.5\}$  delivers nearly identical performance for all variable/horizon pairs of the quarterly exercise. This is likely attributable to macro data having a factor structure. If  $S_j$  is "not available" for a split when it would in fact maximize fit locally, there is another strongly correlated  $S_{j'}$  ready for the task. For instance, if the unemployment rate is discarded by `mtry`, then there are more than 20 other labor indicators that can possibly substitute for it. If those 20 variables are all a noisy representation of the same latent variable the model wants to split on, then the probability of having none to split with at a given point is  $\left(1 - \frac{\text{mtry}}{\#\text{regressors}}\right)^{20} \approx 0$ .
- `Trend Push`: Some minorities may end up being underrepresented as a result of `mtry`'s discriminating action. While there are 20+ labor indicators in the data base, there is only one trend. Since exogenous change should most certainly not be underrepresented, its

"personalized" probability of inclusion can be pushed beyond what `mtry` suggests.

- `Subsampling Rate`: is set at 75%.

A scaled down quarterly forecasting exercise was conducted to see whether tuning any of those could help. Precisely, horizons 1, 2, and 4 quarters were considered and models (ARRF,FA-ARRF,VARRF) were estimated once at the beginning of the OOS period (2002). Tuning parameters were optimized targeting 1998-2002 data as an artificial test set. Possible values were  $RWR \in \{0, 0.5, 0.95\}$ ,  $RL \in \{0.1, 0.5\}$ ,  $mtry \in \{0.2, 0.33, 0.5\}$  and  $min.node.size \in \{10, 40\}$ . It is found that results are largely invariant to pre-optimized HPs. As mentioned earlier, what matters most in the linear part. It is observed that optimizing tuning parameters can help reduce marginally RMSEs of MRFs that were sometimes struggling (like VARRF). Results are available upon request.

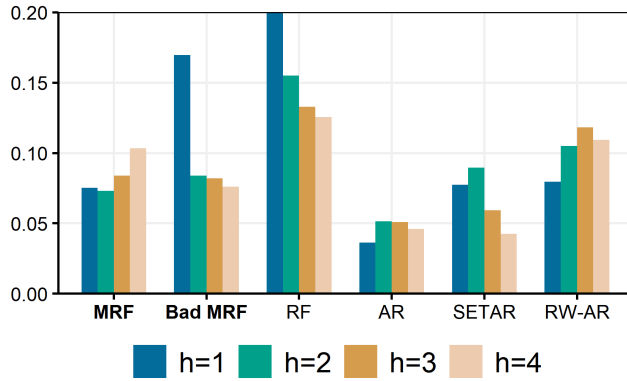
---

**Algorithm 1** How the key tuning parameters enter MRF, and other practical aspects

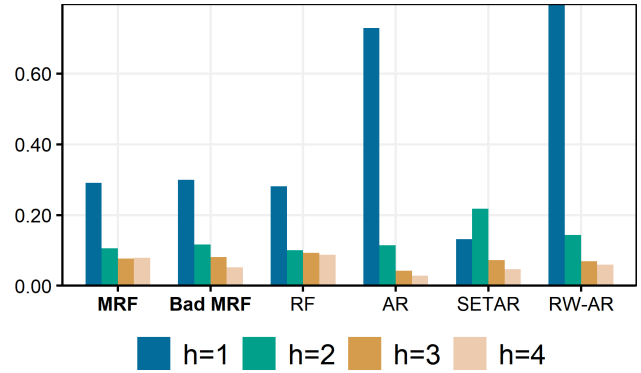
---

- 1: Draw blocks of some size (8 for quarterly, 24 monthly), that makes for `Subsampling Rate`% of the sample. To simply get the mean prediction, 100 trees are usually more than enough. To get credible regions to stabilize, 200-300 trees are typically needed.
  - 2:
    - For each subsample: run (2) recursively on that sample given  $\lambda$  and  $\zeta$  values until each (potential) parent nodes are smaller than `Minimal Node Size`.
    - A total of `mtry` predictors are considered at each splitting step  $\mathcal{J}^-$  is randomly picked out of  $\mathcal{J}$ . Those probabilities are all  $1/\dim(\mathcal{J})$  by default. `Trend Push` pushes that of the trend further if judged appropriate for a given data set.
    - When evaluating potential splits, discard those that would not meet MRF's requirements on resulting children nodes.
    - This outputs one tree structure  $\mathcal{T}$ .
  - 3: When inputted with new observations of  $X_t$  and  $S_t$ , each tree produces a forecast. MRF forecast is the mean of the those.
  - 4: Same goes for  $\beta_t$ : each tree predicts its own  $\beta_t$  out-of-sample and the posterior mean is the average of all those.
  - 5: In-sample  $\beta_t$ 's need an extra step: only draws that did not use observation  $t$  to construct the tree (that is, for which  $t$  was left out of the subsample) are used to characterize the distribution of  $\beta_t$ .
-

