# How is Machine Learning Useful for Macroeconomic Forecasting?

PHILIPPE **GOULET COULOMBE**

MAXIME **LEROUX**

DALIBOR **STEVANOVIC**

STÉPHANE **SURPRENANT**

**2019s-22**

# How is Machine Learning Useful for Macroeconomic Forecasting?

*Philippe Goulet Coulombe, Maxime Leroux, Dalibor Stevanovic, Stéphane Surprenant*

**Série Scientifique**
*Scientific Series*

**Montréal**
**Octobre/October 2019**

# How is Machine Learning Useful for Macroeconomic Forecasting? [*]

*Philippe Goulet Coulombe [†], Maxime Leroux [‡],*
*Dalibor Stevanovic [§], Stéphane Surprenant [**]*

## Abstract/Résumé

We move beyond *Is Machine Learning Useful for Macroeconomic Forecasting?* by adding the *how*. The current forecasting literature has focused on matching specific variables and horizons with a particularly successful algorithm. To the contrary, we study the usefulness of the underlying features driving ML gains over standard macroeconometric methods. We distinguish four so-called features (nonlinearities, regularization, cross-validation and alternative loss function) and study their behavior in both the data-rich and data-poor environments. To do so, we design experiments that allow to identify the "treatment" effects of interest. We conclude that (i) nonlinearity is the true game changer for macroeconomic prediction, (ii) the standard factor model remains the best regularization, (iii) K-fold cross-validation is the best practice and (iv) the $L_2$ is preferred to the e-insensitive in-sample loss. The forecasting gains of nonlinear techniques are associated with high macroeconomic uncertainty, financial stress and housing bubble bursts. This suggests that Machine Learning is useful for macroeconomic forecasting by mostly capturing important nonlinearities that arise in the context of uncertainty and financial frictions.

**Keywords/Mots-clés:** Machine Learning, Big Data, Forecasting

**JEL Codes/Codes JEL:** C53, C55, E37

---

[†] Corresponding Author: gouletc@sas.upenn.edu. Department of Economics, UPenn.
[‡] Université du Québec à Montréal.
[§] Corresponding Author: dstevanovic.econ@gmail.com. Département des sciences économiques, UQAM.
[**] Université du Québec à Montréal.

# 1 Introduction

The intersection of Machine Learning (ML) with econometrics has become an important research landscape in economics. ML has gained prominence due to the availability of large data sets, especially in microeconomic applications (Belloni et al., 2017; Athey, 2019). Despite the growing interest in ML, understanding the properties of ML procedures when they are applied to predict macroeconomic outcomes remains a difficult challenge.[1] Nevertheless, that very understanding is an interesting econometric research endeavor *per se*. It is more appealing to applied econometricians to upgrade a standard framework with a subset of specific insights rather than to drop everything altogether for an off-the-shelf ML model.

Despite appearances, ML has a long history in macroeconometrics (see Lee et al. (1993); Kuan and White (1994); Swanson and White (1997); Stock and Watson (1999); Trapletti et al. (2000); Medeiros et al. (2006)). However, only recently did the field of macroeconomic forecasting experience an overwhelming (and succesful) surge in the number of studies applying ML methods,[2] while works such as Joseph (2019) and Zhao and Hastie (2019) contribute to their interpretability. However, the vast catalogue of tools, often evaluated with few models and forecasting targets, creates a large conceptual space, much of which remains to be explored. To map that large space without getting lost in it, we move beyond the coronation of a single winning model and its subsequent interpretation. Rather, we conduct a meta-analysis of many ML products by projecting them in their "characteristic" space. Then, we provide a direct assessment of which characteristics matter and which do not.

More precisely, we aim to answer the following question: What are the key features of ML modeling that improve the macroeconomic prediction? In particular, no clear attempt

---

[1]The linear techniques have been extensively examined since Stock and Watson (2002b,a). Kotchoni et al. (2019) compare more than 30 forecasting models, including factor-augmented and regularized regressions. Giannone et al. (2018) study the relevance of sparse modeling in various economic prediction problems.

[2]Moshiri and Cameron (2000); Nakamura (2005); Marcellino (2008) use neural networks to predict inflation and Cook and Smalter Hall (2017) explore deep learning. Sermpinis et al. (2014) apply support vector regressions, while Diebold and Shin (2019) propose a LASSO-based forecast combination technique. Ng (2014), Döpke et al. (2017) and Medeiros et al. (2019) improve forecast accuracy with random forests and boosting, while Yousuf and Ng (2019) use boosting for high-dimensional predictive regressions with time varying parameters. Others compare machine learning methods in horse races (Ahmed et al., 2010; Stock and Watson, 2012b; Li and Chen, 2014; Kim and Swanson, 2018; Smeekes and Wijler, 2018; Chen et al., 2019; Milunovich, 2020).

has been made at understanding why one algorithm might work while another does not. We address this question by designing an *experiment* to identify important characteristics of machine learning and big data techniques. The exercise consists of an extensive pseudo-out-of-sample forecasting horse race between many models that differ with respect to the four main features: nonlinearity, regularization, hyperparameter selection and loss function. To control for the big data aspect, we consider data-poor and data-rich models, and administer those *patients* one particular ML *treatment* or combinations of them. Monthly forecast errors are constructed for five important macroeconomic variables, five forecasting horizons and for almost 40 years. Then, we provide a straightforward framework to identify which of them are actual game changers for macroeconomic forecasting.

The main results can be summarized as follows. First, the ML nonparametric nonlinearities constitute the most salient feature as they improve substantially the forecasting accuracy for all macroeconomic variables in our exercise, especially when predicting at long horizons. Second, in the big data framework, alternative regularization methods (Lasso, Ridge, Elastic-net) do not improve over the factor model, suggesting that the factor representation of the macroeconomy is quite accurate as a means of dimensionality reduction.

Third, the hyperparameter selection by K-fold cross-validation (CV) and the standard BIC (when possible) do better on average than any other criterion. This suggests that ignoring information criteria when opting for more complicated ML models is not harmful. This is also quite convenient: K-fold is the built-in CV option in most standard ML packages. Fourth, replacing the standard in-sample quadratic loss function by the $\bar{\epsilon}$-insensitive loss function in Support Vector Regressions (SVR) is not useful, except in very rare cases. The latter finding is a direct by-product of our strategy to disentangle treatment effects. In accordance with other empirical results (Sermpinis et al., 2014; Colombo and Pelagatti, 2020), in absolute terms, SVRs do perform well – even if they use a loss at odds with the one used for evaluation. However, that performance is a mixture of the attributes of both nonlinearities (via the kernel trick) and an alternative loss function. Our results reveal that this change in the loss function has detrimental effects on performance in terms of both mean squared errors and absolute

errors. Fifth, the marginal effect of big data is positive and significant, and improves as the forecast horizon grows. The robustness analysis shows that these results remain valid when: (i) the absolute loss is considered; (ii) quarterly targets are predicted; (iii) the exercise is re-conducted with a large Canadian data set.

The evolution of economic uncertainty and financial conditions are important drivers of the NL treatment effect. ML nonlinearities are particularly useful: (i) when the level of macroeconomic uncertainty is high; (ii) when financial conditions are tight and (iii) during housing bubble bursts. The effects are bigger in the case of data-rich models, which suggests that combining nonlinearity with factors made of many predictors is an accurate way to capture complex macroeconomic relationships.

These results give a clear recommendation for practitioners. For most cases, start by reducing the dimensionality with principal components and then augment the standard diffusion indices model by a ML nonlinear function approximator of your choice. That recommendation is conditional on being able to keep overfitting in check. To that end, if cross-validation must be applied to hyperparameter selection, the best practice is the standard K-fold.

These novel empirical results also complement a growing theoretical literature on ML with dependent observations. As Alquier et al. (2013) points out, much of the work in statistical learning has focus on the cross-section setting where the assumption of independent draws is more plausible. Nevertheless, some theoretical guarantees exist in the time series context. Mohri and Rostamizadeh (2010) provide generalization bounds for Support Vector Machines and Regressions, and Kernel Ridge Regression under the assumption of a stationary joint distribution of predictors and target variable. Kuznetsov and Mohri (2015) generalize some of those results to non-stationary distributions and non-mixing processes. However, as the macroeconomic time series framework is characterized by short samples and structural instability, our exercise contributes to the general understanding of machine learning properties in the context of time series modeling and forecasting.

In the remainder of this paper, we first present the general prediction problem with machine learning and big data. Section 3 describes the four important features of machine learn-

ing methods. Section 4 presents the empirical setup, section 5 discusses the main results, followed by section 6 that aims to open the black box. Section 7 concludes. Appendices A, B, C and D contain respectively: tables with overall performance; robustness of treatment analysis; additional results and robustness of nonlinearity analysis. The supplementary material contains the following appendices: results for absolute loss, results with quarterly US data, results with monthly Canadian data, description of CV techniques and technical details on forecasting models.

## 2   Making Predictions with Machine Learning and Big Data

Machine learning methods are meant to improve our predictive ability especially when the "true" model is unknown and complex. To illustrate this point, let $y_{t+h}$ be the variable to be predicted $h$ periods ahead (target) and $Z_t$ the $N_Z$-dimensional vector of predictors made out of $H_t$, the set of all the inputs available at time $t$. Let $g^*(Z_t)$ be the true model and $g(Z_t)$ a functional (parametric or not) form selected by the practitioner. In addition, denote $\hat{g}(Z_t)$ and $\hat{y}_{t+h}$ the fitted model and its forecast. The forecast error can be decomposed as

$$y_{t+h} - \hat{y}_{t+h} = \underbrace{g^*(Z_t) - g(Z_t)}_{\text{approximation error}} + \underbrace{g(Z_t) - \hat{g}(Z_t)}_{\text{estimation error}} + e_{t+h}. \tag{1}$$

The intrinsic error $e_{t+h}$ is not shrinkable, while the estimation error can be reduced by adding more data. The approximation error is controlled by the functional estimator choice. While it can be potentially minimized by using flexible functions, it also rise the risk of overfitting and a judicious regularization is needed to control this risk. This problem can be embedded in the general prediction setup from Hastie et al. (2009)

$$\min_{g \in \mathcal{G}} \{ \hat{L}(y_{t+h}, g(Z_t)) + pen(g; \tau) \}, \quad t = 1, \dots, T. \tag{2}$$

This setup has four main features:

1. $\mathcal{G}$ is the space of possible functions $g$ that combine the data to form the prediction. In particular, the interest is how much nonlinearities can we allow for in order to reduce the approximation error in (1)?

2. $pen()$ is the regularization penalty limiting the flexibility of the function $g$ and hence

5

controlling the overfitting risk. This is quite general and can accommodate Bridge-type penalties and dimension reduction techniques.

3. $\tau$ is the set of hyperparameters including those in the penalty and the approximator $g$. The usual problem is to choose the best data-driven method to optimize $\tau$.

4. $\hat{L}$ is the loss function that defines the optimal forecast. Some ML models feature an in-sample loss function different from the standard $l_2$ norm.

Most of (supervised) machine learning consists of a combination of those ingredients and popular methods like linear (penalized) regressions can be obtained as special cases of (2).

## 2.1 Predictive Modeling

We consider the *direct* predictive modeling in which the target is projected on the information set, and the forecast is made directly using the most recent observables. This is opposed to *iterative* approach where the model recursion is used to simulate the future path of the variable.[3] Also, the direct approach is the standard practice for in ML applications.

We now define the forecast objective given the variable of interest $Y_t$. If $Y_t$ is stationary, we forecast its level $h$ periods ahead:

$$y_{t+h}^{(h)} = y_{t+h}, \tag{3}$$

where $y_t \equiv \ln Y_t$ if $Y_t$ is strictly positive. If $Y_t$ is I(1), then we forecast the average growth rate over the period $[t+1, t+h]$ (Stock and Watson, 2002b). We shall therefore define $y_{t+h}^{(h)}$ as:

$$y_{t+h}^{(h)} = (1/h)\ln(Y_{t+h}/Y_t). \tag{4}$$

In order to avoid a cumbersome notation, we use $y_{t+h}$ instead of $y_{t+h}^{(h)}$ in what follows. In addition, all the predictors in $Z_t$ are assumed to be covariance stationary.

## 2.2 *Data-Poor* versus *Data-Rich* Environments

Large time series panels are now widely constructed and used for macroeconomic analysis. The most popular is FRED-MD monthly panel of US variables constructed by McCracken

---

[3]Marcellino et al. (2006) conclude that the direct approach provides slightly better results but does not dominate uniformly across time and series. See Chevillon (2007) for a survey on multi-step forecasting.

and Ng (2016).[4] Unfortunately, the performance of standard econometric models tends to deteriorate as the dimensionality of data increases. Stock and Watson (2002b) first proposed to solve the problem by replacing the high-dimensional predictor set by common factors.[5]

On other hand, even though the machine learning models do not require big data, they are useful to perform variable selection and digest large information sets to improve the prediction. Therefore, in addition to treatment effects in terms of characteristics of forecasting models, we will also interact those with the width of the sample. The data-poor, defined as $H_t^-$, will only contain a finite number of lagged values of the target, while the data-rich panel, defined as $H_t^+$ will also include a large number of exogenous predictors. Formally,

$$H_t^- \equiv \{y_{t-j}\}_{j=0}^{p_y} \quad \text{and} \quad H_t^+ \equiv \left[ \{y_{t-j}\}_{j=0}^{p_y}, \{X_{t-j}\}_{j=0}^{p_f} \right]. \tag{5}$$

The analysis we propose can thus be summarized in the following way. We will consider two standard models for forecasting.

1. The $H_t^-$ model is the *autoregressive direct* (AR) model, which is specified as:

$$y_{t+h} = c + \rho(L)y_t + e_{t+h}, \quad t = 1, \ldots, T, \tag{6}$$

where $h \geq 1$ is the forecasting horizon. The only hyperparameter in this model is $p_y$, the order of the lag polynomial $\rho(L)$.

2. The $H_t^+$ workhorse model is the autoregression augmented with diffusion indices (ARDI) from Stock and Watson (2012b):

$$y_{t+h} = c + \rho(L)y_t + \beta(L)F_t + e_{t+h}, \quad t = 1, \ldots, T \tag{7}$$

$$X_t = \Lambda F_t + u_t \tag{8}$$

where $F_t$ are $K$ *consecutive static* factors, and $\rho(L)$ and $\beta(L)$ are lag polynomials of orders $p_y$ and $p_f$ respectively. The feasible procedure requires an estimate of $F_t$ that is usually obtained by principal component analysis (PCA).

Then, we will take these models as two different types of "patients" and will administer them

---

[4]Fortin-Gagnon et al. (2020) have recently proposed similar data for Canada.