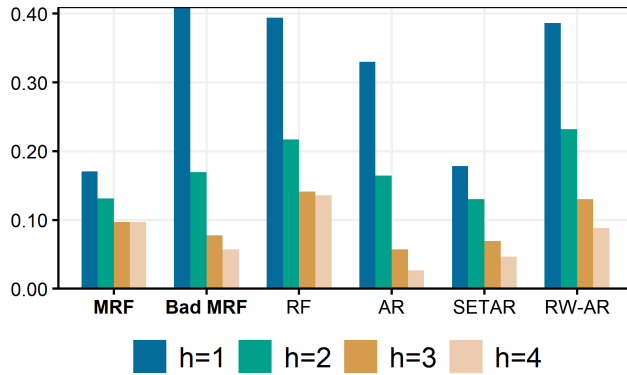
## A.5 Additional Figures and Tables



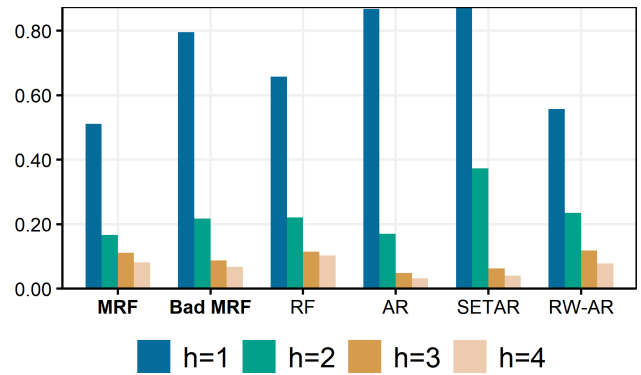
(a) DGP 1



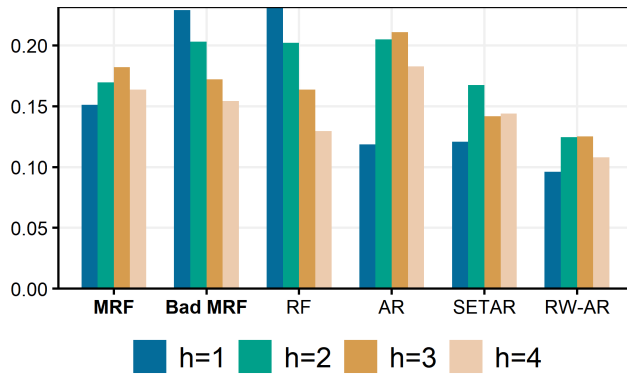
(b) DGP 2



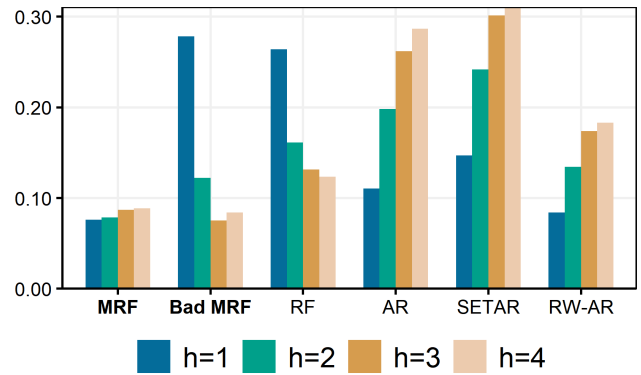
(c) DGP 3



(d) DGP 4



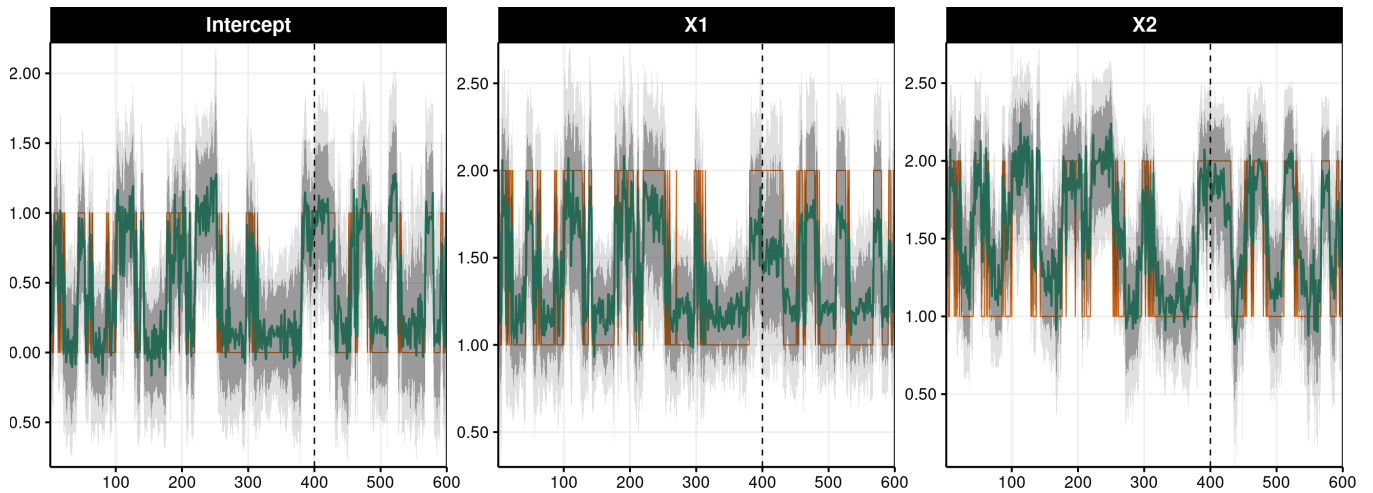
(e) DGP 5



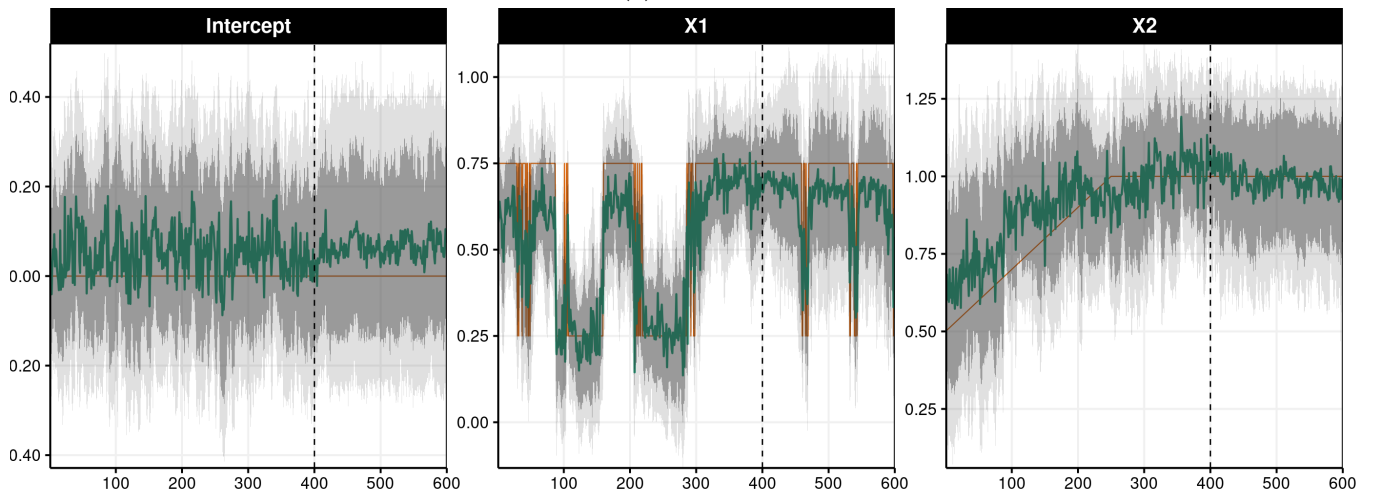
(f) DGP 6

Figure 15: Investigation of the consequences of  $X_t$ 's misspecification, as exemplified by "Bad ARRF". Instead of the first two lags of  $y_t$ ,  $X_t$  is replaced by randomly generated *iid* (normal) variables. Total number of simulations is 50, and the total number of squared errors is thus 2000.

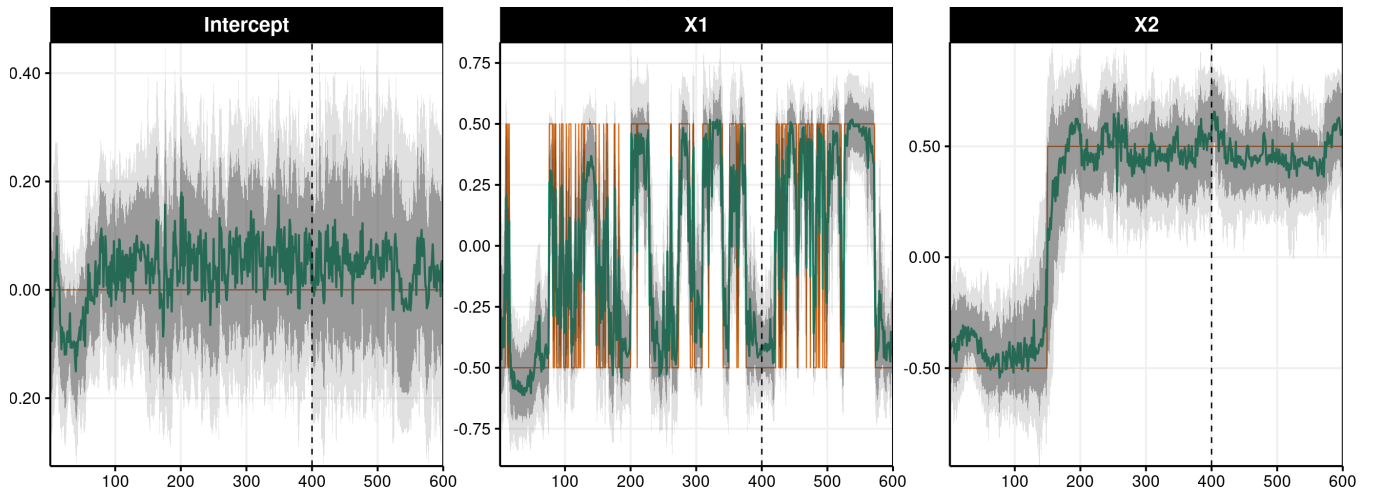




(a) DGP 1

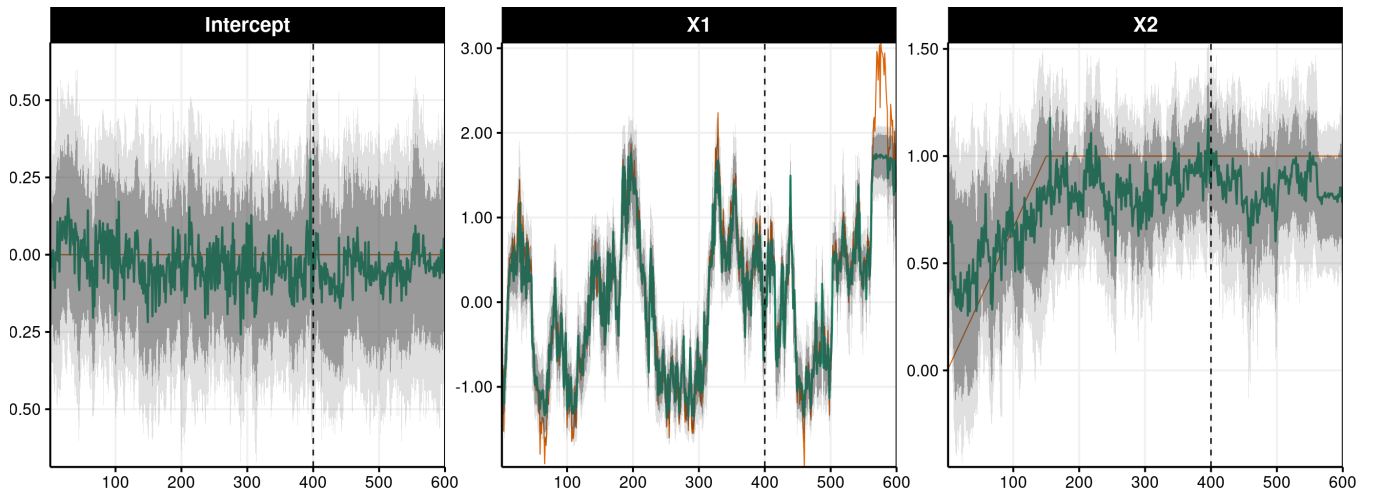


(b) DGP 2

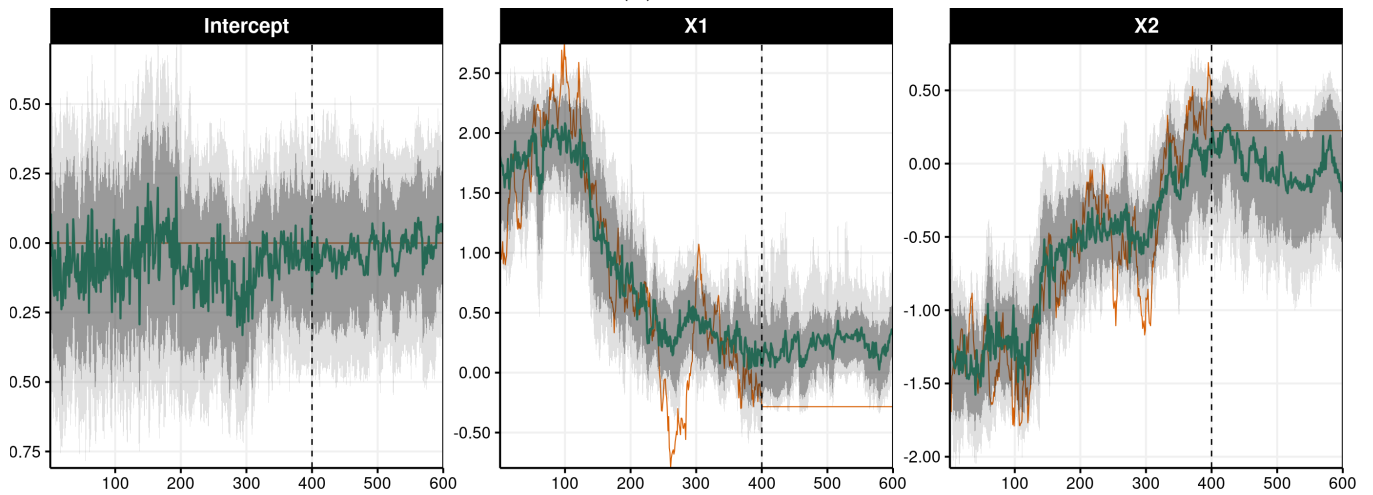


(c) DGP 3

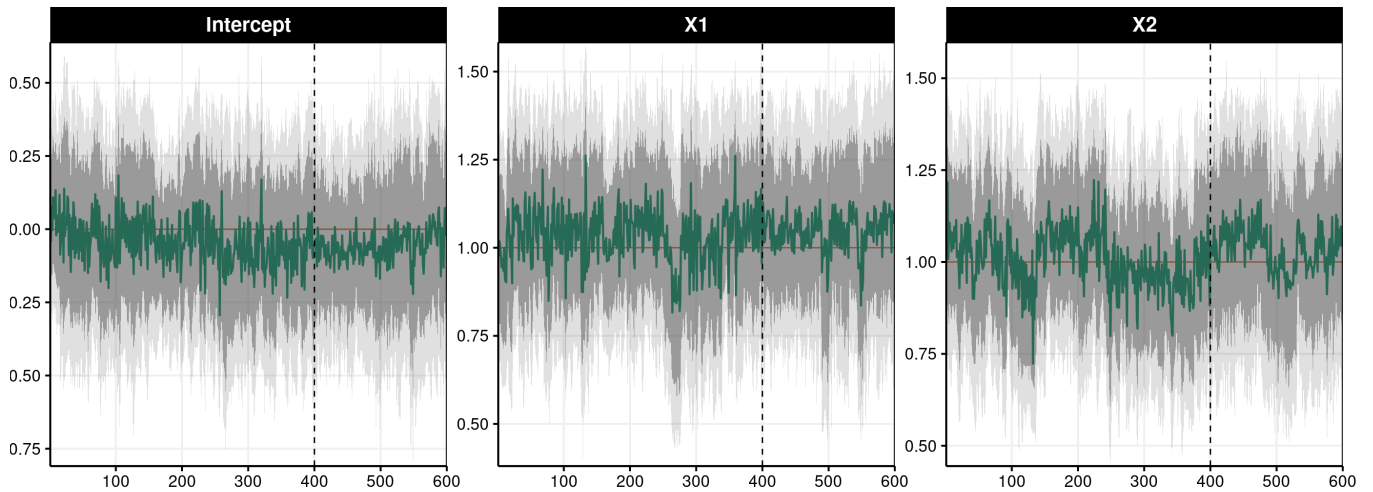
Figure 16: The grey bands are the 68% and 90% credible region. After the blue line is the hold-out sample. Green line is the posterior mean and orange is the truth. The plots include only the first 400 observations for visual convenience.



(d) DGP 4



(e) DGP 5



(f) DGP 6

Figure 16: (Continued) The grey bands are the 68% and 90% credible region. After the blue line is the hold-out sample. Green line is the posterior mean and orange is the truth. The plots include only the first 400 observations for visual convenience.

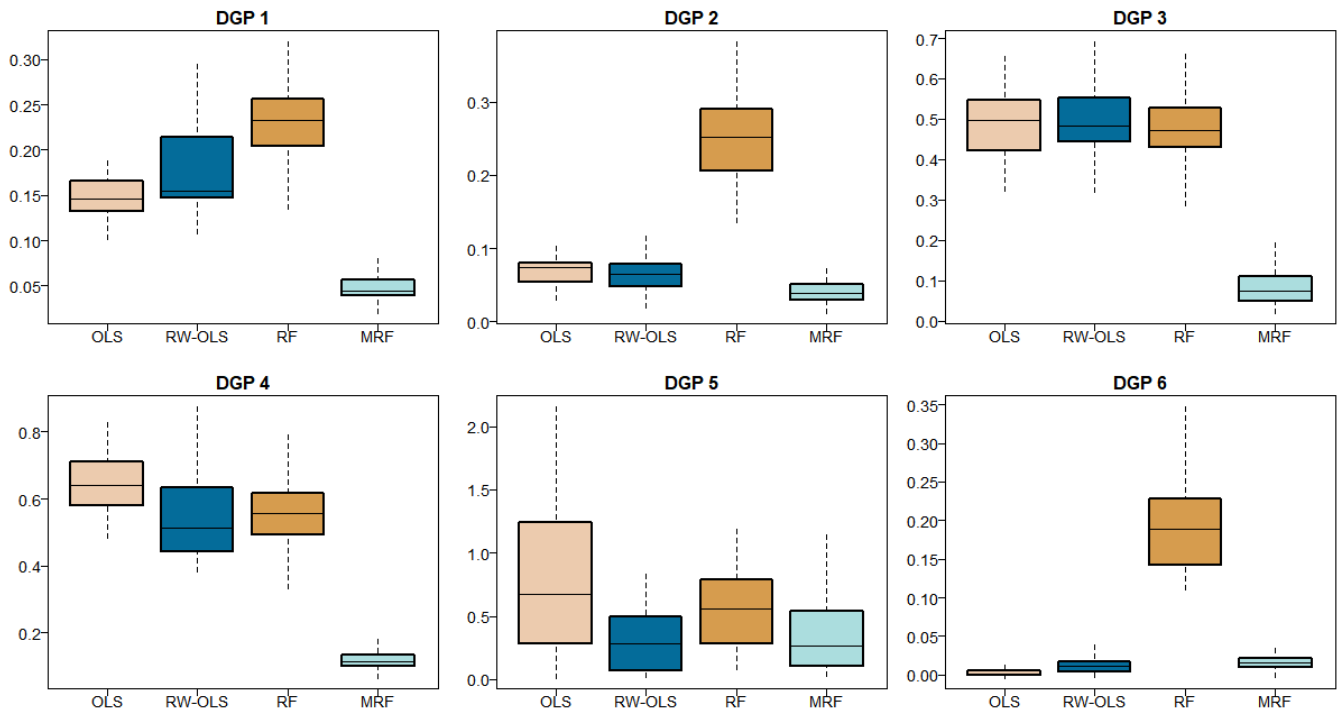
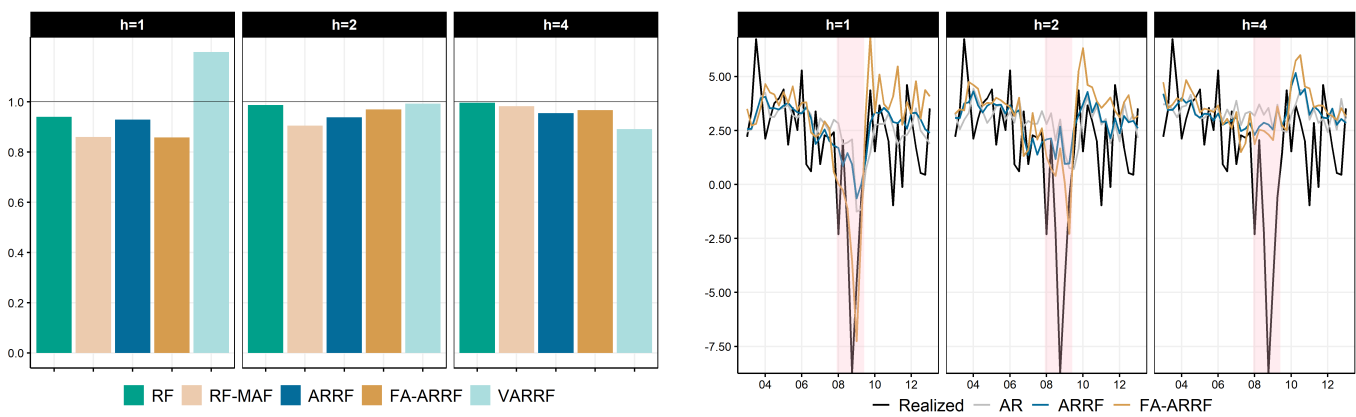


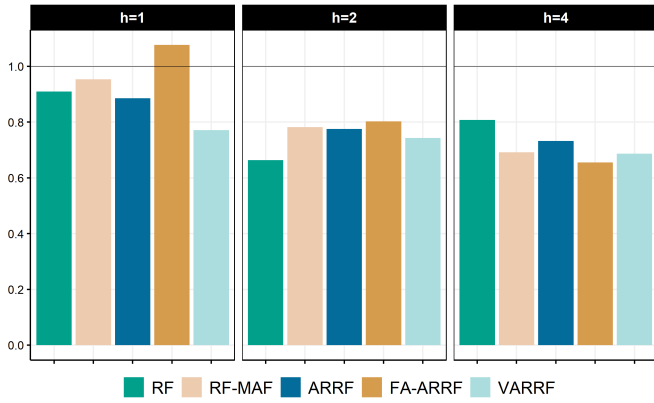
Figure 17: The distribution of RMSE dis-improvements with respect to the oracle's forecast for 4 models: OLS, Rolling-Window OLS, plain RF, MRF. 50 simulations of 750 OOS forecasts each.