[5]Another way to approach the dimensionality problem is to use Bayesian methods. Indeed, some of our Ridge regressions will look like a direct version of a Bayesian VAR with a Litterman (1979) prior. Giannone et al. (2015) have shown that an hierarchical prior can lead the BVAR to perform as well as a factor model.

one particular ML treatment or combinations of them. That is, we will upgrade these models with one or many features of ML and evaluate the gains/losses in both environments. From the perspective of the machine learning literature, equation (8) motivates the use of PCA as a form of feature engineering. Although more sophisticated methods have been used[6], PCA remains popular (Uddin et al., 2018). As we insist on treating models as symmetrically as possible, we will use the same feature transformations throughout such that our nonlinear models, such as Kernel Ridge Regression, will introduce nonlinear transformations of lagged target values as well as of lagged values of the principal components. Hence, our nonlinear models postulate that a sparse set of latent variables impact the target in a flexible way.[7]

## 2.3 Evaluation

The objective of this paper is to disentangle important characteristics of the ML prediction algorithms when forecasting macroeconomic variables. To do so, we design an *experiment* that consists of a pseudo-out-of-sample (POOS) forecasting horse race between many models that differ with respect to the four main features above, i.e., nonlinearity, regularization, hyperparameter selection and loss function. To create variation around those *treatments*, we will generate forecast errors from different models associated to each feature.

To test this paper's hypothesis, suppose the following model for forecasting errors

$$e^2_{t,h,v,m} = \alpha_m + \psi_{t,v,h} + v_{t,h,v,m} \tag{9a}$$

$$\alpha_m = \alpha'_F \mathbf{1} + \eta_m \tag{9b}$$

where $e^2_{t,h,v,m}$ are squared prediction errors of model $m$ for variable $v$ and horizon $h$ at time $t$. $\psi_{t,v,h}$ is a fixed effect term that demeans the dependent variable by "forecasting target", that is a combination of $t$, $v$ and $h$. $\alpha_F$ is a vector of $\alpha_{\mathcal{G}}$, $\alpha_{pen()}$, $\alpha_\tau$ and $\alpha_{\hat{L}}$ terms associated to each feature. We re-arrange equation (9) to obtain

$$e^2_{t,h,v,m} = \alpha'_F \mathbf{1} + \psi_{t,v,h} + u_{t,h,v,m}. \tag{10}$$

---

[6]The autoencoder method of Gu et al. (2020a) can be seen as a form of feature engineering, just as the independent components used in conjunction with SVR in Lu et al. (2009). The interested reader may also see Hastie et al. (2009) for a detailed discussion of the use of PCA and related method in machine learning.

[7]We omit considering a VAR as an additional option. VAR iterative approach to produce $h$-step-ahead predictions is not comparable with the direct forecasting used with ML models.

$H_0$ is now $\alpha_f = 0 \quad \forall f \in F = [\mathcal{G}, \text{ pen}(), \tau, \hat{L}]$. In other words, the null is that there is no predictive accuracy gain with respect to a base model that does not have this particular feature.[8] By interacting $\alpha_F$ with other fixed effects or variables, we can test many hypotheses about the heterogeneity of the "ML treatment effect." To get interpretable coefficients, we define $R^2_{t,h,v,m} \equiv 1 - \frac{e^2_{t,h,v,m}}{\frac{1}{T}\sum_{t=1}^{T}(y_{v,t+h}-\bar{y}_{v,h})^2}$ and run

$$R^2_{t,h,v,m} = \dot{\alpha}'_F \mathbf{1} + \dot{\psi}_{t,v,h} + \dot{u}_{t,h,v,m}. \tag{11}$$

While (10) has the benefit of connecting directly with the specification of a Diebold and Mariano (1995) test, the transformation of the regressand in (11) has two main advantages justifying its use. First and foremost, it provides standardized coefficients $\dot{\alpha}_F$ interpretable as marginal improvements in OOS-$R^2$'s. In contrast, $\alpha_F$ are a unit- and series-dependant marginal increases in MSE. Second, the $R^2$ approach has the advantage of standardizing *ex-ante* the regressand and removing an obvious source of $(v, h)$-driven heteroskedasticity.

While the generality of (10) and (11) is appealing, when investigating the heterogeneity of specific partial effects, it will be much more convenient to run specific regressions for the multiple hypothesis we wish to test. That is, to evaluate a feature $f$, we run

$$\forall m \in \mathcal{M}_f: \quad R^2_{t,h,v,m} = \dot{\alpha}_f + \dot{\phi}_{t,v,h} + \dot{u}_{t,h,v,m} \tag{12}$$

where $\mathcal{M}_f$ is defined as the set of models that differs only by the feature under study $f$. An analogous evaluation setup has been considered in Carriero et al. (2019).

## 3 Four Features of ML

In this section we detail the forecasting approaches that create variations for each characteristic of machine learning prediction problem defined in (2).

### 3.1 Feature 1: Nonlinearity

Although linearity is popular in practice, if the data generating process (DGP) is complex, using linear $g$ introduces approximation error as shown in (1). As a solution, ML proposes an apparatus of nonlinear functions able to estimate the true DGP, and thus reduces the approx-

---

[8]If we consider two models that differ in one feature and run this regression for a specific $(h, v)$ pair, the t-test on coefficients amounts to Diebold and Mariano (1995) – conditional on having the proper standard errors.

imation error. We focus on applying the Kernel trick and random forests to our two baseline models to see if the nonlinearities they generate will lead to significant improvements.[9]

### 3.1.1 Kernel Ridge Regression

A simple way to make predictive regressions (6) and (7) nonlinear is to adopt a generalized linear model with multivariate functions of predictors (e.g. spline series expansions). However, this rapidly becomes overparameterized, so we opt for the Kernel trick (KT) to avoid computing all possible interactions and higher order terms. It is worth noting that Kernel Ridge Regression (KRR) has several implementation advantages. It has a closed-form solution that rules out convergence problems associated with models trained with gradient descent. It is also fast to implement since it implies inverting a $T$x$T$ matrix at each step.

To show how KT is implemented in our benchmark models, suppose a Ridge regression direct forecast with generic regressors $Z_t$

$$\min_{\beta} \sum_{t=1}^{T} (y_{t+h} - Z_t \beta)^2 + \lambda \sum_{k=1}^{K} \beta_k^2.$$

The solution to that problem is $\hat{\beta} = (Z'Z + \lambda I_k)^{-1} Z'y$. By the representer theorem of Smola and Schölkopf (2004), $\beta$ can also be obtained by solving the dual of the convex optimization problem above. The dual solution for $\beta$ is $\hat{\beta} = Z'(ZZ' + \lambda I_T)^{-1} y$. This equivalence allows to rewrite the conditional expectation in the following way:

$$\hat{E}(y_{t+h}|Z_t) = Z_t \hat{\beta} = \sum_{i=1}^{t} \hat{\alpha}_i \langle Z_i, Z_t \rangle$$

where $\hat{\alpha} = (ZZ' + \lambda I_T)^{-1} y$ is the solution to the dual Ridge Regression problem.

Suppose now we approximate a general nonlinear model $g(Z_t)$ with basis functions $\phi()$

$$y_{t+h} = g(Z_t) + \varepsilon_{t+h} = \phi(Z_t)' \gamma + \varepsilon_{t+h}.$$

---

[9]A popular approach to model nonlinearity is deep learning. However, since we re-optimize our models recursively in a POOS, selecting an accurate network architecture by cross-validation is practically infeasible. In addition to optimize numerous neural net hyperparameters (such as the number of hidden layers and neurons, activation function, etc.), our forecasting models also require careful input selection (number of lags and number of factors in case of data-rich). An alternative is to fix ex-ante a variety of networks as in Gu et al. (2020b), but this would potentially benefit other models that are optimized over time. Still, since few papers have found similar predictive ability of random forests and neural nets (Gu et al., 2020a; Joseph, 2019), we believe that considering random forests and Kernel trick is enough to properly identify the ML nonlinear treatment. Nevertheless, we have conducted a robustness analysis with feed-forward neural networks and boosted trees. The results are presented in Appendix D.

The so-called Kernel trick is the fact that there exist a reproducing kernel $K()$ such that

$$\hat{E}(y_{t+h}|Z_t) = \sum_{i=1}^{t} \hat{\alpha}_i \langle \phi(Z_i), \phi(Z_t) \rangle = \sum_{i=1}^{t} \hat{\alpha}_i K(Z_i, Z_t).$$

This means we do not need to specify the numerous basis functions, a well-chosen kernel implicitly replicates them. This paper will use the standard radial basis function (RBF) kernel

$$K_\sigma(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

where $\sigma$ is a tuning parameter to be chosen by cross-validation. This choice of kernel is motivated by its good performance in macroeconomic forecasting as reported in Sermpinis et al. (2014) and Exterkate et al. (2016). The advantage of the kernel trick is that, by using the corresponding $Z_t$, we can easily make our data-rich or data-poor model nonlinear. For instance, in the case of the factor model, we can apply it to the regression equation to implicitly estimate

$$y_{t+h} = c + g(Z_t) + \varepsilon_{t+h}, \tag{13}$$

$$Z_t = \left[ \{y_{t-j}\}_{j=0}^{p_y}, \{F_{t-j}\}_{j=0}^{p_f} \right], \tag{14}$$

$$X_t = \Lambda F_t + u_t. \tag{15}$$

In terms of implementation, this means extracting factors via PCA and then getting

$$\hat{E}(y_{t+h}|Z_t) = K_\sigma(Z_t, Z)(K_\sigma(Z_t, Z) + \lambda I_T)^{-1} y_t. \tag{16}$$

The final set of tuning parameters for such a model is $\tau = \{\lambda, \sigma, p_y, p_f, n_f\}$.

### 3.1.2 Random Forests

Another way to introduce nonlinearity in the estimation of the predictive equation (7) is to use regression trees instead of OLS. The idea is to split sequentially the space of $Z_t$, as defined in (14) into several regions and model the response by the mean of $y_{t+h}$ in each region. The process continues according to some stopping rule. The details of the recursive algorithm can be found in Hastie et al. (2009). Then, the tree regression forecast has the following form:

$$\hat{f}(Z) = \sum_{m=1}^{M} c_m I_{(Z \in R_m)}, \tag{17}$$

where $M$ is the number of terminal nodes, $c_m$ are node means and $R_1, ..., R_M$ represents a partition of feature space. In the diffusion indices setup, the regression tree would estimate a

nonlinear relationship linking factors and their lags to $y_{t+h}$. Once the tree structure is known, it can be related to a linear regression with dummy variables and their interactions.

While the idea of obtaining nonlinearities via decision trees is intuitive and appealing – especially for its interpretability potential, the resulting prediction is usually plagued by high variance. The recursive tree fitting process is (i) unstable and (ii) prone to overfitting. The latter can be partially addressed by the use of pruning and related methodologies (Hastie et al., 2009). Notwithstanding, a much more successful (and hence popular) fix was proposed in Breiman (2001): Random Forests. This consists in growing many trees on subsamples (or nonparametric bootstrap samples) of observations. Further randomization of underlying trees is obtained by considering a random subset of regressors for each potential split.[10] The main hyperparameter to be selected is the number of variables to be considered at each split. The forecasts of the estimated regression trees are then averaged together to make one single "ensemble" prediction of the targeted variable.[11]

## 3.2 Feature 2: Regularization

In this section we will only consider models where dimension reduction is needed, which are the models with $H_t^+$. The traditional shrinkage method used in macroeconomic forecasting is the ARDI model that consists of extracting principal components of $X_t$ and to use them as data in an ARDL model. Obviously, this is only one out of many ways to compress the information contained in $X_t$ to run a well-behaved regression of $y_{t+h}$ on it.[12]

In order to create identifying variations for $pen()$ treatment, we need to generate multiple different shrinkage schemes. Some will also blend in selection, some will not. The alternative shrinkage methods will all be special cases of the Elastic Net (EN) problem:

$$\min_{\beta} \sum_{t=1}^{T} (y_{t+h} - Z_t\beta)^2 + \lambda \sum_{k=1}^{K} \left( \alpha|\beta_k| + (1-\alpha)\beta_k^2 \right) \tag{18}$$

where $Z_t = B(H_t)$ is some transformation of the original predictive set $X_t$. $\alpha \in [0,1]$ and

---

[10]Only using a bootstrap sample of observations would be a procedure called Bagging – for Bootstrap Aggregation. Also selecting randomly regressors has the effect of decorrelating the trees and hence boosting the variance reduction effect of averaging them.

[11]In this paper, we consider 500 trees, which is usually more than enough to get a stabilized prediction (that will not change with the addition of another tree).

[12]De Mol et al. (2008) compares Lasso, Ridge and ARDI and finds that forecasts are very much alike.

$\lambda > 0$ can either be fixed or found via CV. By using different $B$ operators, we can generate shrinkage schemes. Also, by setting $\alpha$ to either 1 or 0 we generate LASSO and Ridge Regression respectively. All these possibilities are reasonable alternatives to the traditional factor hard-thresholding procedure that is ARDI.

Each type of shrinkage in this section will be defined by the tuple $S = \{\alpha, B()\}$. To begin with the most straightforward dimension, for a given $B$, we will evaluate the results for $\alpha \in \{0, \hat{\alpha}_{CV}, 1\}$. For instance, if $B$ is the identity mapping, we get in turns the LASSO, EN and Ridge shrinkage. We now detail different $pen()$ resulting when we vary $B()$ for a fixed $\alpha$.

1. **(Fat Regression)**: First, we consider the case $B_1() = I()$. That is, we use the entirety of the untransformed high-dimensional data set. The results of Giannone et al. (2018) point in the direction that specifications with a higher $\alpha$ should do better, that is, sparse models do worse than models where every regressor is kept but shrunk to zero.

2. **(Big ARDI)** Second, $B_2()$ corresponds to first rotating $X_t \in \mathbb{R}^N$ so that we get $N$-dimensional uncorrelated $F_t$. Note here that contrary to the ARDI approach, we do not select factors recursively, we keep them all. Hence, $F_t$ has exactly the same span as $X_t$. Comparing LASSO and Ridge in this setup will allow to verify whether sparsity emerges in a rotated space.

3. **(Principal Component Regression)** A third possibility is to rotate $H_t^+$ rather than $X_t$ and still keep all the factors. $H_t^+$ includes all the relevant preselected lags. If we were to just drop the $F_t$ using some hard-thresholding rule, this would correspond to Principal Component Regression (PCR). Note that $B_3() = B_2()$ only when no lags are included.

Hence, the tuple $S$ has a total of 9 elements. Since we will be considering both POOS-CV and K-fold CV for each of these models, this leads to a total of 18 models.[13]

To see clearly through all of this, we describe where the benchmark ARDI model stands in this setup. Since it uses a hard thresholding rule that is based on the eigenvalues ordering, it cannot be a special case of the Elastic Net problem. While it uses $B_2$, we would need to set

---

[13] Adaptive versions (in the sense of Zou (2006)) of the 9 models were also considered but gave either similar or deteriorated results with respect to their plain counterparts.

$\lambda = 0$ and select $F_t$ *a priori* with a hard-thresholding rule. The closest approximation in this EN setup would be to set $\alpha = 1$ and fix the value of $\lambda$ to match the number of consecutive factors selected by an information criteria directly in the predictive regression (7).

## 3.3   Feature 3: Hyperparameter Optimization

The conventional wisdom in macroeconomic forecasting is to either use AIC or BIC and compare results. The prime reason for the popularity of CV is that it can be applied to any model, including those for which the derivation of an information criterion is impossible.[14]

It is not obvious that CV should work better only because it is "out of sample" while AIC and BIC are "in sample". All model selection methods are actually approximations to the OOS prediction error that relies on different assumptions that are sometime motivated by different theoretical goals. Also, it is well known that asymptotically, these methods have similar behavior.[15] Hence, it is impossible *a priori* to think of one model selection technique being the most appropriate for macroeconomic forecasting.

For samples of small to medium size encountered in macro, the question of which one is optimal in the forecasting sense is inevitably an empirical one. For instance, Granger and Jeon (2004) compared AIC and BIC in a generic forecasting exercise. In this paper, we will compare AIC, BIC and two types of CV for our two baseline models. The two types of CV are relatively standard. We will first use POOS CV and then K-fold CV. The first one will always behave correctly in the context of time series data, but may be quite inefficient by only using the end of the training set. The latter is known to be valid only if residual autocorrelation is absent from the models as shown in Bergmeir et al. (2018). If it were not to be the case, then we should expect K-fold to underperform. The specific details of the implementation of both CVs is discussed in the section D of the supplementary material.

The contributions of this section are twofold. First, it will shed light on which model

---

[14] Abadie and Kasy (2019) show that hyperparemeter tuning by CV performs uniformly well in high-dimensional context.

[15] Hansen and Timmermann (2015) show equivalence between test statistics for OOS forecasting performance and in-sample Wald statistics. For instance, one can show that Leave-one-out CV (a special case of K-fold) is asymptotically equivalent to the Takeuchi Information criterion (TIC), Claeskens and Hjort (2008). AIC is a special case of TIC where we need to assume in addition that all models being considered are at least correctly specified. Thus, under the latter assumption, Leave-one-out CV is asymptotically equivalent to AIC.

selection method is most appropriate for typical macroeconomic data and models. Second, we will explore how much of the gains/losses of using ML can be attributed to widespread use of CV. Since most nonlinear ML models cannot be easily tuned by anything other than CV, it is hard for the researcher to disentangle between gains coming from the ML method itself or just the way it is tuned.[16] Hence, it is worth asking the question whether some gains from ML are simply coming from selecting hyperparameters in a different fashion using a method whose assumptions are more in line with the data at hand. To investigate that, a natural first step is to look at our benchmark macro models, AR and ARDI, and see if using CV to select hyperparameters gives different selected models and forecasting performances.

## 3.4 Feature 4: Loss Function

Until now, all of our estimators use a quadratic loss function. Of course, it is very natural for them to do so: the quadratic loss is the measure used for out-of-sample evaluation. Thus, someone may legitimately wonder if the fate of the SVR is not sealed in advance as it uses an in-sample loss function which is inconsistent with the out-of-sample performance metric. As we will discuss later after the explanation of the SVR, there are reasons to believe the alternative (and mismatched) loss function can help. As a matter of fact, SVR has been successfully applied to forecasting financial and macroeconomic time series.[17] An important question remains unanswered: are the good results due to kernel-based non-linearities or to the use of an alternative loss-function?

We provide a strategy to isolate the marginal effect of the SVR's $\bar{\epsilon}$-insensitive loss function which consists in, perhaps unsurprisingly by now, estimating different variants of the same model. We considered the Kernel Ridge Regression earlier. The latter only differs from the Kernel-SVR by the use of different in-sample loss functions. This identifies directly the effect of the loss function, for nonlinear models. Furthermore, we do the same exercise for linear

---

[16]Zou et al. (2007) show that the number of remaining parameters in the LASSO is an unbiased estimator of the degrees of freedom and derive LASSO-BIC and LASSO-AIC criteria. Considering these as well would provide additional evidence on the empirical debate of CV vs IC.

[17]See for example, Lu et al. (2009), Choudhury et al. (2014), Patel et al. (2015a), Patel et al. (2015b), Yeh et al. (2011) and Qu and Zhang (2016) for financial forecasting. See Sermpinis et al. (2014) and Zhang et al. (2010) macroeconomic forecasting.

models: comparing a linear SVR to the plain ARDI. To sum up, to isolate the "treatment effect" of a different in-sample loss function, we consider: (1) the linear SVR with $H_t^-$; (2) the linear SVR with $H_t^+$; (3) the RBF Kernel SVR with $H_t^-$; and (4) the RBF Kernel SVR with $H_t^+$.

What follows is a bird's-eye overview of the underlying mechanics of the SVR. As it was the case for the Kernel Ridge regression, the SVR estimator approximates the function $g \in G$ with basis functions. We opted to use the $\epsilon$-SVR variant which implicitly defines the size $2\bar{\epsilon}$ of the insensitivity tube of the loss function. The $\epsilon$-SVR is defined by:

$$\min_{\gamma} \frac{1}{2}\gamma'\gamma + C\left[\sum_{t=1}^{T}(\xi_t + \xi_t^*)\right]$$

$$s.t. \begin{cases} y_{t+h} - \gamma'\phi(Z_t) - \alpha \leq \bar{\epsilon} + \xi_t \\ \gamma'\phi(Z_t) + \alpha - y_{t+h} \leq \bar{\epsilon} + \xi_t^* \\ \xi_t, \xi_t^* \geq 0. \end{cases}$$

Where $\xi_t, \xi_t^*$ are slack variables, $\phi()$ is the basis function of the feature space implicitly defined by the kernel used and $T$ is the size of the sample used for estimation. $C$ and $\bar{\epsilon}$ are hyperparameters. Additional hyperparameters vary depending on the choice of a kernel. In case of the RBF kernel, a scale parameter $\sigma$ also has to be cross-validated. Associating Lagrange multipliers $\lambda_j, \lambda_j^*$ to the first two types of constraints, Smola and Schölkopf (2004) show that we can derive the dual problem out of which we would find the optimal weights $\gamma = \sum_{j=1}^{T}(\lambda_j - \lambda_j^*)\phi(Z_j)$ and the forecasted values

$$\hat{E}(y_{t+h}|Z_t) = \hat{c} + \sum_{j=1}^{T}(\lambda_j - \lambda_j^*)\phi(Z_j)\phi(Z_j) = \hat{c} + \sum_{j=1}^{T}(\lambda_j - \lambda_j^*)K(Z_j, Z_t). \tag{19}$$

Let us now turn to the resulting loss function of such a problem. For the $\epsilon$-SVR, the penalty is given by:

$$P_{\bar{\epsilon}}(\epsilon_{t+h|t}) := \begin{cases} 0 & if \quad |e_{t+h}| \leq \bar{\epsilon} \\ |e_{t+h}| - \bar{\epsilon} & otherwise \end{cases}.$$

For other estimators, the penalty function is quadratic $P(e_{t+h}) := e_{t+h}^2$. Hence, for our other estimators, the rate of the penalty increases with the size of the forecasting error, whereas it is constant and only applies to excess errors in the case of the $\epsilon$-SVR. Note that this insen-

sitivity has a nontrivial consequence for the forecasting values. The Karush-Kuhn-Tucker conditions imply that only support vectors, i.e. points lying outside the insensitivity tube, will have nonzero Lagrange multipliers and contribute to the weight vector.

As discussed briefly earlier, given that SVR forecasts will eventually be evaluated according to a quadratic loss, it is reasonable to ask why this alternative loss function isn't trivially suboptimal. Smola et al. (1998) show that the optimal size of $\bar{\epsilon}$ is a linear function of the underlying noise, with the exact relationship depending on the nature of the data generating process. This idea is not at odds with Gu et al. (2020a) using the Huber Loss for asset pricing with ML (where outliers seldomly happen in-sample) or Colombo and Pelagatti (2020) successfully using SVR to forecast (notoriously noisy) exchange rates. Thus, while SVR can work well in macroeconomic forecasting, it is unclear which feature between the nonlinearity and $\bar{\epsilon}$-insensitive loss has the primary influence on its performance.

To sum up, the table 1 shows a list of all forecasting models and highlights their relationship with each of four features discussed above. The computational details for every model in this list are available in section E in the supplementary material.

## 4    Empirical setup

This section presents the data and the design of the pseudo-of-sample experiment used to generate the treatment effects above.

### 4.1    Data

We use historical data to evaluate and compare the performance of all the forecasting models described previously. The dataset is FRED-MD, available at the Federal Reserve of St-Louis's web site. It contains 134 monthly US macroeconomic and financial indicators observed from 1960M01 to 2017M12. Since many of them are usually very persistent or not stationary, we follow McCracken and Ng (2016) in the choice of transformations in order to achieve stationarity.[18] Even though the universe of time series available at FRED is huge, we stick to FRED-MD for several reasons. First, we want to have the test set as long as possible

---

[18]Alternative data transformations in the context of ML modeling are used in Goulet Coulombe et al. (2020).

Table 1: List of all forecasting models

| Models | Feature 1: selecting the function $g$ | Feature 2: selecting the regularization | Feature 3: optimizing hyperparameters $\tau$ | Feature 4: selecting the loss function |
| --- | --- | --- | --- | --- |
| *Data-poor models* | | | | |
| AR,BIC | Linear | | BIC | Quadratic |
| AR,AIC | Linear | | AIC | Quadratic |
| AR,POOS-CV | Linear | | POOS CV | Quadratic |
| AR,K-fold | Linear | | K-fold CV | Quadratic |
| RRAR,POOS-CV | Linear | Ridge | POOS CV | Quadratic |
| RRAR,K-fold | Lineal | Ridge | K-fold CV | Quadratic |
| RFAR,POOS-CV | Nonlinear | | POOS CV | Quadratic |
| RFAR,K-fold | Nonlinear | | K-fold CV | Quadratic |
| KRRAR,POOS-CV | Nonlinear | Ridge | POOS CV | Quadratic |
| KRRAR,K-fold | Nonlinear | Ridge | K-fold CV | Quadratic |
| SVR-AR,Lin,POOS-CV | Linear | | POOS CV | $\bar{\varepsilon}$-insensitive |
| SVR-AR,Lin,K-fold | Linear | | K-fold CV | $\bar{\varepsilon}$-insensitive |
| SVR-AR,RBF,POOS-CV | Nonlinear | | POOS CV | $\bar{\varepsilon}$-insensitive |
| SVR-AR,RBF,K-fold | Nonlinear | | K-fold CV | $\bar{\varepsilon}$-insensitive |
| *Data-rich models* | | | | |
| ARDI,BIC | Linear | PCA | BIC | Quadratic |
| ARDI,AIC | Linear | PCA | AIC | Quadratic |
| ARDI,POOS-CV | Linear | PCA | POOS CV | Quadratic |
| ARDI,K-fold | Linear | PCA | K-fold CV | Quadratic |
| RRARDI,POOS-CV | Linear | Ridge-PCA | POOS CV | Quadratic |
| RRARDI,K-fold | Linear | Ridge-PCA | K-fold CV | Quadratic |
| RFARDI,POOS-CV | Nonlinear | PCA | POOS CV | Quadratic |
| RFARDI,K-fold | Nonlinear | PCA | K-fold CV | Quadratic |
| KRRARDI,POOS-CV | Nonlinear | Ridge-PCR | POOS CV | Quadratic |
| KRRARDI,K-fold | Nonlinear | Ridge-PCR | K-fold CV | Quadratic |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | Linear | EN | POOS CV | Quadratic |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | Linear | EN | K-fold CV | Quadratic |
| $(B_1, \alpha = 1)$,POOS-CV | Linear | Lasso | POOS CV | Quadratic |
| $(B_1, \alpha = 1)$,K-fold | Linear | Lasso | K-fold CV | Quadratic |
| $(B_1, \alpha = 0)$,POOS-CV | Linear | Ridge | POOS CV | Quadratic |
| $(B_1, \alpha = 0)$,K-fold | Linear | Ridge | K-fold CV | Quadratic |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | Linear | EN-PCA | POOS CV | Quadratic |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | Linear | EN-PCA | K-fold CV | Quadratic |
| $(B_2, \alpha = 1)$,POOS-CV | Linear | Lasso-PCA | POOS CV | Quadratic |
| $(B_2, \alpha = 1)$,K-fold | Linear | Lasso-PCA | K-fold CV | Quadratic |
| $(B_2, \alpha = 0)$,POOS-CV | Linear | Ridge-PCA | POOS CV | Quadratic |
| $(B_2, \alpha = 0)$,K-fold | Linear | Ridge-PCA | K-fold CV | Quadratic |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | Linear | EN-PCR | POOS CV | Quadratic |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | Linear | EN-PCR | K-fold CV | Quadratic |
| $(B_3, \alpha = 1)$,POOS-CV | Linear | Lasso-PCR | POOS CV | Quadratic |
| $(B_3, \alpha = 1)$,K-fold | Linear | Lasso-PCR | K-fold CV | Quadratic |
| $(B_3, \alpha = 0)$,POOS-CV | Linear | Ridge-PCR | POOS CV | Quadratic |
| $(B_3, \alpha = 0)$,K-fold | Linear | Ridge-PCR | K-fold CV | Quadratic |
| SVR-ARDI,Lin,POOS-CV | Linear | PCA | POOS CV | $\bar{\varepsilon}$-insensitive |
| SVR-ARDI,Lin,K-fold | Linear | PCA | K-fold CV | $\bar{\varepsilon}$-insensitive |
| SVR-ARDI,RBF,POOS-CV | Nonlinear | PCA | POOS CV | $\bar{\varepsilon}$-insensitive |
| SVR-ARDI,RBF,K-fold | Nonlinear | PCA | K-fold CV | $\bar{\varepsilon}$-insensitive |

*Note: PCA stands for Principal Component Analysis, EN for Elastic Net regularizer, PCR for Principal Component Regression.*

since most of the variables do not start early enough.Second, most of the timely available series are disaggregated components of the variables in FRED-MD. Hence, adding them alters the estimation of common factors (Boivin and Ng, 2006), and induces too much collinearity for Lasso performance (Fan and Lv, 2010). Third, it is the standard high-dimensional dataset that has been extensively used in the macroeconomic literature.