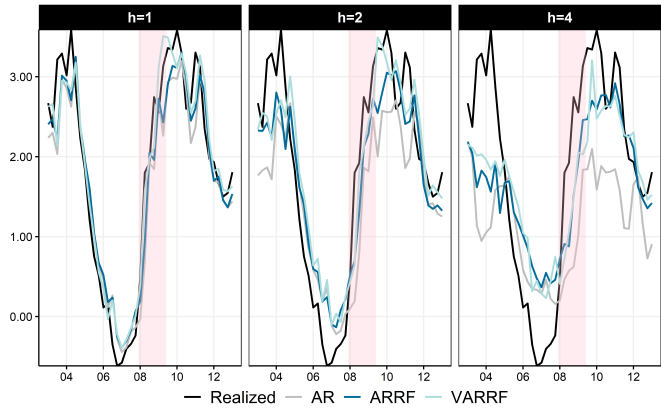
(a)  $RMSE_{GDP,h,m} / RMSE_{GDP,h,AR}$

(b) A look at forecasts

Figure 18: GDP results in detail

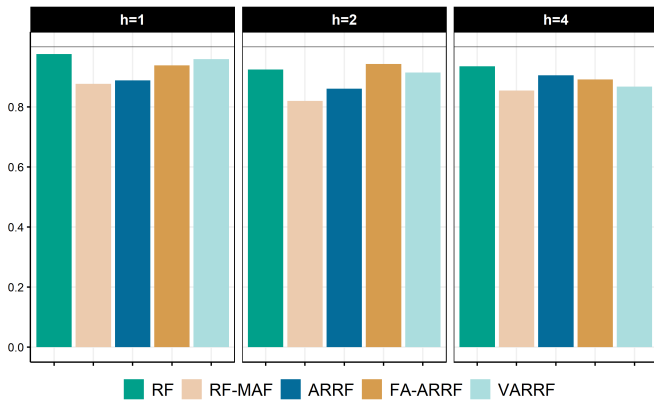


(a)  $RMSE_{SPREAD,h,m} / RMSE_{SPREAD,h,AR}$

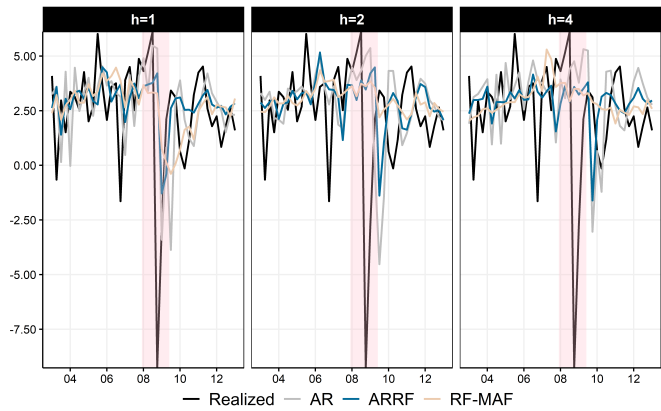


(b) A look at forecasts

Figure 19: SPREAD results in detail



(a)  $RMSE_{INF,h,m} / RMSE_{INF,h,AR}$



(b) A look at forecasts

Figure 20: INF results in detail

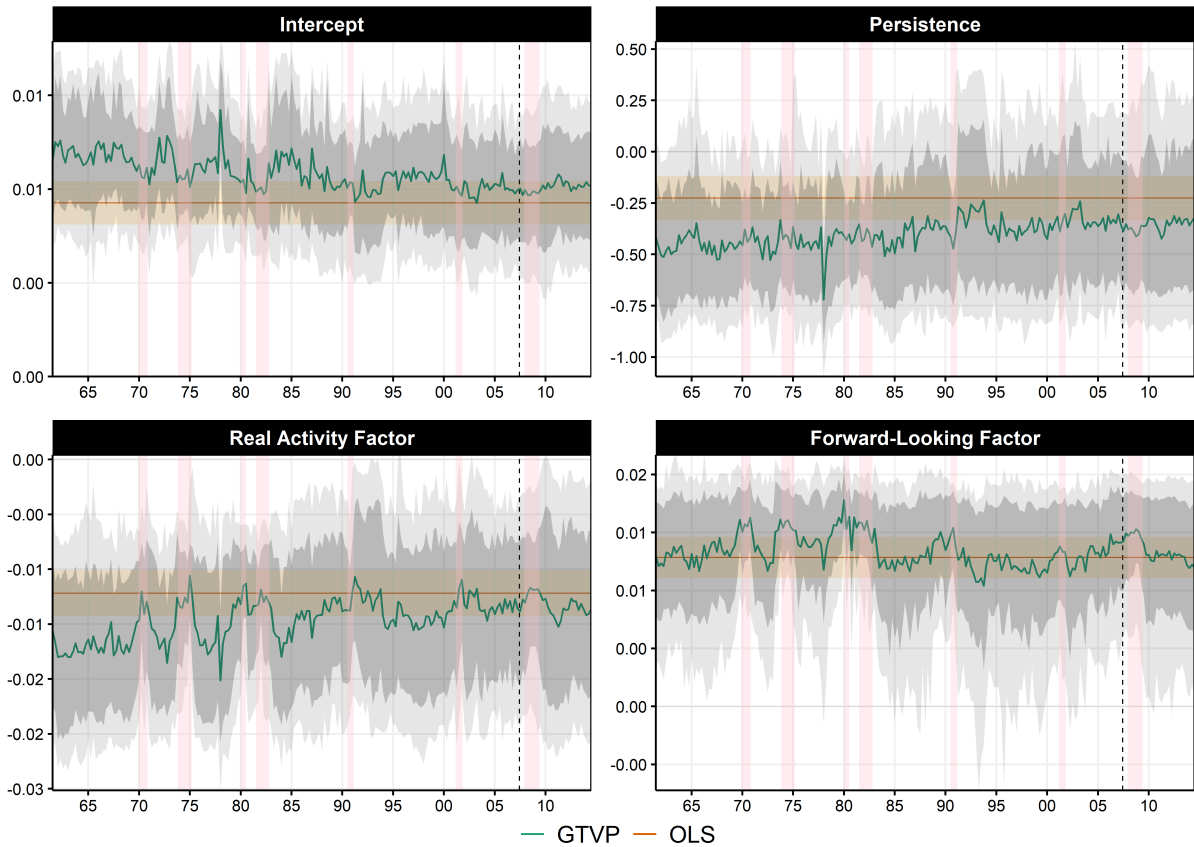


Figure 21: GTVPs of the one-quarter ahead GDP forecast. Persistence is defined as  $\phi_{1,t} + \phi_{2,t}$ . The grey bands are the 68% and 90% credible region. The pale orange region is the OLS coefficient  $\pm$  one standard error. The vertical dotted blue line is the end of the training sample. Pink shading corresponds to NBER recessions.

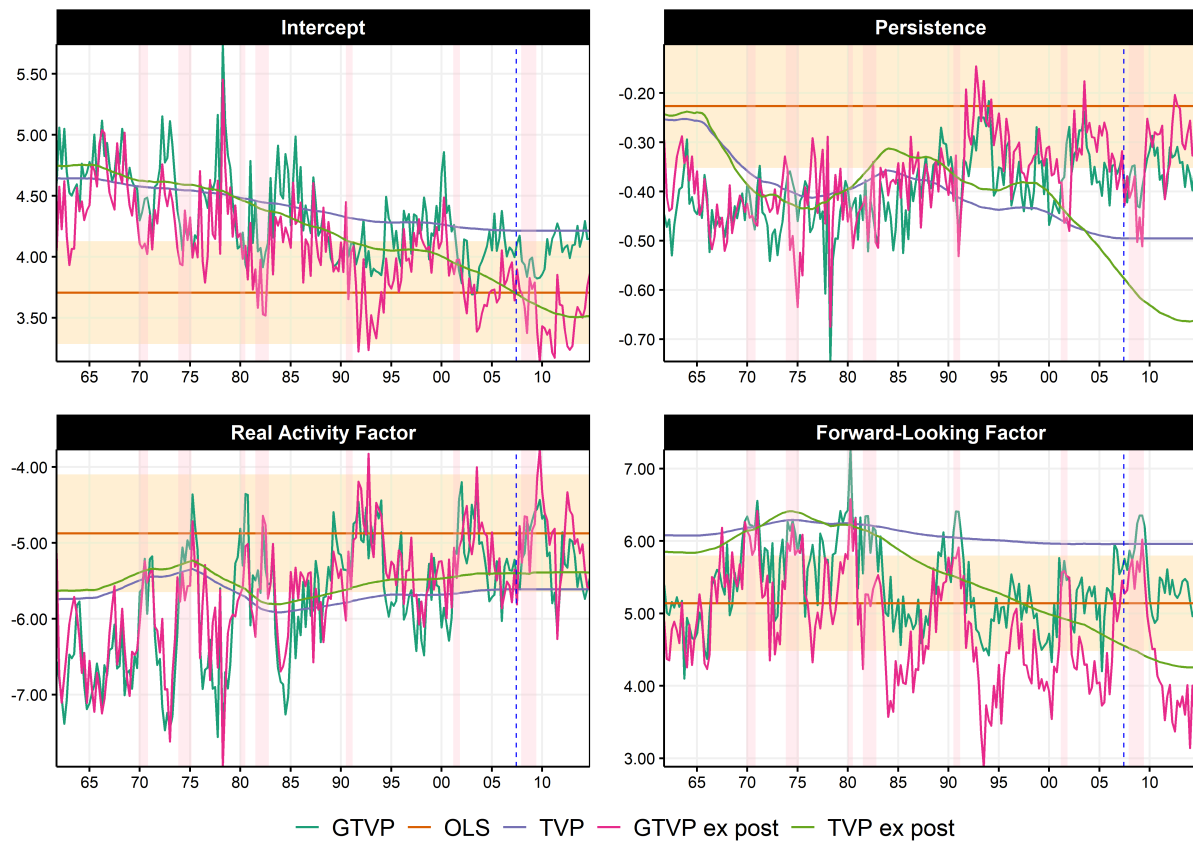
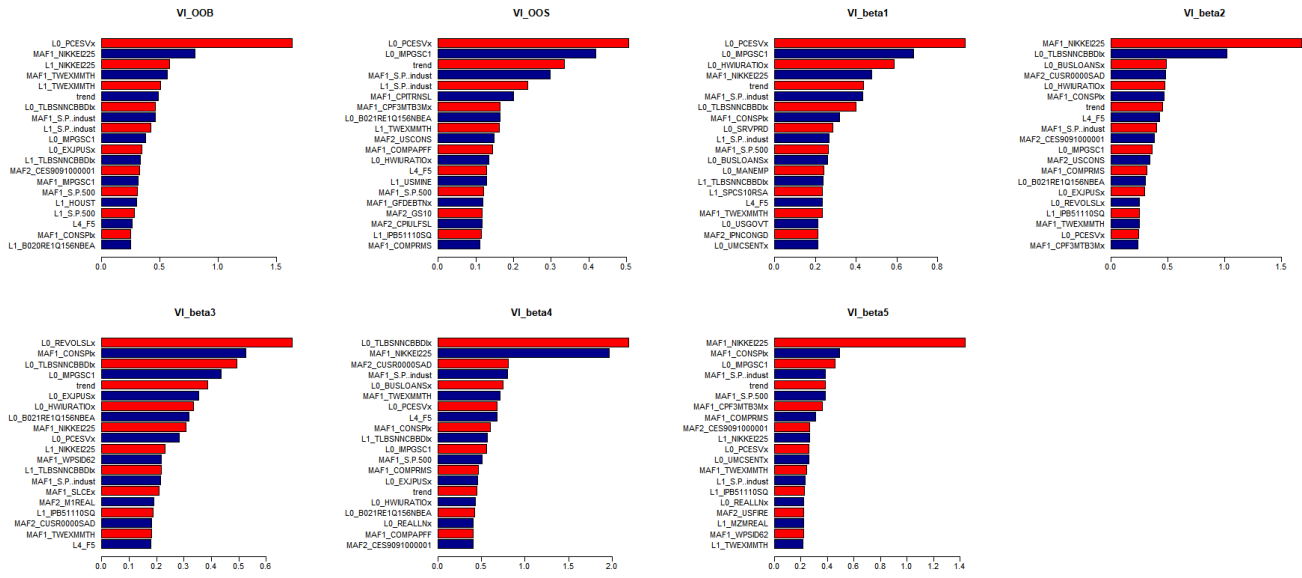
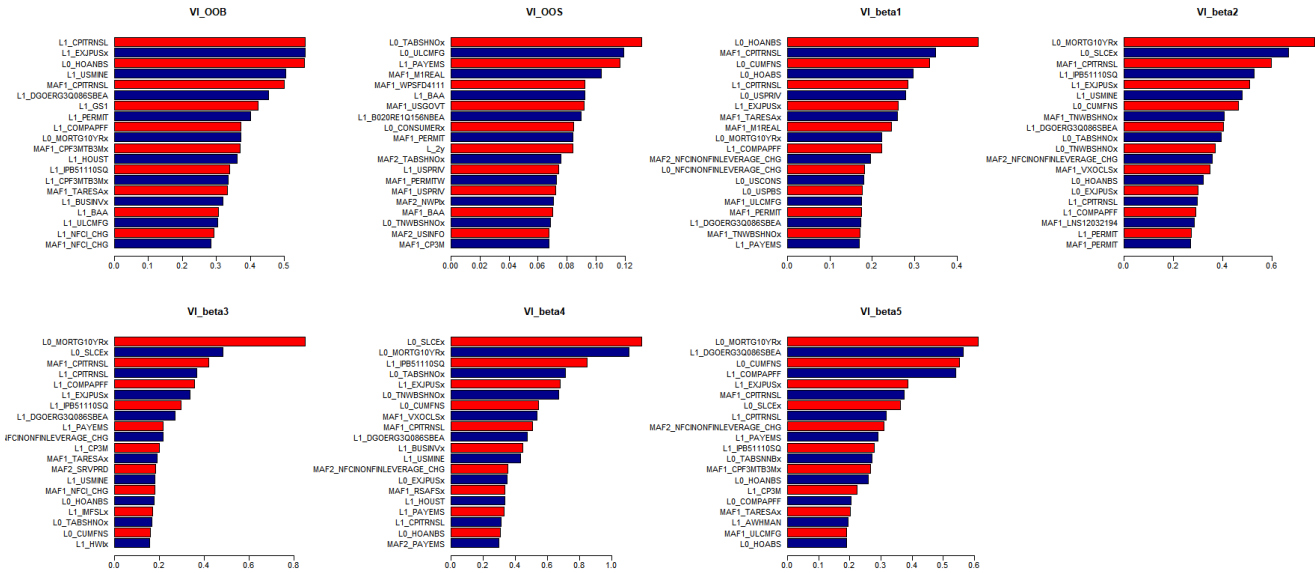