## 4.2 Variables of Interest

We focus on predicting five representative macroeconomic indicators of the US economy: Industrial Production (INDPRO), Unemployment rate (UNRATE), Consumer Price Index (INF), difference between 10-year Treasury Constant Maturity rate and Federal funds rate (SPREAD) and housing starts (HOUST). INDPRO, CPI and HOUST are assumed $I(1)$ so we forecast the average growth rate as in equation (4). UNRATE is considered $I(1)$ and we target the average change as in (4) but without logs. SPREAD is $I(0)$ and the target is as in (3).[19]

## 4.3 Pseudo-Out-of-Sample Experiment Design

The pseudo-out-of-sample period is 1980M01 - 2017M12. The forecasting horizons considered are 1, 3, 9, 12 and 24 months. Hence, there are 456 evaluation periods for each horizon. All models are estimated recursively with an expanding window as means of erring on the side of including more data so as to potentially reduce the variance of more flexible models.[20]

Hyperparameter optimization is done with in-sample criteria (AIC and BIC) and two types of CV (POOS and K-fold). The in-sample selection is standard, we fix the upper bounds for the set of HPs. For the POOS CV, the validation set consists of last 25% of the in-sample. In case of K-fold CV, we set $k = 5$. We re-optimize hyperparameters every two years. This isn't uncommon for computationally demanding studies.[21] It is also reasonable to assume that optimal hyperparameters would not be terribly affected by expanding the training set with observations that account for 2-3% of the new training set size. The information on upper / lower bounds and grid search for HPs for every model is available in section E in the

---

[19]The US CPI is sometimes modeled as $I(2)$ due to the possible stochastic trend in inflation rate in 70's and 80's, see (Stock and Watson, 2002b). Since in our test set the the inflation is mostly stationary, we treat the price index as $I(1)$, as in Medeiros et al. (2019). We have compared the mean squared predictive errors of best models under $I(1)$ and $I(2)$ alternatives, and found that errors are minimized when predicting the inflation rate directly.

[20]The alternative is obviously that of a rolling window, which could be more robust to issues of model instability. These are valid concerns and have motivated tests and methods for taking them into account (see for example, Pesaran and Timmermann (2007); Pesaran et al. (2013); Inoue et al. (2017); Boot and Pick (2020)), an adequate evaluation lies beyond the scope of this paper. Moreover, as noted in Boot and Pick (2020), the number of relevant breaks may be much smaller than previously thought.

[21]Sermpinis et al. (2014), for example, split their out-of-sample into four year periods and update both hyperparameters and model parameter estimates every 4 years. Likewise, Teräsvirta (2006) selected the number of lagged values to be included in nonlinear autoregressive models once and for all at the start of the POOS.

supplementary material.

## 4.4 Forecast Evaluation Metrics

Following a standard practice in the forecasting literature, we evaluate the quality of our point forecasts using the root Mean Square Prediction Error (MSPE). Diebold and Mariano (1995) (DM) procedure is used to test the predictive accuracy of each model against the reference (ARDI,BIC). We also implement the Model Confidence Set (MCS), (Hansen et al., 2011), that selects the subset of best models at a given confidence level. These metrics measure the overall predictive performance and classify models according to DM and MCS tests. Regression analysis from section 2.3 is used to estimate the treatment effect of each ML ingredient.

# 5   Results

We present the results in several ways. First, for each variable, we summarize tables containing the relative root MSPEs (to AR,BIC model) with DM and MCS outputs, for the whole pseudo-out-of-sample and NBER recession periods. Second, we evaluate the marginal effect of important features of ML using regressions described in section 2.3.

## 5.1   Overall Predictive Performance

Tables 4 - 8, in the appendix A, summarize the overall predictive performance in terms of root MSPE relative to the reference model AR,BIC. The analysis is done for the full out-of-sample as well as for NBER recessions (i.e., when the target belongs to a recession episode). This address two questions: is ML already useful for macroeconomic forecasting and when?[22]

In case of industrial production, table 4 shows that principal component regressions $B_2$ and $B_3$ with Ridge and Lasso penalty respectively are the best at short-run horizons of 1 and 3 months. The kernel ridge ARDI with POOS CV is best for $h = 9$, while its autoregressive counterpart with K-fold minimizes the MSPE at the one-year horizon. Random forest ARDI, the alternative nonlinear approximator, outperforms the reference model by 11% for $h = 24$.

---

[22]The knowledge of the models that have performed best historically during recessions is of interest for practitioners. If the probability of recession is high enough at a given period, our results can provide an ex-ante guidance on which model is likely to perform best in such circumstances.

During recessions, the ARDI with CV is the best for 1, 3 and 9 months ahead, while the nonlinear SVR-ARDI minimizes the MSPE at the one-year horizon. The ridge regression ARDI is the best for $h = 24$. Ameliorations with respect to AR,BIC are much larger during economic downturns, and the MCS selects fewer models.

Results for the unemployment rate, table 5, highlight the performance of nonlinear models especially for longer horizons. Improvements with respect to the AR,BIC model are bigger for both full OOS and recessions. MCSs are narrower than in case of INDPRO. A similar pattern is observed during NBER recessions. Table 6 summarizes results for the Spread. Nonlinear models are generally the best, combined with data-rich predictors' set.

For inflation, table 7 shows that the kernel ridge autoregressive model with K-fold CV is the best for 3, 9 and 12 months ahead, while the nonlinear SVR-ARDI optimized with K-fold CV reduces the MSPE by more than 20% at two-year horizon. Random forest models are very resilient, as in Medeiros et al. (2019), but generally outperformed by KRR form of nonlinearity. During recessions, the fat regression models ($B_1$) are the best at short horizons, while the ridge regression ARDI with K-fold dominates for $h = 9, 12, 24$. Housing starts, in table 8, are best predicted with nonlinear data-rich models for almost all horizons.

Overall, using data-rich models and nonlinear $g$ functions improve macroeconomic prediction. Their marginal contribution depends on the state of the economy.

## 5.2   Disentangling ML Treatment Effects

The results in the previous section does not easily allow to disentangle the marginal effects of important ML features – as presented in section 3. Therefore, we turn to the regression analysis described in section 2.3. In what follows, [X, NL, SH, CV and LF] stand for data-rich, nonlinearity, alternative shrinkage, cross-validation and loss function features respectively.

Figure 1 shows the distribution of $\grave{\alpha}_F^{(h,v)}$ from equation (11) done by $(h, v)$ subsets. Hence, here we allow for heterogeneous treatment effects according to 25 different targets. This figure highlights by itself the main findings of this paper. **First**, ML nonlinearities improve substantially the forecasting accuracy in almost all situations. The effects are positive and

**Figure 1:** This figure plots the distribution of $\hat{\alpha}_F^{(h,v)}$ from equation (11) done by $(h, v)$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^2$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. Finally, variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. As an example, we clearly see that the partial effect of $X$ on the $R^2$ of INF increases drastically with the forecasted horizon $h$. SEs are HAC. These are the 95% confidence bands.

significant for all horizons in case of INDPRO and SPREAD, and for most of the cases when predicting UNRATE, INF and HOUST. The improvements of the nonlinearity treatment reach up to 23% in terms of pseudo-$R^2$. This is in contrast with previous literature that did not find substantial forecasting power from nonlinear methods, see for example Stock and Watson (1999). In fact, the ML nonlinearity is highly flexible and well disciplined by a careful regularization, and thus can solve the general overfitting problem of standard nonlinear models (Teräsvirta, 2006). This is also in line with the finding in Gu et al. (2020b) that nonlinearities (from ML models) can help predicting financial returns.

**Second**, alternative regularization means of dimensionality reduction do not improve on average over the standard factor model, except few cases. Choosing sparse modeling can decrease the forecast accuracy by up to 20% of the pseudo-$R^2$ which is not negligible. Interestingly, Gu et al. (2020b) also reach similar conclusions that dense outperforms sparse in the context of applying ML to returns.

**Third**, the average effect of CV appears not significant. However, as we will see in sec-

tion 5.2.3, the averaging in this case hides some interesting and relevant differences between K-fold and POOS CVs. **Fourth**, on average, dropping the standard in-sample squared-loss function for what the SVR proposes is not useful, except in very rare cases. **Fifth** and lastly, the marginal benefits of data-rich models ($X$) seems roughly to increase with horizons for every variable-horizon pair, except for few cases with spread and housing. Note that this is almost exactly like the picture we described for NL. Indeed, visually, it seems like the results for $X$ are a compressed-range version of NL that was translated to the right. Seeing NL models as data augmentation via basis expansions, we conclude that for predicting macroeconomic variables, we need to augment the AR($p$) model with more regressors either created from the lags of the dependent variable itself or coming from additional data. The possibility of joining these two forces to create a "data-filthy-rich" model is studied in section 5.2.1.

It turns out these findings are somewhat robust as graphs included in the appendix section B show. ML treatment effects plots of very similar shapes are obtained for data-poor models only (figure 11), data-rich models only (figure 12) and recessions / expansions periods (figures 13 and 14). It is important to notice that nonlinearity effect is not only present during recession periods, but it is even more important during expansions.[23] The only exception is the data-rich feature that has negative and significant effects for housing starts prediction when we condition on the last 20 years of the forecasting exercise (figure 15).

Figure 2 aggregates by $h$ and $v$ in order to clarify whether variable or horizon heterogeneity matters most. Two facts detailed earlier are now quite easy to see. For both $X$ and NL, the average marginal effects roughly increase in $h$. In addition, it is now clear that all the variables benefit from both additional information and nonlinearities. Alternative shrinkage is least harmful for inflation and housing, and at short horizons. Cross-validation has negative and sometimes significant impacts, while the SVR loss function is often damaging.

Supplementary material contains additional results. Section A shows the results obtained using the absolute loss. The importance of each feature and the way it behaves according to the variable/horizon pair is the same. Finally, sections B and C show results for two similar

---

[23]This suggests that our models behave relatively similarly over the business cycle and that our analysis does not suffer from undesirable forecast ranking due to extreme events as pointed out in Lerch et al. (2017).

## Distribution of averaged ML Treatment Effects

**Figure 2:** This figure plots the distribution of $\dot{\alpha}_F^{(v)}$ and $\dot{\alpha}_F^{(h)}$ from equation (11) done by $h$ and $v$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^2$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. However, in this graph, $v-$specific heterogeneity and $h-$specific heterogeneity have been integrated out in turns. SEs are HAC. These are the 95% confidence bands.

exercises. The first consider quarterly US data where we forecast the average growth rate of GDP, consumption, investment and disposable income, and the PCE inflation. The results are consistent with the findings obtained in the main body of this paper. In the second, we use a large Canadian monthly dataset and forecast the same target variables for Canada. Results are qualitatively in line with those on US data, except that NL effect is smaller in size.

In what follows we break down averages and run specific regressions as in (12) to study how homogeneous are the $\dot{\alpha}_F$'s reported above.

### 5.2.1 Nonlinearities

Figure 3 suggests that nonlinearities can be very helpful at forecasting all the five variables in the data-rich environment. The marginal effects of random forests and KRR are almost never statistically different for data-rich models, except for inflation combined with data-rich, suggesting that the common NL feature is the driving force. However, this is not the case for data-poor models where the kernel-type nonlinearity shows significant improvements

**Figure 3:** This figure compares the two NL models averaged over all horizons. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.



**Figure 4:** This figure compares the two NL models averaged over all variables. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

for all variables, while the random forests have positive impact on predicting INDPRO and inflation, but decrease forecasting accuracy for the rest of the variables.

Figure 4 suggests that nonlinearities are in general more useful for longer horizons in data-rich environment while the KRR can be harmful for a very short horizon. Note again that both nonlinear models follow the same pattern for data-rich models with random forest often being better (but never statistically different from KRR). For data-poor models, it is KRR that has a (statistically significant) growing advantage as $h$ increases. Seeing NL models as data augmentation via some basis expansions, we can join the two facts together to conclude that the need for a complex and "data-filthy-rich" model arises for predicting macroeconomic variables at longer horizons. Similar conclusions are obtained with neural networks and boosted trees as shown in figures 20 and 21 in Appendix D.

Figure 17 in the appendix C plots the cumulative and 3-year rolling window root MSPE for linear and nonlinear data-poor and data-rich models, for $h = 12$, as well as Giacomini and Rossi (2010) fluctuation test for those alternatives. The cumulative root MSPE clearly shows the positive impact on forecast accuracy of both nonlinearities and data-rich environment for all series except INF. The rolling window depicts the changing level of forecast accuracy. For all series except the SPREAD, there is a common cyclical behavior with two relatively similar peaks (1981 and 2008 recessions), as well as a drop in MSPE during the Great Moderation period. Fluctuation tests confirm the important role of nonlinear and data-rich models.

For CPI inflation at horizons of 3, 9 and 12 months, Random Forests perform distinctively well. In both its data-poor and data-rich incarnations, the algorithm is included in the superior model set of Hansen et al. (2011) and significantly outperforms the AR-BIC benchmark according to the DM test. This result can help shed some light on long standing issues in the inflation forecasting literature. A consensus emerged that nonlinear models in-sample good performance does not materialize out-of-sample (Marcellino, 2008; Stock and Watson, 2009).[24] In contrast, we found – as in Medeiros et al. (2019), that Random Forests are a particularly potent tool to forecast CPI inflation. One possible explanation is that previous studies

---

[24]Concurrently, simple benchmarks such as a random walk or moving averages emerged as surprisingly hard to beat (Atkeson and Ohanian, 2001; Stock and Watson, 2009; Kotchoni et al., 2019).

**Figure 5:** This figure compares models of section 3.2 averaged over all variables and horizons. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. The base models are ARDIs specified with POOS-CV and KF-CV respectively. SEs are HAC. These are the 95% confidence bands.

suffer from overfitting (Marcellino, 2008) while Random Forests are arguably completely immune from it (Goulet Coulombe, 2020), all this while retaining relevant nonlinearities. In that regard, it is noted that INF is the only target where KRR performance does not match that of Random Forests in the data rich environment. In the data-poor case, roles are reversed. Unlike most other targets, it seems the type of NL being used matters for inflation. Nonetheless, ML generally appears to be useful for inflation forecasting by providing better-behaved non-parametric nonlinearities than what was considered by the older literature.

### 5.2.2 Regularization

Figure 5 shows that the ARDI reduces dimensionality in a way that certainly works well with economic data: all competing schemes do at most as good on average. It is overall safe to say that on average, all shrinkage schemes give similar or lower performance, which is in line with conclusions from Stock and Watson (2012b) and Kim and Swanson (2018), but contrary to Smeekes and Wijler (2018). No clear superiority for the Bayesian versions of some of these models was also documented in De Mol et al. (2008). This suggests that the factor model view of the macroeconomy is quite accurate in the sense that when we use it as a

means of dimensionality reduction, it extracts the most relevant information to forecast the relevant time series. This is good news. The ARDI is the simplest model to run and results from the preceding section tells us that adding nonlinearities to an ARDI can be quite helpful.