Figure 22: GDP equation  $\beta_t$ 's obtained with different techniques. Persistence is defined as  $\phi_{1,t} + \phi_{2,t}$ . TVPs estimated with a ridge regression as in Goulet Coulombe (2020a) and the parameter volatility is tuned with k-fold cross-validation. Ex Post TVP means using the full sample for estimation and tuning as opposed to only using pre-2002 data as for GTVPs. The pale orange region is the OLS coefficient  $\pm$  one standard error. Pink shading corresponds to NBER recessions.

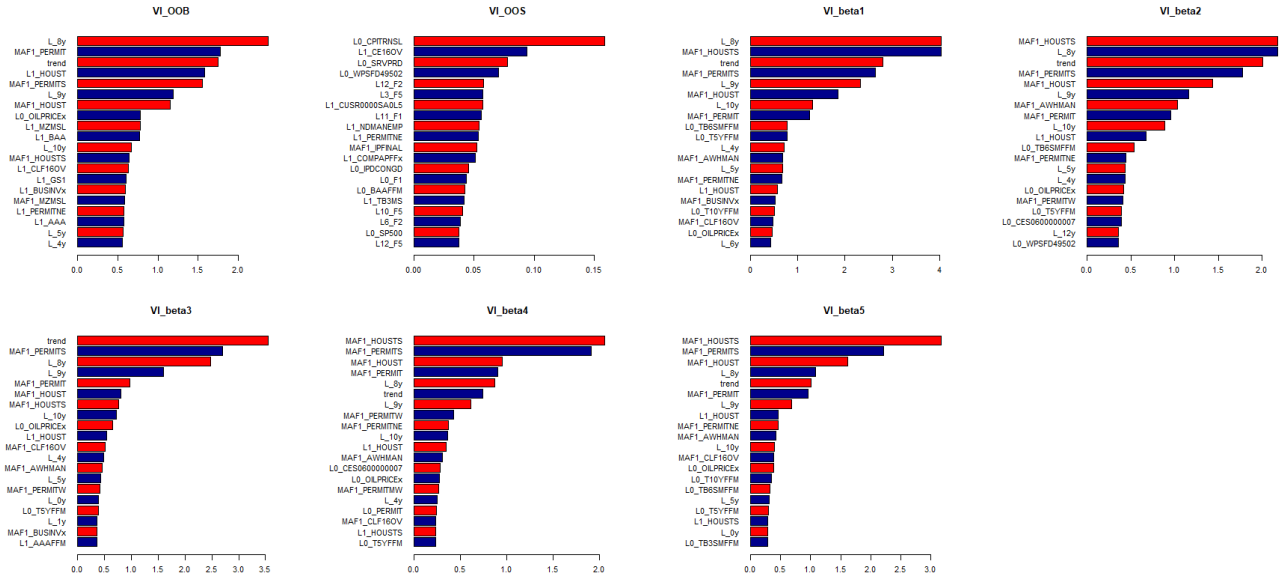


(a) GDP horizon 1

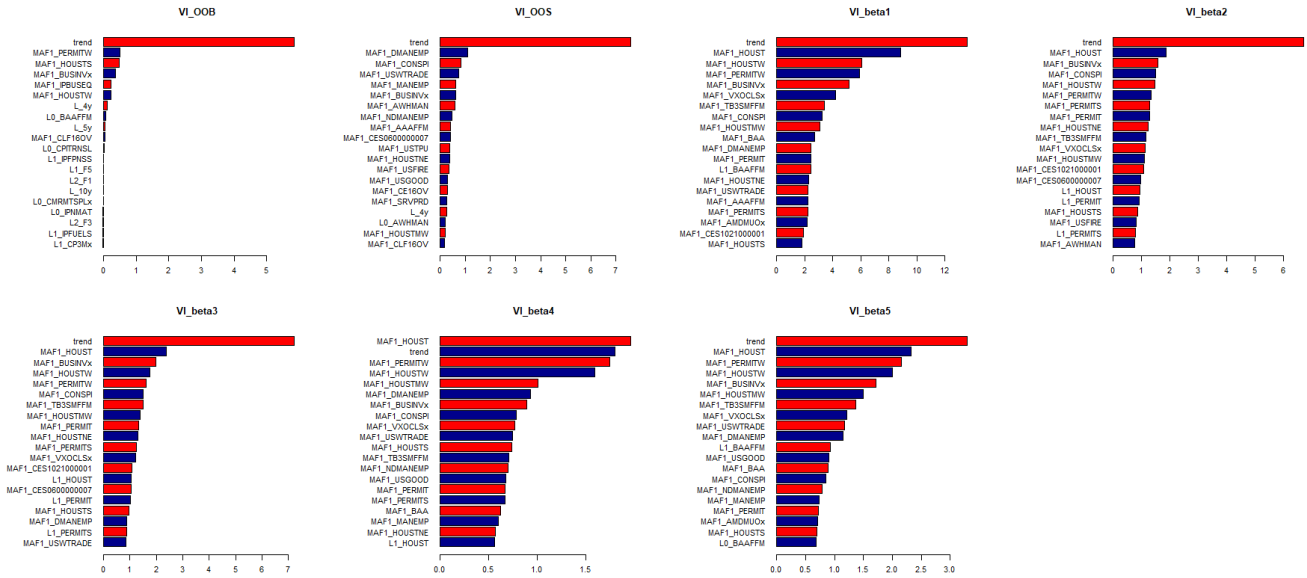


(b) UR horizon 1

Figure 23: 20 most important series according to the various variable importance (VI) criteria. Units are relative RMSE gains (in percentage) from including the specific predictor in the forest part.  $VI_{OOB}$  means VI for the out-of-bag criterion.  $VI_{OOS}$  is using the hold-out sample.  $VI_{\beta}$  is an out-of-bag measure of how much  $\beta_{t,k}$  varies by withdrawing a certain predictor.



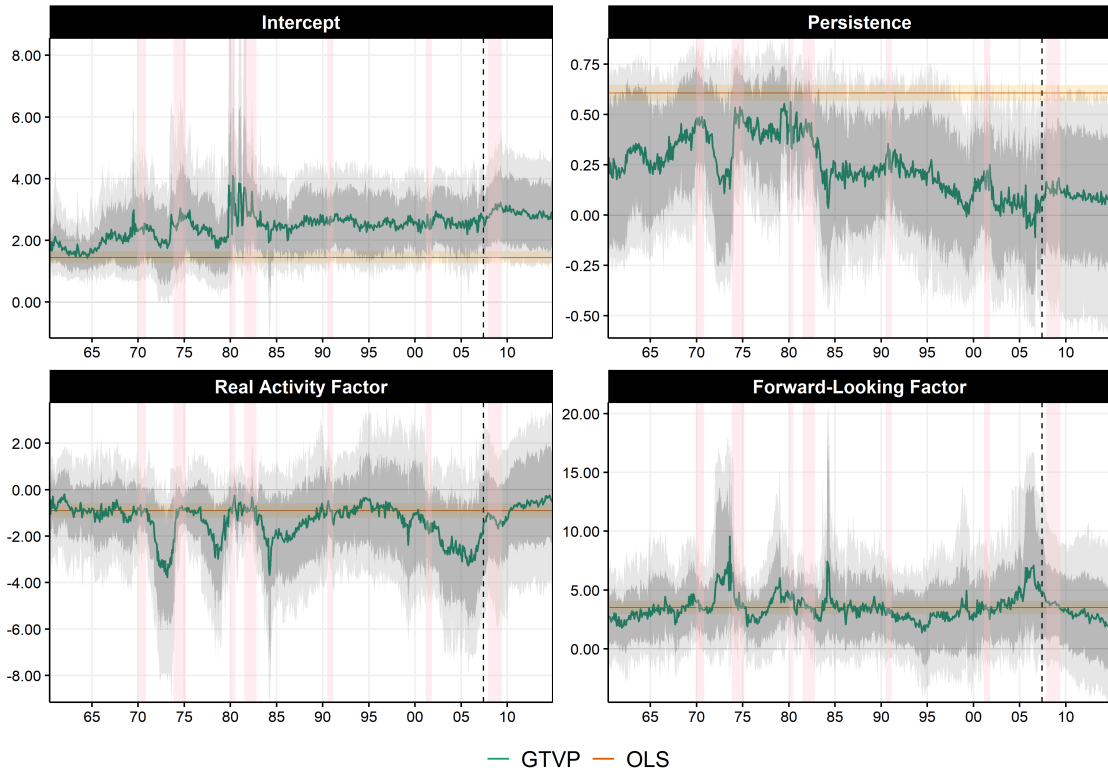
(a) One month ahead inflation forecast



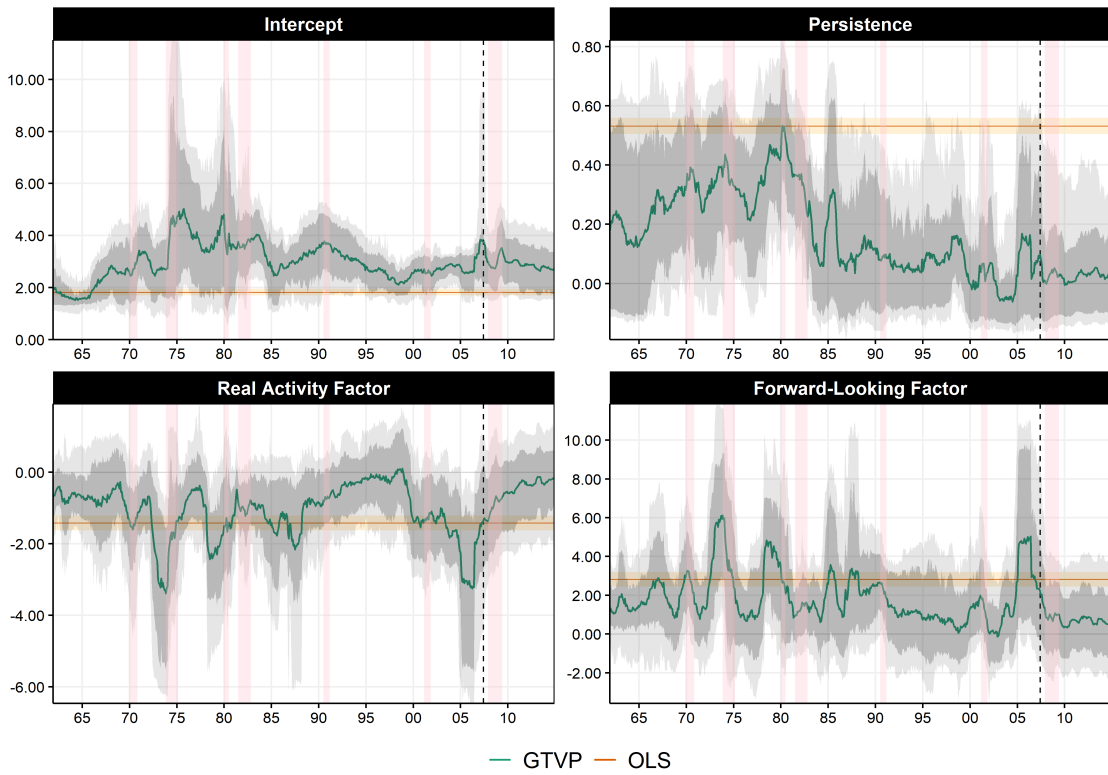
(b) Average inflation over the next 12 months

Figure 24: 20 most important series according to the various variable importance (VI) criteria. Units are relative RMSE gains (in percentage) from including the specific predictor in the forest part.  $VI_{OOB}$  means VI for the out-of-bag criterion.  $VI_{OOS}$  is using the hold-out sample.  $VI_{\beta}$  is an out-of-bag measure of how much  $\beta_{t,k}$  varies by withdrawing a certain predictor.



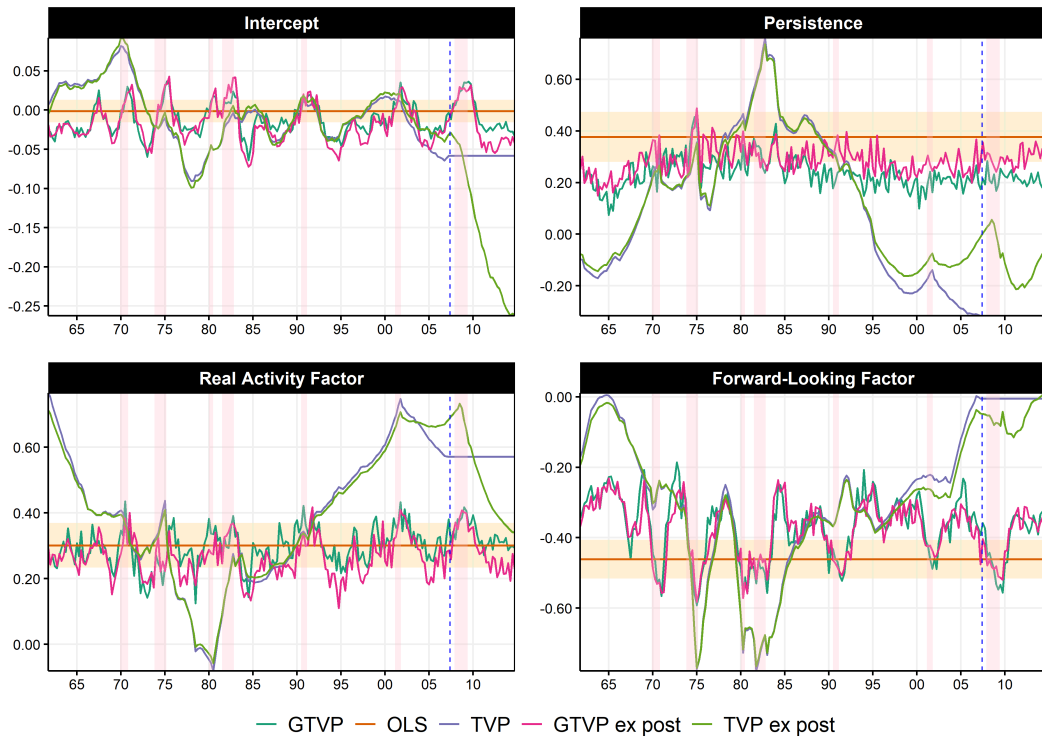


(a) One-month ahead

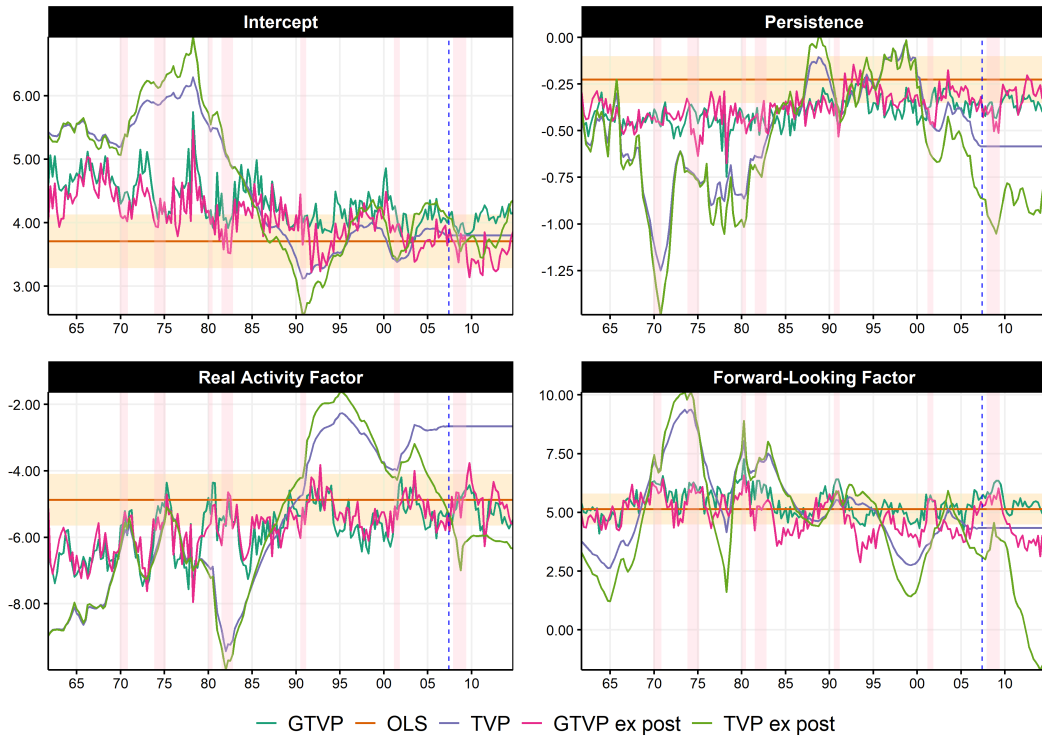


(b) 12-months ahead

Figure 25: GTVPs of monthly inflation forecast. The grey bands are the 68% and 90% credible regions. The pale orange region is the OLS coefficient  $\pm$  one standard error. The vertical dotted line is the end of the training sample. Pink shading corresponds to NBER recessions.