Obviously, the deceiving behavior of alternative shrinkage methods does not mean there are no interesting $(h, v)$ cases where using a different dimensionality reduction has significant benefits as discussed in section 5.1 and Smeekes and Wijler (2018). Furthermore, LASSO and Ridge can still be useful to tackle specific time series problems (other than dimensionality reduction), as shown with time-varying parameters in Coulombe (2019).

### 5.2.3 Hyperparameter Optimization

Figure 6 shows how many regressors are kept by different selection methods in the case of ARDI. As expected, BIC is in general the lower envelope of each of these graphs. Both cross-validations favor larger models, especially when combined with Ridge regression. We remark a common upward trend for all model selection methods in case of INDPRO and UNRATE. This is not the case for inflation where large models have been selected in 80's and most recently since 2005. In case of HOUST, there is a downward trend since 2000's which is consistent with the finding in Figure 15 that data-poor models do better in last 20 years. POOS CV selection is more volatile and selects bigger models for unemployment rate, spread and housing. While K-fold also selects models of considerable size, it does so in a more slowly growing fashion. This is not surprising because K-fold samples from all available data to build the CV criterion: adding new data points only gradually change the average. POOS CV is a shorter window approach that offers flexibility against structural hyperparameters change at the cost of greater variance and vulnerability of rapid regime changes in the data.

We know that different model selection methods lead to quite different models, but what about their predictions? First, let us note that changes in OOS-$R^2$ are much smaller in magnitude for CV (as can be seen easily in figures 1 and 2) than for other studied ML treatment effects. Nevertheless, table 2 tells many interesting tales. The models included in the regressions are the standard linear ARs and ARDIs (that is, excluding the Ridge versions) that have all been tuned using BIC, AIC, POOS CV and CV-KF. First, we see that overall, only POOS

**Figure 6:** This figure shows the number of regressors in linear ARDI models. Results averaged across horizons.

CV is distinctively worse, especially in data-rich environment, and that AIC and CV-KF are not significantly different from BIC on average. For data-poor models and during recessions, AIC and CV-KF are being significantly better than BIC in downturns, while CV-KF seems harmless. The state-dependent effects are not significant in data-rich environment. Hence,

for that class of models, we can safely opt for either BIC or CV-KF. Assuming some degree of external validity beyond that model class, we can be reassured that the quasi-necessity of leaving ICs behind when opting for more complicated ML models is not harmful.

Table 2: CV comparison

|  | (1) All | (2) Data-rich | (3) Data-poor | (4) Data-rich | (5) Data-poor |
|---|---|---|---|---|---|
| CV-KF | -0.0380 | -0.314 | 0.237 | -0.494 | -0.181 |
|  | (0.800) | (0.711) | (0.411) | (0.759) | (0.438) |
| CV-POOS | -1.351 | -1.440* | -1.262** | -1.069 | -1.454*** |
|  | (0.800) | (0.711) | (0.411) | (0.759) | (0.438) |
| AIC | -0.509 | -0.648 | -0.370 | -0.580 | -0.812 |
|  | (0.800) | (0.711) | (0.411) | (0.759) | (0.438) |
| CV-KF * Recessions |  |  |  | 1.473 | 3.405** |
|  |  |  |  | (2.166) | (1.251) |
| CV-POOS * Recessions |  |  |  | -3.020 | 1.562 |
|  |  |  |  | (2.166) | (1.251) |
| AIC * Recessions |  |  |  | -0.550 | 3.606** |
|  |  |  |  | (2.166) | (1.251) |
| Observations | 91200 | 45600 | 45600 | 45600 | 45600 |

Standard errors in parentheses. $^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$

We now consider models that are usually tuned by CV and compare the performance of the two CVs by horizon and variables. Since we are now pooling multiple models, including all the alternative shrinkage models, if a clear pattern only attributable to a certain CV existed, it would most likely appear in figure 7. What we see are two things. First, CV-KF is at least as good as POOS CV on average for almost all variables and horizons, irrespective of the informational content of the regression. The exceptions are HOUST in data-rich and INF in data-poor frameworks, and the two-year horizon with large data. Figure 8's message has the virtue of clarity. POOS CV's failure is mostly attributable to its poor record in recessions periods for the first three variables at any horizon. Note that this is the same subset of variables that benefits from adding in more data ($X$) and nonlinearities as discussed in 5.2.1.

By using only recent data, POOS CV will be more robust to gradual structural change but will perhaps have an Achilles heel in regime switching behavior. If the optimal hyperparameters are state-dependent, then a switch from expansion to recession at time $t$ can be quite

## CV-KF performance relative to CV-POOS



**Figure 7:** This figure compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

## CV-KF performance relative to CV-POOS



**Figure 8:** This figure compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

harmful. K-fold, by taking the average over the whole sample, is less immune to such problems. Since results in 5.1 point in the direction that smaller models are better in expansions and bigger models in recessions, the behavior of CV and how it picks the effective complexity of the model can have an effect on overall predictive ability. This is exactly what we see in figure 8: POOS CV is having a hard time in recessions with respect to K-fold.

### 5.2.4 Loss Function

In this section, we investigate whether replacing the $l_2$ norm as an in-sample loss function for the SVR machinery helps in forecasting. We again use as baseline models ARs and ARDIs trained by the same corresponding CVs. The very nature of this ML feature is that the model is less sensible to extreme residuals, thanks to the $l_1$ norm outside of the $\bar{\epsilon}$-insensitivity tube. We first compare linear models in figure 9. Clearly, changing the loss function is generally harmful and that is mostly due to recessions period. However, in expansions, the linear SVR is better on average than a standard ARDI for UNRATE and SPREAD, but these small gains are clearly offset (on average) by the huge recession losses.

The SVR is usually used in its nonlinear form. We hereby compare KRR and SVR-NL to study whether the loss function effect could reverse when a nonlinear model is considered. Comparing these models makes sense since they both use the same kernel trick (with an RBF kernel). Hence, like linear models of figure 9, models in figure 10 only differ by the use of a different loss function $\hat{L}$. It turns out conclusions are exactly the same as for linear models with the negative effects being slightly smaller in nonlinear world. There are few exceptions: inflation rate and one month ahead horizon during recessions. Furthermore, figures 18 and 19 in the appendix C confirm that these findings are valid for both the data-rich and the data-poor environments.

By investigating these results more in depth using tables 4 - 8, we see an emerging pattern. First, SVR sometimes does very good (best model for UNRATE at horizon 3 months) but underperforms for many targets – in its AR or ARDI form. When it does perform well compared to the benchmark, it is more often than not outshined marginally by the KRR version. For instance, in table 5, linear and nonlinear SVR-Kfold provide respectively reductions

**Figure 9:** This graph displays the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in both recession and expansion periods. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

of 17% and 13% in RMSPE over the benchmark for UNRATE at horizon 9 months. However, analogous KRR and Random Forest similarly do so. Moreover, for targets for which SVR fails, the two models it is compared to in order to extract $\alpha_{\hat{L}}$, KRR or the AR/ARDI, have a more stable (good) record. Hence, on average nonlinear SVR is much worse than KRR and the linear SVR is also inferior to the plain ARDI. This explains the clear-cut results reported in this section: if the SVR wins, it is rather for its use of the kernel trick (nonlinearities) than an alternative in-sample loss function.

These results point out that an alternative $\hat{L}$ like the $\bar{\epsilon}$-insensitive loss function is not the most salient feature ML has to offer for macroeconomic forecasting. From a practical point of view, our results indicate that, on average, one can obtain the benefits of SVR and more by considering the much simpler KRR. This is convenient since obtaining the KRR forecast is a matter of less than 10 lines of codes implying the most straightforward form of linear algebra. In contrast, obtaining the SVR solution can be a serious numerical enterprise.

**Figure 10:** This graph displays the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in both recession and expansion periods. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

# 6    When are the ML Nonlinearities Important?

In this section we aim to explain some of the heterogeneity of ML treatment effects by interacting them in equation (12) with few macroeconomic variables $\xi_t$ that have been used to explain main sources of observed nonlinear macroeconomic fluctuations. We focus on NL feature only given its importance for both macroeconomic prediction and modeling.

The first element in $\xi_t$ is the Chicago Fed adjusted national financial conditions index (ANFCI). Adrian et al. (2019) find that lower quantiles of GDP growth are time varying and are predictable by tighter financial conditions, suggesting that higher order approximations are needed in general equilibrium models with financial frictions. In addition, Beaudry et al. (2018) build on the observation that recessions are preceded by accumulations of business, consumer and housing capital, while Beaudry et al. (2020) add nonlinearities in the estimation part of a model with financial frictions and household capital accumulation. Therefore, we add to the list the house price growth (HOUSPRICE), measured by the S&P/Case-Shiller U.S. National Home Price Index. The goal is then to test whether financial conditions and

34

capital buildups interact with the nonlinear ML feature, and if they could explain its superior performance in macroeconomic forecasting.

Uncertainty is also related to nonlinearity in macroeconomic modeling (Bloom, 2009). Benigno et al. (2013) provide a second-order approximation solution for a model with time-varying risk that has its own effect on endogenous variables. Gorodnichenko and Ng (2017) find evidence on volatility factors that are persistent and load on the housing sector, while Carriero et al. (2018) estimate uncertainty and its effects in a large nonlinear VAR model. Hence, we include the Macro Uncertainty from Jurado et al. (2015) (MACROUNCERT).[25]

Then we add measures of sentiments: University of Michigan Consumer Expectations (UMCSENT) and Purchasing Managers Index (PMI). Angeletos and La'O (2013) and Benhabib et al. (2015) have suggested that waves of pessimism and optimism play an important role in generating (nonlinear) macroeconomic fluctuations. In the case of Benhabib et al. (2015), optimal decisions based on sentiments produce multiple self-fulfilling rational expectations equilibria. Consequently, including measures of sentiment in $\xi_t$ aims to test if this channel plays a role for nonlinearities in macro forecasting. Standard monetary VAR series are used as controls: UNRATE, PCE inflation (PCEPI) and one-year treasury rate (GS1).[26]

Interactions are formed with $\xi_{t-h}$ to measure its impact when the forecast is made. This is of interest for practitioners as it indicates which macroeconomic conditions favor nonlinear ML forecast modeling. Hence, this expands the equation (12) to

$$\forall m \in \mathcal{M}_{NL}: \quad R^2_{t,h,v,m} = \dot{\alpha}_{NL} + \dot{\gamma}I(m \in NL)\xi_{t-h} + \dot{\phi}_{t,v,h} + \dot{u}_{t,h,v,m}$$

where $\mathcal{M}_{NL}$ is defined as the set of models that differs only by the use of NL.

The results are presented in table 3. The first column shows regression coefficients for $h = \{9, 12, 24\}$, since nonlinearity has been found more important for longer horizons. The second column average across all horizons, while the third presents the results for data-rich models only. The last column shows the heterogeneity of NL treatments during last 20 years.

Results show that macroeconomic uncertainty is a true game changer for ML nonlinearity

---

[25] We did not consider the Economic Policy Uncertainty from Baker et al. (2016) as it starts only from 1985.

[26] We consider GS1 instead of the federal funds rate because of the long zero lower bound period. Time series of elements in $\xi_t$ are plotted in figure 16.

Table 3: Heterogeneity of NL treatment effect

|  | (1) Base | (2) All Horizons | (3) Data-Rich | (4) Last 20 years |
|---|---|---|---|---|
| NL | 8.998*** | 5.808*** | 13.48*** | 19.87*** |
|  | (0.748) | (0.528) | (1.012) | (1.565) |
| HOUSPRICE | -9.668*** | -4.491*** | -11.56*** | -1.219 |
|  | (1.269) | (0.871) | (1.715) | (1.596) |
| ANFCI | 7.244*** | 2.625 | 6.803** | 20.29*** |
|  | (1.881) | (1.379) | (2.439) | (4.891) |
| MACROUNCERT | 17.98*** | 10.28*** | 34.87*** | 9.660*** |
|  | (1.875) | (1.414) | (2.745) | (2.038) |
| UMCSENT | 4.695** | 3.853** | 10.29*** | -3.625 |
|  | (1.768) | (1.315) | (2.294) | (1.922) |
| PMI | 0.0787 | -1.443 | -2.048 | -1.919 |
|  | (1.179) | (0.879) | (1.643) | (1.288) |
| UNRATE | 0.834 | 2.517** | 5.732*** | 8.526*** |
|  | (1.353) | (0.938) | (1.734) | (2.199) |
| GS1 | -14.24*** | -9.500*** | -17.30*** | 2.081 |
|  | (2.288) | (1.682) | (3.208) | (3.390) |
| PCEPI | 5.953* | 6.814** | -1.142 | -6.242 |
|  | (2.828) | (2.180) | (4.093) | (3.888) |
| Observations | 136800 | 228000 | 68400 | 72300 |

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

as it improves its forecast accuracy by 34% in the case of data-rich models. This means that if the macro uncertainty goes from -1 standard deviation to +1 standard deviation from its mean, the expected NL treatment effect (in terms OOS-$R^2$ difference) is 2*34=+68%. Tighter financial conditions and a decrease in house prices are also positively correlated with a higher NL treatment, which supports the findings in Adrian et al. (2019) and Beaudry et al. (2020). It is particularly interesting that the effect of ANFCI reaches 20% during last 20 years, while the impact of uncertainty decreases to less than 10%, emphasizing that the determinant role of financial conditions in recent US macro history is also reflected in our results. Waves of consumer optimism positively affect nonlinearities, especially with data-rich models.

Among control variables, unemployment rate has a positive effect on nonlinearity. As expected, this suggests that the importance of nonlinearities is a cyclical feature. Lower interest rates also improve NL treatment by as much as 17% in the data-rich setup. Higher inflation also leads to stronger gains from ML nonlinearities, but mainly at shorter horizons and for

data-poor models, as suggested by comparing specifications (2) and (3).

These results document clear historical situations where NL consistently helps: (i) when the level of macroeconomic uncertainty is high and (ii) during episodes of tighter financial conditions and housing bubble bursts.[27] Also, we note that effects are often bigger in the case of data-rich models. Hence, allowing nonlinear relationship between factors made of many predictors can capture better the complex relationships that characterize the episodes above.

These findings suggest that ML captures important macroeconomic nonlinearities, especially in the context of financial frictions and high macroeconomic uncertainty. They can also serve as guidance for forecasters that use a portfolio of predictive models: one should put more weight on nonlinear specifications if economic conditions evolve as described above.

## 7 Conclusion

In this paper we have studied important features driving the performance of machine learning techniques in the context of macroeconomic forecasting. We have considered many ML methods in a substantive POOS setup over 38 years for 5 key variables and 5 horizons. We have classified these models by "features" of machine learning: nonlinearities, regularization, cross-validation and alternative loss function. The data-rich and data-poor environments were considered. In order to recover their marginal effects on forecasting performance, we designed a series of experiments that easily allow to identify the treatment effects of interest.