(a) UR equation



(b) GDP equation

Figure 26:  $\beta_t$ 's obtained with different techniques. TVPs estimated with a ridge regression as in Goulet Coulombe (2020a) and the parameter volatility  $\lambda$  is tuned with k-fold cross-validation, then divided by 100. This means the standard deviation of parameters shocks is allowed to be about 10 times higher than what CV recommends. Ex Post TVP means using the full sample for estimation and tuning as opposed to only using pre-2002 data as for GTVPs. The pale orange region is the OLS coefficient  $\pm$  one standard error.

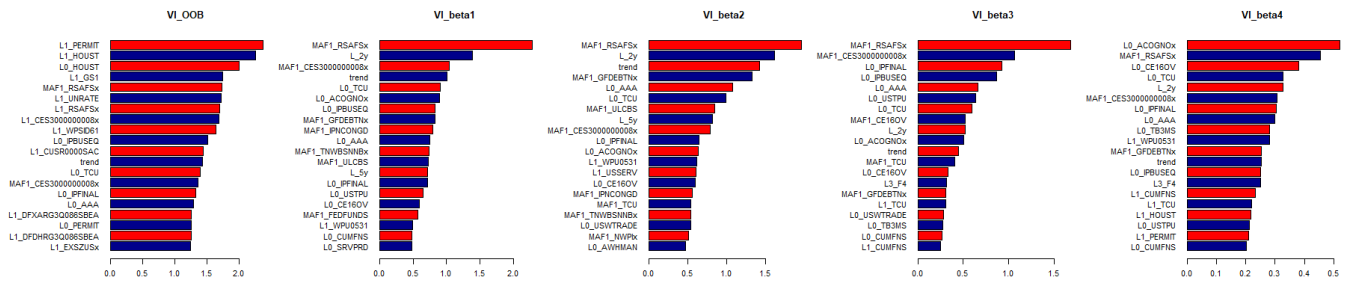


Figure 27: 20 most important series according to the various variable importance (VI) criteria. Units are relative RMSE gains (in percentage) from including the specific predictor in the forest part.  $VI_{OOB}$  means VI for the out-of-bag criterion.  $VI_{OOS}$  is using the hold-out sample.  $VI_{\beta}$  is an out-of-bag measure of how much  $\beta_{t,k}$  varies by withdrawing a certain predictor.

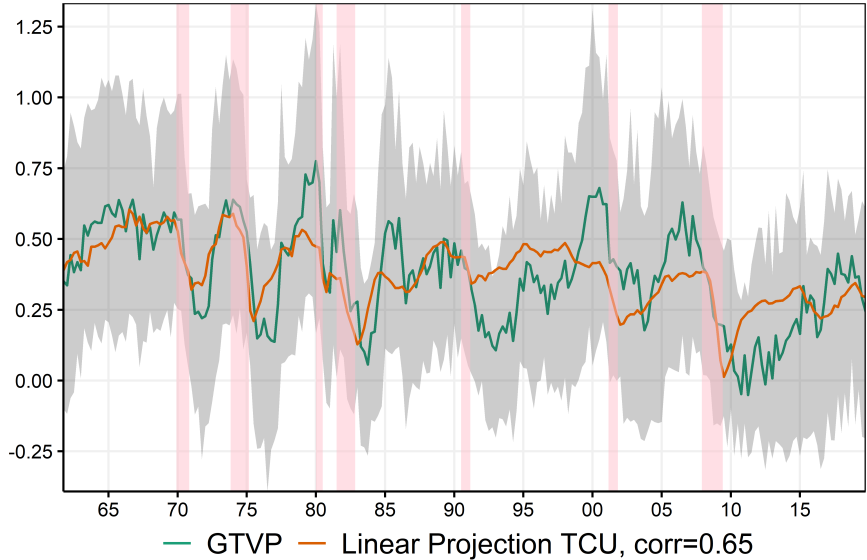


Figure 28:  $\beta_{3,t}$  in (6) with additional controls for supply and monetary policy shocks. Capacity Utilization is still substantially correlated with the inflation-unemployment trade-off. The grey band is the 68% credible region. Pink shading corresponds to NBER recessions.

Table 4: Main Quarterly Results

|               | FA-AR       | LASSO-MAF     | Ridge-MAF | RF            | RF-MAF        | AR+RF         | Tiny RF | FA-ARRF       | ARRF        | Tiny ARRF    | VARRF         | SETAR       | STAR        | TV-AR  |
|---------------|-------------|---------------|-----------|---------------|---------------|---------------|---------|---------------|-------------|--------------|---------------|-------------|-------------|--------|
| <b>GDP</b>    |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 1.02        | 0.96          | 0.89**    | 0.94          | 0.86          | 0.89          | 1.03    | <b>0.86</b>   | 0.93        | 1.04         | 1.20          | 1.01        | 1.03        | 0.99   |
| h=2           | 0.96        | 0.98          | 0.98      | 0.99          | <b>0.91</b>   | 0.93          | 1.01    | 0.97          | 0.94**      | 1.03         | 0.99          | 0.97        | 0.98        | 1.03   |
| h=4           | 1.03        | 0.98          | 0.99***   | 1.00          | 0.98          | 0.99          | 1.03    | 0.97          | 0.95        | 0.98         | <b>0.89</b>   | 0.97***     | 0.96***     | 0.96   |
| h=6           | 1.36        | 0.98          | 0.98      | 0.98          | 1.00          | 1.00          | 1.08    | 1.01          | 0.97        | 0.98         | 1.00          | 0.98        | <b>0.95</b> | 0.98   |
| h=8           | 1.37        | 1.00          | 0.99      | 0.99          | 0.99          | <b>0.96</b>   | 1.15    | 1.06          | 1.00        | 1.01         | 1.04***       | 1.00        | 0.97        | 1.00   |
| <b>UR</b>     |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 0.83        | 0.99          | 0.99      | 1.00          | 0.85*         | 0.84          | 1.24**  | <b>0.72</b>   | 0.90***     | 1.00         | 1.24          | 1.18        | 1.10        | 1.00   |
| h=2           | 0.80        | 0.98          | 0.92*     | 0.98          | 0.85          | 0.84          | 1.15*   | <b>0.76</b>   | 0.90        | 0.96         | 0.89          | 1.03        | 0.97        | 0.99   |
| h=4           | 0.88        | 0.96***       | 0.94**    | 0.96*         | 0.87*         | 0.84*         | 1.37    | <b>0.79</b>   | 0.87        | 0.92         | 0.91          | 1.02        | 1.01        | 1.34   |
| h=6           | 1.18*       | 0.98          | 0.98      | 1.01          | 0.94          | 0.90          | 1.60*   | <b>0.89</b>   | 0.95        | 0.97         | 0.95          | 1.07        | 1.04        | 1.14   |
| h=8           | 1.25        | 0.98          | 1.01      | 1.01          | <b>0.95</b>   | 0.95          | 1.57    | 1.01          | 0.98        | 0.98         | 1.04          | 1.09        | 1.06        | 1.11** |
| <b>SPREAD</b> |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 1.28        | 2.16***       | 0.93      | 0.91          | 0.95          | 0.79**        | 0.96    | 1.08          | 0.89**      | 1.06         | <b>0.77**</b> | 1.51***     | 1.53***     | 0.98   |
| h=2           | 1.13        | 1.20          | 0.77      | <b>0.66**</b> | 0.78          | 0.72***       | 0.93    | 0.80          | 0.78**      | 1.11         | 0.74**        | 1.19        | 1.20        | 1.04   |
| h=4           | 0.86        | 0.95          | 1.01      | 0.81          | 0.69**        | <b>0.61**</b> | 1.48*   | 0.66**        | 0.73**      | 1.07         | 0.69**        | 1.04        | 1.06        | 1.30   |
| h=6           | 1.51        | 0.80*         | 1.13      | 0.98          | 0.80          | 0.80          | 1.43    | <b>0.72**</b> | 0.82        | 1.05         | 0.74*         | 1.03        | 1.06        | 1.19   |
| h=8           | 1.28        | <b>0.76**</b> | 0.96      | 0.92          | 0.83          | 0.89          | 1.36    | 0.82          | 0.88        | 0.99         | 0.85          | 1.11        | 1.14        | 0.99   |
| <b>INF</b>    |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 1.01        | 0.93          | 0.95      | 0.98          | 0.88          | 1.23          | 0.90    | 0.94          | 0.89        | <b>0.87*</b> | 0.96          | 1.05        | 1.00        | 0.93   |
| h=2           | 1.01        | 0.96          | 0.92      | 0.92          | <b>0.82</b>   | 1.00          | 0.88    | 0.94          | 0.86        | 0.87         | 0.91          | 0.86*       | 0.86        | 0.89   |
| h=4           | 1.08        | 0.92          | 0.87      | 0.94          | <b>0.85**</b> | 0.96          | 0.86    | 0.89          | 0.91*       | 0.95*        | 0.87*         | 0.90*       | 0.87*       | 0.91   |
| h=6           | 1.32        | 0.96          | 0.90      | 1.01          | 0.88          | 1.00          | 0.86    | 0.91          | <b>0.85</b> | 0.92**       | 0.87          | 0.94        | 0.89        | 0.98   |
| h=8           | 1.21        | 0.98          | 1.27      | 1.44          | <b>0.88*</b>  | 0.94          | 0.88    | 0.91*         | 0.92        | 0.94         | 0.91*         | 0.96        | 0.92        | 0.98   |
| <b>HOUST</b>  |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 1.13        | 1.04          | 0.94*     | <b>0.92*</b>  | 1.00          | 1.01          | 1.24*** | 1.08          | 0.94**      | 0.95         | 1.09          | 1.01        | 0.99        | 1.00   |
| h=2           | 1.13        | 0.99          | 0.94**    | 0.95*         | 1.01          | 1.02          | 1.10*   | 1.06          | 1.00        | 1.02         | 0.99          | <b>0.94</b> | 0.97        | 1.01   |
| h=4           | 1.11        | 0.98**        | 0.97*     | 0.97          | 1.01          | 1.03          | 1.12    | 1.02          | 1.00        | 1.02         | 1.02          | <b>0.95</b> | 0.96        | 1.08   |
| h=6           | 1.40        | 0.96          | 0.96      | 0.96          | 0.96***       | 1.01          | 1.16    | 0.97***       | 0.99        | 1.00         | 0.98          | <b>0.95</b> | 0.96        | 0.99   |
| h=8           | 1.04        | 0.95          | 0.95      | 0.95          | 0.99          | 1.02          | 1.44    | 0.96          | 0.99        | 1.01         | 1.00          | <b>0.95</b> | 0.95        | 1.03   |
| <b>IR</b>     |             |               |           |               |               |               |         |               |             |              |               |             |             |        |
| h=1           | 1.85        | 1.02          | 1.55      | 1.17          | 1.11          | 0.97          | 0.99    | 1.29          | 0.94        | <b>0.92</b>  | 1.43          | 1.39        | 1.20        | 0.97   |
| h=2           | 1.49        | 0.96          | 1.01      | 1.00          | 0.93          | 0.98          | 1.29*** | 1.22          | 0.93        | <b>0.92</b>  | 1.10          | 1.15        | 1.11        | 1.04   |
| h=4           | <b>0.96</b> | 1.00          | 1.03      | 1.03          | 1.04          | 0.99          | 1.39*   | 0.99          | 0.97        | 1.12         | 0.97          | 1.08        | 1.07        | 1.09   |
| h=6           | 1.87        | 0.95          | 0.99      | 1.00          | <b>0.93</b>   | 0.93          | 1.23*   | 0.98          | 0.95*       | 1.07         | 1.12          | 1.19        | 1.14        | 1.06** |
| h=8           | 1.58        | 0.98          | 1.02      | 1.03          | <b>0.96</b>   | 0.96          | 1.20    | 1.04          | 0.96        | 1.10         | 0.98          | 1.25**      | 1.20**      | 1.06   |