The first result indicates that nonlinearities are the true game changer for the data-rich environment, as they improve substantially the forecasting accuracy for all macroeconomic variables in our exercise and especially when predicting at long horizons. This gives a stark recommendation for practitioners. It recommends for most variables and horizons what is in the end a partially nonlinear factor model – that is, factors are still obtained by PCA. The best of ML (at least of what considered here) can be obtained by simply generating the data for a standard ARDI model and then feed it into a ML nonlinear function of choice. The perfor-

---

[27]Granziera and Sekhposyan (2019) have exploited similar regression setup for model selection and found that 'economic' forecasting models, AR augmented by few macroeconomic indicators, outperform the time series models during turbulent times (recessions, tight financial conditions and high uncertainty).

mance of nonlinear models is magnified during periods of high macroeconomic uncertainty, financial stress and housing bubble bursts. These findings suggest that Machine Learning is useful for macroeconomic forecasting by mostly capturing important nonlinearities that arise in the context of uncertainty and financial frictions.

The second result is that the standard factor model remains the best regularization. Alternative regularization schemes are most of the time harmful. Third, if cross-validation has to be applied to select models' features, the best practice is the standard K-fold. Finally, the standard $L_2$ is preferred to the $\bar{\epsilon}$-insensitive loss function for macroeconomic predictions. We found that most (if not all) the benefits from the use of SVR in fact comes from the nonlinearities it creates via the kernel trick rather than its use of an alternative loss function.

# References

Abadie, A. and Kasy, M. (2019). Choosing among regularized estimators in empirical economics: The risk of machine learning. *Review of Economics and Statistics*, 101(5):743–762.

Adrian, T., Boyarchenko, N., and Giannone, D. (2019). Vulnerable growth. *American Economic Review*, 109(4):1263–1289.

Ahmed, N. K., Atiya, A. F., El Gayar, N., and El-Shishiny, H. (2010). An empirical comparison of machine learning models for time series forecasting. *Econometric Reviews*, 29(5):594–621.

Alquier, P., Li, X., and Wintenberger, O. (2013). Prediction of time series by statistical learning: General losses and fast rates. *Dependence Modeling*, 1(1):65–93.

Angeletos, G.-M. and La'O, J. (2013). Sentiments. *Econometrica*, 81(2):739–779.

Athey, S. (2019). The Impact of Machine Learning on Economics. In Agrawal, A., Gans, J., and Goldfarb, A., editors, *The Economics of Artificial Intelligence: An Agenda*, pages 507–552. University of Chicago Press.

Atkeson, A. and Ohanian, L. E. (2001). Are Phillips Curves Useful for Forecasting Inflation? *Quarterly Review*, 25(1):2–11.

Baker, S. R., Bloom, N., and Davis, S. J. (2016). Measuring Economic Policy Uncertainty. *The Quarterly Journal of Economics*, 131(4):1593–1636.

Beaudry, P., Galizia, D., and Portier, F. (2018). Reconciling Hayek's and Keynes'views of recessions. *Review of Economic Studies*, 85(1):119–156.

Beaudry, P., Galizia, D., and Portier, F. (2020). Putting the cycle back into business cycle analysis. *American Economic Review*, 110(1):1–47.

Belloni, A., Chernozhukov, V., Fernandes-Val, I., and Hansen, C. B. (2017). Program Evaluation and Causal Inference With High-Dimensional Data. *Econometrica*, 85(1):233–298.

Benhabib, J., Wang, P., and Wen, Y. (2015). Sentiments and Aggregate Demand Fluctuations. *Econometrica*, 83(2):549–585.

Benigno, G., Benigno, P., and Nisticò, S. (2013). Second-order approximation of dynamic models with time-varying risk. *Journal of Economic Dynamics and Control*, 37(7):1231–1247.

Bergmeir, C. and Benítez, J. M. (2012). On the use of cross-validation for time series predictor evaluation. *Information Sciences*, 191:192–213.

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 and Data Analysis*, 120:70–83.

Bloom, N. (2009). The Impact of Uncertainty Shocks. *Econometrica*, 77(3):623–685.

Boivin, J. and Ng, S. (2006). Are More Data Always Better for Factor Analysis? *Journal of Econometrics*, 132(1):169–194.

Boot, T. and Pick, A. (2020). Does Modeling a Structural Break Improve Forecast Accuracy? *Journal of Econometrics*, 215(1):35–59.

Bordo, M. D., Redish, A., and Rockoff, H. (2015). Why didn't Canada have a banking crisis in 2008 (or in 1930, or 1907, or...)? *Economic History Review*, 68(1):218–243.

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

Carriero, A., Clark, T. E., and Marcellino, M. (2018). Measuring Uncertainty and Its Impact on the Economy. *Review of Economics and Statistics*, 100(5):799–815.

Carriero, A., Galvão, A. B., and Kapetanios, G. (2019). A Comprehensive Evaluation of Macroeconomic Forecasting Methods. *International Journal of Forecasting*, 35(4):1226 – 1239.

Chen, J., Dunn, A., Hood, K., and Batch, A. (2019). Off to the Races : A Comparison of Machine Learning and Alternative Data for Nowcasting of Economic Indicators. In Abraham, K., Jarmin, R. S., Moyer, B., and Shapiro, M. D., editors, *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.

Choudhury, S., Ghosh, S., Bhattacharya, A., Fernandes, K. J., and Tiwari, M. K. (2014). A real time clustering and SVM based price-volatility prediction for optimal trading strategy. *Neurocomputing*, 131:419–426.

Claeskens, G. and Hjort, N. L. (2008). Akaike's Information Criterion. In Claeskens, G. and

Hjort, N. L., editors, *Model Averaging and Model Selection*, chapter 2, pages 22–69.

Colombo, E. and Pelagatti, M. (2020). Statistical learning and exchange rate forecasting. *International Journal of Forecasting*, xxx(xxxx):1–30.

Cook, T. and Smalter Hall, A. (2017). Macroeconomic Indicator Forecasting with Deep Neural Networks.

Coulombe, P. G. (2019). Time-varying Parameters : A Machine Learning Approach.

De Mol, C., Giannone, D., and Reichlin, L. (2008). Forecasting using a large number of predictors: Is Bayesian shrinkage a valid alternative to principal components? *Journal of Econometrics*, 146(2):318–328.

Diebold, F. X. and Mariano, R. S. (1995). Comparing predictive accuracy. *Journal of Business and Economic Statistics*, 13(3):253–263.

Diebold, F. X. and Shin, M. (2019). Machine learning for regularized survey forecast combination: Partially-egalitarian LASSO and its derivatives. *International Journal of Forecasting*, 35(4):1679–1691.

Döpke, J., Fritsche, U., and Pierdzioch, C. (2017). Predicting recessions with boosted regression trees. *International Journal of Forecasting*, 33(4):745–759.

Exterkate, P., Groenen, P. J. F., Heij, C., and van Dijk, D. (2016). Nonlinear forecasting with many predictors using kernel ridge regression. *International Journal of Forecasting*, 32(3):736–753.

Fan, J. and Lv, J. (2010). A Selective Overview of Variable Selection in High Dimensional Feature Space. *Statistica Sinica*, 20(1):101–148.

Fortin-Gagnon, O., Leroux, M., Stevanovic, D., and Surprenant, S. (2020). A Large Canadian Database for Macroeconomic Analysis.

Giacomini, R. and Rossi, B. (2010). Forecast Comparisons in Unstable Environments. *Journal of Applied Econometrics*, 25(4):595 – 620.

Giannone, D., Lenza, M., and Primiceri, G. E. (2015). Prior Selection for Vector Autoregressions. *Review of Economics and Statistics*, 97(2):436–451.

Giannone, D., Lenza, M., and Primiceri, G. E. (2018). Economic Predictions with Big Data: The Illusion of Sparsity.

Gorodnichenko, Y. and Ng, S. (2017). Level and volatility factors in macroeconomic data. *Journal of Monetary Economics*, 91:52–68.

Goulet Coulombe, P. (2020). To Bag is to Prune. *arXiv preprint arXiv:2008.07063*.

Goulet Coulombe, P., Leroux, M., Stevanovic, D., and Surprenant, S. (2020). Macroeconomic

Data Transformations Matter. *arXiv preprint arXiv:2008.01714*.

Granger, C. W. J. and Jeon, Y. (2004). Thick Modeling. *Economic Modelling*, 21(2):323–343.

Granziera, E. and Sekhposyan, T. (2019). Predicting relative forecasting performance: An empirical investigation. *International Journal of Forecasting*, 35(4):1636–1657.

Gu, S., Kelly, B., and Xiu, D. (2020a). Autoencoder asset pricing models. *Journal of Econometrics*, 0(0).

Gu, S., Kelly, B., and Xiu, D. (2020b). Empirical Asset Pricing via Machine Learning. *Review of Financial Studies*, 33(5):2223–2273.

Hansen, P. R., Lunde, A., and Nason, J. M. (2011). The Model Confidence Set. *Econometrica*, 79(2):453–497.

Hansen, P. R. and Timmermann, A. (2015). Equivalence Between Out-of-Sample Forecast Comparisons and Wald Statistics. *Econometrica*, 83(6):2485–2505.

Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Interference, and Prediction*. Springer Science & Business Media, second edition.

Inoue, A., Jin, L., and Rossi, B. (2017). Rolling Window Selection for out-of-sample Forecasting with Time-varying Parameters. *Journal of Econometrics*, 196(1):55–67.

Joseph, A. (2019). Parametric Inference with Universal Function Approximators.

Jurado, K., Ludvigson, S. C., and Ng, S. (2015). Measuring Uncertainty. *American Economic Review*, 105(3):1177–1216.

Kim, H. H. and Swanson, N. R. (2018). Mining big data using parsimonious factor, machine learning, variable selection and shrinkage methods. *International Journal of Forecasting*, 34(2):339–354.

Koenker, R. and Machado, J. A. (1999). Goodness of Fit and Related Inference Processes for Quantile Regression. *Journal of the American Statistical Association*, 94(448):1296–1310.

Kotchoni, R., Leroux, M., and Stevanovic, D. (2019). Macroeconomic forecast accuracy in a data-rich environment. *Journal of Applied Econometrics*, 34(7):1050–1072.

Kuan, C. M. and White, H. (1994). Artificial neural networks: An econometric perspective. *Econometric Reviews*, 13(1).

Kuznetsov, V. and Mohri, M. (2015). Learning theory and algorithms for forecasting non-stationary time series. In *Advances in Neural Information Processing Systems*, pages 541–549.

Lee, T. H., White, H., and Granger, C. W. J. (1993). Testing for neglected nonlinearity in time series models. A comparison of neural network methods and alternative tests. *Journal of Econometrics*, 56(3):269–290.

Lerch, S., Thorarinsdottir, T. L., Ravazzolo, F., and Gneiting, T. (2017). Forecaster's dilemma: Extreme events and forecast evaluation. *Statistical Science*, 32(1):106–127.

Li, J. and Chen, W. (2014). Forecasting macroeconomic time series: LASSO-based approaches and their forecast combinations with dynamic factor models. *International Journal of Forecasting*, 30(4):996–1015.

Litterman, R. B. (1979). Techniques of Forecasting Using Vector Autoregressions.

Lu, C. J., Lee, T. S., and Chiu, C. C. (2009). Financial time series forecasting using independent component analysis and support vector regression. *Decision Support Systems*, 47(2):115–125.

Marcellino, M. (2008). A linear benchmark for forecasting GDP growth and inflation? *Journal of Forecasting*, 27(4):305–340.

Marcellino, M., Stock, J. H., and Watson, M. W. (2006). A comparison of Direct and Iterated Multistep AR Methods for Forecasting Macroeconomic Time Series. *Journal of Econometrics*, 135(1-2):499–526.

McCracken, M. W. and Ng, S. (2016). FRED-MD: A Monthly Database for Macroeconomic Research. *Journal of Business and Economic Statistics*, 34(4):574–589.

Medeiros, M. C., Teräsvirta, T., and Rech, G. (2006). Building Neural Network Models for Time Series: A Statistical Approach. *Journal of Forecasting*.

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 and Economic Statistics*, 0(0):1–45.

Milunovich, G. (2020). Forecasting Australia's real house price index: A comparison of time series and machine learning methods. *Journal of Forecasting*, pages 1–21.

Mohri, M. and Rostamizadeh, A. (2010). Stability bounds for stationary $\phi$-mixing and $\beta$-mixing processes. *Journal of Machine Learning Research*, 11:789–814.

Moshiri, S. and Cameron, N. (2000). Neural network versus econometric models in forecasting inflation. *Journal of Forecasting*, 19(3):201–217.

Nakamura, E. (2005). Inflation forecasting using a neural network. *Economics Letters*, 86(3):373–378.

Ng, S. (2014). Viewpoint: Boosting recessions. *Canadian Journal of Economics*, 47(1):1–34.

Patel, J., Shah, S., Thakkar, P., and Kotecha, K. (2015a). Predicting stock and stock price index movement using Trend Deterministic Data Preparation and machine learning techniques. *Expert Systems with Applications*, 42(1):259–268.

Patel, J., Shah, S., Thakkar, P., and Kotecha, K. (2015b). Predicting stock market index using

fusion of machine learning techniques. *Expert Systems with Applications*, 42(4):2162–2172.

Pesaran, M. H., Pick, A., and Pranovich, M. (2013). Optimal forecasts in the presence of structural breaks. *Journal of Econometrics*, 177(2):134–152.

Pesaran, M. H. and Timmermann, A. (2007). Selection of estimation window in the presence of breaks. *Journal of Econometrics*, 137(1):134–161.

Qu, H. and Zhang, Y. (2016). A new kernel of support vector regression for forecasting high-frequency stock returns. *Mathematical Problems in Engineering*, pages 1–9.

Sermpinis, G., Stasinakis, C., Theofilatos, K., and Karathanasopoulos, A. (2014). Inflation and unemployment forecasting with genetic support vector regression. *Journal of Forecasting*, 33(6):471–487.

Smeekes, S. and Wijler, E. (2018). Macroeconomic forecasting using penalized regression methods. *International Journal of Forecasting*, 34(3):408–430.

Smola, A. J., Murata, N., Schölkopf, B., and Müller, K.-R. (1998). Asymptotically Optimal Choice of $\epsilon$-Loss for Support Vector Machines. In *International Conference on Artificial Neural Networks*, number 2, pages 105–110, London. Springer.

Smola, A. J. and Schölkopf, B. (2004). A Tutorial on Support Vector Regression. *Statistics and Computing*, 14:199–222.

Stock, J. H. and Watson, M. W. (1999). A Comparison of Linear and Nonlinear Univariate Models for Forecasting Macroeconomic Time Series. In Engle, R. F. and White, H., editors, *Cointegration, Causality and Forecasting: A Festschrift for Clive W.J. Granger*, pages 1–44. Oxford University Press, Oxford.