Notes: This table report the root MSPE of the model  $m$  with respect to the root MSPE the AR(4). Best forecast of the row is in bold. Diebold-Mariano test is conducted for each model against the AR(4). "\*\*", "\*\*\*" and "\*\*\*\*" means p-values of below 10%, 5% and 1%.

Table 5: Monthly Results

|               | AR4           | AO-12       | AO-h    | FAAR   | RF          | RF-MAF         | AR+RF          | ARRF        | FA-ARRF        | Tiny ARRF     | VARRF         |
|---------------|---------------|-------------|---------|--------|-------------|----------------|----------------|-------------|----------------|---------------|---------------|
| <b>IP</b>     |               |             |         |        |             |                |                |             |                |               |               |
| h=1           | 1.00          | 1.11*       | 1.14    | 0.96   | 1.03        | <b>0.94*</b>   | 0.97           | 0.99        | 0.96           | 1.02          | 1.02          |
| h=3           | 1.02          | 1.17*       | 1.02    | 0.99   | 1.12        | 0.98           | <b>0.96</b>    | 1.03        | 1.01           | 1.02          | 1.08          |
| h=9           | 1.01          | 1.04        | 1.03    | 1.06   | 1.02        | 1.06           | 1.02           | 1.04        | 1.10           | 1.09          | 1.03          |
| h=12          | 1.01          | 1.00        | 1.00    | 1.05   | 0.99        | 0.97           | <b>0.91</b>    | 0.97        | 1.05           | 1.13          | 0.96          |
| h=24          | 1.00          | <b>0.84</b> | 0.84    | 1.17   | 0.92        | 0.86           | 0.86           | 0.88        | 0.95           | 1.11          | 0.89          |
| <b>UR</b>     |               |             |         |        |             |                |                |             |                |               |               |
| h=1           | 1.01          | 1.03        | 1.09    | 0.95   | 0.97        | <b>0.87***</b> | 0.95           | 0.91***     | 0.90**         | 0.98          | 0.94**        |
| h=3           | 1.00          | 1.10        | 1.05    | 0.86   | 1.05        | <b>0.81***</b> | 0.92           | 0.89**      | 0.82*          | 1.03          | 0.89***       |
| h=9           | 0.99          | 1.11        | 1.10    | 0.92   | 1.02        | 0.96           | <b>0.91</b>    | 0.97        | 0.98           | 1.16*         | 0.97          |
| h=12          | 0.99          | 1.07        | 1.07    | 0.96   | 0.97        | 0.96           | <b>0.91</b>    | 0.99        | 0.94           | 1.17          | 0.96          |
| h=24          | 1.02**        | 1.02        | 1.03    | 1.06   | 0.91*       | 0.84           | <b>0.81</b>    | 0.91        | 0.97           | 1.28          | 0.87          |
| <b>SPREAD</b> |               |             |         |        |             |                |                |             |                |               |               |
| h=1           | 0.99          | 2.88***     | 1.23*** | 1.21** | 3.52***     | 1.07           | <b>0.91***</b> | 0.99        | 0.98           | 0.96          | 0.93**        |
| h=3           | 1.01          | 1.68***     | 1.07    | 1.25   | 1.69***     | 0.82**         | <b>0.81***</b> | 1.06        | 0.85**         | 1.00          | 0.88**        |
| h=9           | 1.01          | 1.36        | 1.27    | 1.06   | 0.94        | 0.73**         | 0.72**         | 0.70***     | <b>0.62***</b> | 1.07          | 0.67***       |
| h=12          | 1.02          | 1.28        | 1.28    | 1.05   | 0.80***     | 0.66***        | <b>0.60***</b> | 0.68***     | 0.65***        | 1.07          | 0.64***       |
| h=24          | 1.03          | 1.34*       | 1.34*   | 0.96   | 0.80*       | 0.70*          | 0.71*          | 0.69**      | <b>0.63***</b> | 0.90          | 0.70**        |
| <b>INF</b>    |               |             |         |        |             |                |                |             |                |               |               |
| h=1           | 1.02          | 1.11*       | 1.18*   | 0.99   | 1.07        | 1.06*          | 1.01           | 0.95        | 0.96           | 0.95          | <b>0.93**</b> |
| h=3           | 1.04          | 1.02        | 1.24*   | 1.04   | 0.93        | <b>0.88</b>    | 1.05           | 0.90        | 0.88           | 0.90          | 0.88          |
| h=9           | 1.07          | 0.92        | 1.01    | 1.16   | 0.86        | 0.78           | 1.15*          | <b>0.72</b> | 0.82           | 0.73          | 0.76          |
| h=12          | 1.09*         | 0.91        | 0.91    | 1.21   | 0.88        | 0.79           | 1.15*          | 0.73        | 0.67           | <b>0.67*</b>  | 0.70          |
| h=24          | 1.04          | 0.90**      | 0.86**  | 1.35   | 1.00        | 1.12           | 1.12           | 0.71        | 0.69           | <b>0.55**</b> | 0.73          |
| <b>HOUST</b>  |               |             |         |        |             |                |                |             |                |               |               |
| h=1           | <b>1.00</b>   | 1.10**      | 1.35*** | 1.07   | 1.08**      | 1.02           | 1.00           | 1.01        | 1.02           | 1.02          | 1.01          |
| h=3           | <b>0.96**</b> | 1.06        | 1.34*** | 1.15   | 1.03        | 1.07           | 1.03           | 1.04        | 1.03           | 1.01          | 1.04          |
| h=9           | <b>0.98</b>   | 1.05        | 1.12    | 1.35   | 0.98        | 1.02           | 1.01           | 1.02        | 1.14           | 1.03          | 1.03          |
| h=12          | 0.98          | 1.05        | 1.05    | 1.32   | <b>0.95</b> | 1.00           | 1.01           | 1.00        | 1.12           | 1.11          | 1.03          |
| h=24          | 0.95          | 1.09        | 1.07    | 1.17   | <b>0.87</b> | 0.94           | 0.95           | 1.00        | 1.15           | 1.23          | 1.06          |

Notes: This table report the root MSPE of the model  $m$  with respect to the root MSPE the AR(4). Best forecast of the row is in bold. Diebold-Mariano test is for each model against the AR(4). "\*", "\*\*" and "\*\*\*" means p-values of below 10%, 5% and 1%. "AO- $i$ " means  $i$ -months moving average forecasts à la [Atkeson et al. \(2001\)](#).