Stock, J. H. and Watson, M. W. (2002a). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association*, 97(460):1167–1179.

Stock, J. H. and Watson, M. W. (2002b). Macroeconomic forecasting using diffusion indexes. *Journal of Business and Economic Statistics*, 20(2):147–162.

Stock, J. H. and Watson, M. W. (2009). Phillips Curve Inflation Forecasts. In Fuhrer, J., Kodrzycki, Y. K., Sneddon Little, J., and Olivei, G. P., editors, *Understanding Inflation and the Implication for Monetary Policy*, chapter 3, pages 99–202. MIT Press, Cambridge, Massachusetts.

Stock, J. H. and Watson, M. W. (2012a). Disentangling the Channels of the 2007–09 Recession. *Brookings Papers on Economic Activity*, (1):81–156.

Stock, J. H. and Watson, M. W. (2012b). Generalized Shrinkage Methods for Forecasting Using Many Predictors. *Journal of Business and Economic Statistics*, 30(4):481–493.

Swanson, N. R. and White, H. (1997). A Model Selection Approach To Real-Time Macroe-

conomic Forecasting Using Linear Models And Artificial Neural Networks. *The Review of Economics and Statistics*, 79(4):540–550.

Tashman, L. J. (2000). Out-of-sample Tests of Forecasting Accuracy: An Analysis and Review. *International Journal of Forecasting*, 16(4):437–450.

Teräsvirta, T. (2006). Forecasting Economic Variables with Nonlinear Models. In Granger, C. W. J. and Elliott, G., editors, *Handbook of Economic Forecasting*, chapter 8, pages 413–457. Elsevier.

Trapletti, A., Leisch, F., and Hornik, K. (2000). Stationary and Integrated Autoregressive Neural Network Processes. *Neural Computation*, 12(10):2427–2450.

Uddin, M. F., Lee, J., Rizvi, S., and Hamada, S. (2018). Proposing enhanced feature engineering and a selection model for machine learning processes. *Applied Sciences (Switzerland)*, 8(4):1–32.

Yeh, C. Y., Huang, C. W., and Lee, S. J. (2011). A multiple-kernel support vector regression approach for stock market price forecasting. *Expert Systems with Applications*, 38(3):2177–2186.

Yousuf, K. and Ng, S. (2019). Boosting High Dimensional Predictive Regressions with Time Varying Parameters.

Zhang, X. R., Hu, L. Y., and Wang, Z. S. (2010). Multiple kernel support vector regression for economic forecasting. *2010 International Conference on Management Science and Engineering, ICMSE 2010*, (70872025):129–134.

Zhao, Q. and Hastie, T. (2019). Causal Interpretations of Black-Box Models. *Journal of Business and Economic Statistics*, 0(0):1–19.

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

Zou, H., Hastie, T., and Tibshirani, R. (2007). On the "degrees of freedom" of the lasso. *Annals of Statistics*, 35(5):2173–2192.

# A  Detailed Overall Predictive Performance

## Table 4: Industrial Production: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=3 | h=9 | h=12 | h=24 | h=1 | h=3 | h=9 | h=12 | h=24 |
| Data-poor ($H_t^-$) models | | | | | | | | | | |
| AR,BIC (RMSPE) | 0.0765 | **0.0515** | **0.0451** | **0.0428** | **0.0344** | 0.127 | 0.1014 | 0.0973 | 0.0898 | 0.0571 |
| AR,AIC | 0.991* | **1.000** | **0.999** | **1.000** | **1.000** | 0.987* | 1.000 | 1.000 | 1.000 | 1.000 |
| AR,POOS-CV | 0.999 | 1.021*** | **0.985*** | **1.001** | 1.032* | 1.01 | 1.023*** | 0.988* | 1.000 | 1.076** |
| AR,K-fold | 0.991* | **1.000** | **0.987*** | **1.000** | 1.033* | 0.987* | 1.000 | 0.992* | 1.000 | 1.078** |
| RRAR,POOS-CV | 1.003 | 1.041** | **0.989** | **0.993*** | **1.002** | 1.039** | 1.083** | 0.991 | 0.993 | 1.016** |
| RRAR,K-fold | **0.988**** | **1.000** | **0.991** | **1.001** | 1.027 | 0.992 | 1.007** | 0.995 | 1.001** | 1.074** |
| RFAR,POOS-CV | 0.995 | 1.045 | **0.985** | 0.955 | **0.991** | 1.009 | 1.073 | 0.902*** | 0.890** | 0.983 |
| RFAR,K-fold | 0.995 | **1.020** | **0.960** | **0.930**** | **0.983** | 0.999 | 1.013 | 0.894*** | 0.887*** | 0.970* |
| KRR-AR,POOS-CV | 1.023 | 1.09 | **0.980** | 0.944 | **0.982** | 1.117 | 1.166* | 0.896** | 0.853*** | 0.903*** |
| KRR,AR,K-fold | **0.947**** | **0.937**** | **0.936** | _**0.910***_ | **0.959** | 0.922** | 0.902** | 0.835*** | **0.799*** | 0.864*** |
| SVR-AR,Lin,POOS-CV | 1.134*** | 1.226*** | 1.114*** | 1.132*** | **0.952*** | 1.186** | 1.285*** | 1.079** | 1.034*** | 0.893*** |
| SVR-AR,Lin,K-fold | 1.069* | 1.159** | 1.055** | 1.042*** | 1.016*** | 1.268*** | 1.319*** | 1.067*** | 1.035*** | 1.013*** |
| SVR-AR,RBF,POOS-CV | 0.999 | 1.061*** | **1.020** | 1.048 | **0.980** | 1.062* | 1.082*** | 0.876*** | 0.941*** | 0.930*** |
| SVR-AR,RBF,K-fold | **0.978*** | **1.004** | 1.080* | 1.193** | 1.017*** | 0.992 | 1.009 | 0.989 | 1.016*** | 1.012*** |
| Data-rich ($H_t^+$) models | | | | | | | | | | |
| ARDI,BIC | **0.946*** | **0.991** | 1.037 | **1.004** | **0.968** | **0.801*** | **0.807*** | 0.887** | 0.833*** | 0.784*** |
| ARDI,AIC | **0.959*** | **0.968** | 1.017 | 0.998 | **0.943** | 0.840*** | **0.803*** | 0.844** | **0.798*** | **0.768*** |
| ARDI,POOS-CV | 0.994 | **1.015** | **0.984** | 0.968 | **0.966** | 0.896*** | _**0.698***_ | _**0.773***_ | 0.777*** | 0.812*** |
| ARDI,K-fold | **0.940*** | **0.977** | 1.013 | 0.982 | **0.912*** | _**0.787***_ | 0.812*** | 0.841** | 0.808** | 0.762*** |
| RRARDI,POOS-CV | **0.994** | **1.032** | **0.987** | 0.973 | **0.948** | 0.908** | **0.725*** | 0.793*** | 0.778*** | 0.861** |
| RRARDI,K-fold | **0.943**** | **0.977** | **0.986** | **0.990** | **0.921** | 0.847** | **0.718*** | **0.794*** | 0.796*** | _**0.702***_ |
| RFARDI,POOS-CV | **0.948**** | **0.991** | **0.951** | **0.919*** | **0.899**** | 0.865** | **0.802*** | 0.837*** | 0.782*** | 0.819*** |
| RFARDI,K-fold | **0.953**** | **1.016** | **0.957** | **0.924*** | _**0.890**_** | 0.889*** | 0.864* | 0.846*** | 0.803*** | 0.767*** |
| KRR-ARDI,POOS-CV | 1.038 | **1.016** | _**0.921***_ | 0.934 | **0.959** | 1.152* | 1.021 | 0.847*** | 0.814*** | 0.886** |
| KRR,ARDI,K-fold | **0.971** | **0.983** | **0.923*** | **0.914*** | **0.959** | 1.006 | 0.983 | 0.827*** | 0.793*** | 0.848*** |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | 1.014 | **1.001** | 1.023 | 0.996 | **0.946** | 1.067 | 0.956 | 0.979 | 0.916** | 0.855*** |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | **0.957**** | **0.952** | 1.029 | 1.046 | 1.051 | 0.908** | 0.856*** | 0.874** | **0.816*** | 0.890* |
| $(B_1, \alpha = 1)$,POOS-CV | **0.971*** | **1.013** | 1.067* | **1.020** | **0.955** | 0.991 | 0.889 | 1.01 | 0.935* | 0.880** |
| $(B_1, \alpha = 1)$,K-fold | **0.957**** | **0.952** | **1.029** | 1.046 | 1.051 | 0.908** | 0.856*** | 0.874** | **0.816*** | 0.890* |
| $(B_1, \alpha = 0)$,POOS-CV | 1.047 | 1.112** | **1.021** | 1.051 | **0.969** | 1.134* | 1.182** | 0.997 | 1.005 | 0.821*** |
| $(B_1, \alpha = 0)$,K-fold | 1.025 | 1.056* | 1.065 | 1.082 | 1.052 | 1.032 | 0.974 | 0.923 | 0.929 | 0.847*** |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | 1.061 | **0.968** | **0.975** | 0.999 | **0.923*** | 1.237 | 0.810*** | 0.889*** | 0.904** | 0.869** |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | 1.098 | **0.949** | **0.993** | 0.974 | **0.970** | 1.332 | **0.801*** | 0.896** | 0.851*** | **0.756*** |
| $(B_2, \alpha = 1)$,POOS-CV | **0.973** | 1.045 | **1.012** | **1.023** | **0.920*** | 1.034 | 1.033 | 0.997 | 0.957 | 0.839*** |
| $(B_2, \alpha = 1)$,K-fold | **0.956*** | 1.022 | 1.032 | **1.025** | **0.990** | 0.961 | 0.935 | 0.959 | 0.913** | 0.809*** |
| $(B_2, \alpha = 0)$,POOS-CV | _**0.933***_ | **0.955** | 0.972 | 0.937 | **0.913*** | 0.902** | **0.781*** | 0.904** | 0.840*** | 0.807*** |
| $(B_2, \alpha = 0)$,K-fold | **0.937*** | **0.927*** | **0.961** | **0.927** | **0.959** | 0.871*** | **0.787*** | 0.858*** | **0.775*** | 0.776*** |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | 0.980 | **0.994** | 1.016 | 1.05 | **0.952** | 1.032 | 0.95 | 0.957 | 0.97 | 0.861*** |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | **0.973*** | **0.946*** | 1.042 | **0.948** | **0.997** | 1.016 | 0.916** | 0.938 | 0.825*** | 0.827*** |
| $(B_3, \alpha = 1)$,POOS-CV | **0.969*** | 1.053 | 1.053 | 1.080* | **0.956** | 0.972 | 0.946 | 1.002 | 1.014 | 0.906** |
| $(B_3, \alpha = 1)$,K-fold | **0.946*** | _**0.913***_ | **0.994** | 0.976 | 1.01 | 0.924** | 0.829*** | 0.888* | **0.803*** | 0.822*** |
| $(B_3, \alpha = 0)$,POOS-CV | **0.976** | 1.049 | 1.04 | 1.063 | **0.973** | 1.034 | 1.061 | 0.997 | 0.932* | 0.846*** |
| $(B_3, \alpha = 0)$,K-fold | 0.981 | 1.01 | 1.03 | **1.011** | **0.985** | 1.002 | 0.997 | 0.95 | **0.826*** | 0.787*** |
| SVR-ARDI,Lin,POOS-CV | **0.989** | 1.165** | 1.216** | 1.193** | **1.034** | 0.915* | 0.900** | 1.006 | 0.862** | **0.778*** |
| SVR-ARDI,Lin,K-fold | 1.109** | 1.367*** | **1.024** | **1.038** | 1.028 | 1.129 | 1.133 | **0.776*** | **0.808*** | **0.726*** |
| SVR-ARDI,RBF,POOS-CV | **0.968*** | **0.986** | 1.100* | 0.960 | **0.936*** | 0.958 | 0.900* | 0.873** | _**0.760***_ | 0.820*** |
| SVR-ARDI,RBF,K-fold | **0.951*** | **0.946** | **0.993** | 0.952 | 1.001 | 0.860** | **0.793*** | 0.806*** | 0.777*** | 0.791*** |

*Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while \*\*\*. \*\*. \* stand for 1%. 5% and 10% significance of Diebold-Mariano test.*

# Table 5: Unemployment rate: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=3 | h=9 | h=12 | h=24 | h=1 | h=3 | h=9 | h=12 | h=24 |
| Data-poor $(H_t^-)$ models | | | | | | | | | | |
| AR,BIC (RMSPE) | 1.9578 | 1.1905 | 1.0169 | 1.0058 | 0.869 | 2.5318 | 2.0826 | 1.8823 | 1.7276 | 1.0562 |
| AR,AIC | 0.991 | 0.984 | 0.988 | 0.993*** | 1.000 | 0.958 | 0.960** | 0.984* | 1.000 | 1.000 |
| AR,POOS-CV | 0.988 | 0.999 | 1.002 | 0.995 | 0.987 | 0.978 | 0.980** | 0.996 | 0.998 | 1.04 |
| AR,K-fold | 0.994 | 0.984 | 0.989 | 0.986*** | 0.991 | 0.956* | 0.960** | 0.998 | 1.000 | 1.038 |
| RRAR,POOS-CV | 0.989 | 1.000 | 1.002 | 0.990* | 0.972** | 0.984 | 0.988* | 0.997 | 0.991* | 1.001 |
| RRAR,K-fold | 0.988 | 0.982* | 0.983* | 0.989** | 0.999 | 0.963 | 0.971* | 0.992 | 0.995 | 1.033 |
| RFAR,POOS-CV | 0.983 | 0.995 | 0.968 | 1.000 | 1.002 | 0.989 | 1.003 | 0.929** | 0.951** | 0.994 |
| RFAR,K-fold | 0.98 | 0.985 | 0.979 | 1.006 | 0.99 | 0.985 | 0.972 | 0.896*** | 0.943* | 0.983 |
| KRR-AR,POOS-CV | 0.99 | 1.04 | **0.882*** | **0.889*** | 0.876*** | 1.04 | 1.116 | 0.843*** | 0.883*** | 0.904** |
| KRR,AR,K-fold | **0.940*** | **0.910*** | **0.878*** | **0.869*** | 0.852*** | 0.847*** | 0.838*** | **0.788*** | **0.798*** | 0.908** |
| SVR-AR,Lin,POOS-CV | 1.028 | 1.133** | 1.130*** | 1.108*** | 1.174*** | 1.065* | 1.274*** | 1.137*** | 1.094*** | 1.185*** |
| SVR-AR,Lin,K-fold | 0.993 | 1.061** | 1.068*** | 1.045*** | 1.013*** | 1.062** | 1.108*** | 1.032** | 1.011 | 1.018*** |
| SVR-AR,RBF,POOS-CV | 1.019 | 1.094* | 1.029 | 1.076** | 1.01 | 1.097** | 1.247** | 1.047* | 1.034*** | 1.112* |
| SVR-AR,RBF,K-fold | 0.997 | 1.011 | 1.078** | 1.053* | 0.993 | 1.026 | 1.009 | 1.058 | 1.023 | 0.985 |
| Data-rich $(H_t^+)$ models | | | | | | | | | | |
| ARDI,BIC | **0.937**** | **0.893**** | 0.938 | 0.939 | 0.875*** | **0.690*** | 0.715*** | **0.798*** | 0.782*** | **0.783*** |
| ARDI,AIC | **0.933**** | **0.878*** | 0.928 | 0.953 | 0.893** | **0.720*** | 0.719*** | **0.798*** | 0.799*** | **0.787*** |
| ARDI,POOS-CV | **0.924*** | **0.913* | 0.957 | 0.925* | **0.856*** | <u>**0.686***</u> | **0.676*** | 0.840** | **0.737*** | **0.777*** |
| ARDI,K-fold | **0.935**** | **0.895**** | 0.929 | 0.93 | 0.915** | **0.696*** | 0.697*** | **0.801*** | 0.807*** | **0.787*** |
| RRARDI,POOS-CV | **0.924*** | **0.896* | 0.968 | 0.946 | 0.870*** | **0.711*** | <u>**0.635***</u> | 0.849** | **0.768*** | **0.767*** |
| RRARDI,K-fold | **0.940**** | **0.899**** | 0.946 | 0.931* | 0.908** | **0.755** | **0.681*** | **0.803*** | 0.790*** | **0.753*** |
| RFARDI,POOS-CV | **0.934*** | 0.945 | **0.857*** | **0.842*** | **0.763*** | **0.724*** | 0.769*** | **0.718*** | **0.734*** | **0.722*** |
| RFARDI,K-fold | **0.932*** | **0.897*** | **0.873** | **0.854*** | **0.785*** | **0.749*** | 0.742*** | **0.731*** | <u>**0.720***</u> | <u>**0.710***</u> |
| KRR-ARDI,POOS-CV | **0.959* | **0.961 | **0.839*** | <u>**0.813***</u> | **0.804*** | 1.01 | 1.017 | **0.748*** | **0.732*** | 0.828*** |
| KRR,ARDI,K-fold | **0.938*** | **0.907** | <u>**0.827***</u> | **0.817*** | **0.795*** | 0.925 | 0.933 | **0.785*** | **0.729*** | **0.814*** |
| $(B_1.\alpha=\hat{\alpha})$,POOS-CV | 0.979 | **0.945 | 0.976 | 0.953 | 0.913*** | 1.049 | 0.899* | 0.933 | 0.910* | 0.871*** |
| $(B_1.\alpha=\hat{\alpha})$,K-fold | 0.971 | **0.925** | **0.867*** | 0.919* | 0.925* | 0.787*** | 0.848*** | 0.840*** | 0.839*** | 0.829** |
| $(B_1.\alpha=1)$,POOS-CV | 0.947*** | **0.937* | 0.962 | 0.922** | 0.889*** | 0.857** | 0.789*** | 0.888** | 0.860*** | 0.915* |
| $(B_1.\alpha=1)$,K-fold | 0.971 | **0.925** | **0.867*** | 0.919* | 0.925* | 0.787*** | 0.848*** | 0.840*** | 0.839*** | 0.829** |
| $(B_1.\alpha=0)$,POOS-CV | 1.238** | 1.319** | 1.021 | 1.07 | 1.01 | 1.393* | 1.476* | 0.979 | 0.972 | **0.764*** |
| $(B_1.\alpha=0)$,K-fold | 1.246** | 0.994 | 1.062* | 1.077* | 1.018 | 1.322 | 0.963 | 0.991 | 0.933 | 0.802*** |
| $(B_2,\alpha=\hat{\alpha})$,POOS-CV | <u>**0.907***</u> | **0.918** | 0.926* | 0.936* | 0.911** | **0.756*** | 0.767*** | 0.869** | 0.832*** | 0.808*** |
| $(B_2,\alpha=\hat{\alpha})$,K-fold | **0.917*** | **0.900*** | 0.915* | 0.931 | 0.974 | **0.728*** | 0.777*** | 0.829*** | **0.738*** | **0.713*** |
| $(B_2,\alpha=1)$,POOS-CV | **0.914*** | 0.955 | 1.057 | 1.011 | 0.883*** | 0.810*** | 0.830*** | 1.029 | 0.952 | 0.795*** |
| $(B_2,\alpha=1)$,K-fold | 0.97 | **0.901** | 0.991 | 0.983 | 0.918** | 0.837** | 0.754*** | 0.903 | 0.833*** | **0.753*** |
| $(B_2,\alpha=0)$,POOS-CV | **0.908*** | **0.893*** | 0.991 | 0.922* | 0.889*** | 0.781** | 0.769*** | 0.915 | 0.786*** | **0.788*** |
| $(B_2,\alpha=0)$,K-fold | **0.949** | **0.898*** | 0.908** | 0.906** | 0.967 | 0.875 | 0.777*** | 0.817*** | **0.756*** | **0.741*** |
| $(B_3,\alpha=\hat{\alpha})$,POOS-CV | **0.949** | **0.888*** | 0.952 | 0.943 | 0.874*** | 0.933 | 0.843*** | 0.886** | 0.829*** | 0.827*** |
| $(B_3,\alpha=\hat{\alpha})$,K-fold | **0.937** | **0.910*** | **0.882** | 0.923* | 0.921** | 0.836* | 0.831*** | 0.868*** | 0.839*** | 0.795*** |
| $(B_3,\alpha=1)$,POOS-CV | **0.929*** | **0.921** | 0.958 | 0.983 | 0.884*** | 0.812** | 0.771*** | 0.864** | 0.851** | 0.845*** |
| $(B_3,\alpha=1)$,K-fold | 0.968 | 0.941* | **0.861*** | 0.907* | 0.943 | 0.808** | 0.806*** | 0.832*** | 0.873** | **0.736*** |
| $(B_3,\alpha=0)$,POOS-CV | **0.948** | 0.974 | 0.994 | 1.066 | 0.946* | 0.979 | 1.03 | 0.956 | 0.877** | 0.799*** |
| $(B_3,\alpha=0)$,K-fold | 0.969 | **0.918*** | 0.983 | 0.998 | 0.945* | 0.963 | 0.901* | 0.957 | 0.912* | **0.730*** |
| SVR-ARDI,Lin,POOS-CV | **0.960* | 1.041 | 1.072 | 0.929 | 1.028 | 0.872 | 0.858* | 0.941 | 0.809*** | **0.779*** |
| SVR-ARDI,Lin,K-fold | **0.959* | <u>**0.873***</u> | **0.838*** | 0.926 | 0.946 | 0.801** | 0.791*** | **0.756*** | **0.800** | 0.872* |
| SVR-ARDI,RBF,POOS-CV | **0.966 | **0.995 | 1.016 | 0.957 | 0.872*** | 0.938 | 0.859* | 0.937 | 0.786*** | **0.777** |
| SVR-ARDI,RBF,K-fold | **0.943**** | 0.958 | **0.871** | 0.911* | 0.930* | **0.769*** | 0.796*** | **0.770*** | **0.763*** | **0.787*** |

*Note: The numbers represent the relative, with respect to AR,BIC model, root MSPE. Models retained in model confidence set are in bold, the minimum values are underlined, while \*\*\*, \*\*, \* stand for 1%, 5% and 10% significance of Diebold-Mariano test.*

# Table 6: Term spread: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=3 | h=9 | h=12 | h=24 | h=1 | h=3 | h=9 | h=12 | h=24 |
| Data-poor ($H_t^-$) models | | | | | | | | | | |
| AR,BIC (RMSPE) | **6.4792** | 12.8246 | **16.3575** | 20.0828 | 22.2091 | **13.3702** | 23.16 | 23.5697 | 31.597 | 23.0842 |
| AR,AIC | **1.002*** | 0.998 | **1.053*** | 1.034** | 1.041** | **1.002** | 1.001 | 1.034 | 0.993 | 0.972 |
| AR,POOS-CV | **1.055*** | 1.139* | **1.000** | 0.969 | 1.040** | **1.041** | 1.017 | 0.895* | 0.857* | 0.972 |
| AR,K-fold | **1.001** | 1.000 | **1.003** | 0.979 | 1.038* | **1.002** | 0.998 | 0.911 | 0.890* | 0.983 |
| RRAR,POOS-CV | **1.055*** | 1.142* | **1.004** | 0.998 | 1.016 | **1.036** | 1.014 | 0.899 | 0.966 | 0.945** |
| RRAR,K-fold | **1.044*** | 0.992 | **1.027** | 0.96 | 1.015 | **1.024** | 0.982 | 0.959 | 0.795** | 0.957* |
| RFAR,POOS-CV | **0.997** | 0.886 | 1.125*** | 1.019 | 1.107** | **0.906** | **0.816** | 1.039 | 0.747** | 1.077** |
| RFAR,K-fold | **0.991** | 0.941 | 1.136*** | 1.011 | 1.084** | **0.909** | 0.823 | 1.023 | 0.764* | 1.038 |
| KRR-AR,POOS-CV | 1.223** | 0.881 | <u>0.949</u> | 0.888** | 0.945* | **1.083** | 0.702 | **0.788*** | 0.758*** | 0.948 |
| KRR,AR,K-fold | 1.141 | 0.983 | **1.098** | 0.999 | 1.048 | **0.999** | 0.737 | 0.833* | **0.663** | **0.924** |
| SVR-AR,Lin,POOS-CV | 1.158** | 1.326*** | **1.071*** | 1.045 | 1.045 | 1.111* | 1.072 | 0.894* | 0.828* | 0.967 |
| SVR-AR,Lin,K-fold | 1.191** | 1.056 | **1.018** | 0.963 | 0.993 | 1.061 | 1.009 | 0.886** | 0.845** | 0.916*** |
| SVR-AR,RBF,POOS-CV | **1.006** | 1.039 | **1.050*** | 0.951 | 0.969 | **0.964** | 0.902 | 0.876* | 0.761** | **0.864*** |
| SVR-AR,RBF,K-fold | **0.985** | 0.911 | **1.038** | 0.946 | 0.933** | **0.990** | 0.737 | 0.851** | 0.747* | 0.968 |
| Data-rich ($H_t^+$) models | | | | | | | | | | |
| ARDI,BIC | **0.953** | 0.971 | **0.979** | 0.93 | **0.892*** | **0.921** | 0.9 | **0.790*** | **0.633*** | 1.049 |
| ARDI,AIC | **0.970** | 0.956 | **1.019** | 0.944 | 0.917** | **0.929** | 0.867 | **0.814*** | **0.647*** | 1.076 |
| ARDI,POOS-CV | **0.954** | 1.015 | **1.067** | 0.991 | 0.915** | **0.912** | 0.92 | 0.958 | 0.769** | 1.087 |
| ARDI,K-fold | **0.991** | 1.026 | **1.001** | 0.928 | 0.939 | **0.958** | 0.967 | 0.812*** | **0.662*** | 1.041 |
| RRARDI,POOS-CV | **0.936** | 0.994 | **1.078** | 0.991 | 0.964 | **0.896** | 0.850 | 0.952 | 0.784** | 1.092 |
| RRARDI,K-fold | **1.015** | 0.992 | **1.018** | 0.934 | 0.981 | **0.978** | 0.899 | 0.881* | **0.635*** | 1.163* |
| RFARDI,POOS-CV | **0.988** | <u>0.830*</u> | **0.957** | 0.873** | 0.921** | <u>0.804</u> | 0.691 | **0.785*** | **0.606*** | 0.985 |
| RFARDI,K-fold | **1.010** | 0.883 | **0.997** | 0.909 | 0.935** | 0.808 | 0.778 | 0.827** | **0.626*** | 0.97 |
| KRR-ARDI,POOS-CV | 1.355** | 0.898 | **0.993** | 0.856** | 0.884*** | **0.861** | <u>0.682*</u> | **0.772*** | 0.621** | **0.905*** |
| KRR,ARDI,K-fold | 1.382*** | 0.96 | **0.974** | <u>0.827**</u> | 0.862*** | **0.858** | 0.684* | <u>**0.754***</u> | **0.569*** | 0.912* |
| ($B_1, \alpha = \hat{\alpha}$),POOS-CV | 1.114 | 1.06 | 1.126*** | 1.021 | 0.866*** | **1.009** | 0.981 | 1.02 | **0.701** | 1.012 |
| ($B_1, \alpha = \hat{\alpha}$),K-fold | **1.089** | 1.149** | 1.199** | 1.106* | 0.969 | **1.001** | 1.041 | 0.885 | 0.767** | 0.941 |
| ($B_1, \alpha = 1$),POOS-CV | 1.125* | 1.115 | 1.172*** | 1.072 | <u>0.844***</u> | 1.071 | 1.006 | 1.033 | 0.833 | 0.96 |
| ($B_1, \alpha = 1$),K-fold | **1.089** | 1.149** | 1.199** | 1.106* | 0.969 | **1.001** | 1.041 | 0.885 | 0.767** | 0.941 |
| ($B_1, \alpha = 0$),POOS-CV | 1.173** | 1.312** | 1.176*** | 1.088 | 0.978 | 1.089 | 1.065 | 0.981 | 0.799 | 0.966 |
| ($B_1, \alpha = 0$),K-fold | 1.163* | 1.059 | **1.069** | 0.929 | 0.921** | **1.041** | 0.869 | **0.810** | 0.729** | **0.880*** |
| ($B_2, \alpha = \hat{\alpha}$),POOS-CV | **1.025** | 0.993 | **1.101** | 1.028 | 0.897*** | **0.918** | 0.908 | 1.02 | **0.651*** | 0.989 |
| ($B_2, \alpha = \hat{\alpha}$),K-fold | **0.976** | 0.954 | **1.098*** | 1.059 | 0.935* | **0.931** | 0.875 | 0.938 | 0.779* | 0.952 |
| ($B_2, \alpha = 1$),POOS-CV | 1.062 | 0.968 | **1.125*** | 1.049 | 0.926*** | **0.897** | 0.855 | 1.058 | 0.79 | 1.001 |
| ($B_2, \alpha = 1$),K-fold | **0.980** | 0.938 | **1.130*** | 1.01 | 0.950* | **0.948** | 0.858 | 0.976 | 0.679** | 1.001 |
| ($B_2, \alpha = 0$),POOS-CV | 1.118* | 1.082 | **1.097*** | 1.008 | 0.901*** | **1.004** | 0.919 | 1.008 | **0.669*** | 1.016 |
| ($B_2, \alpha = 0$),K-fold | 1.102 | 0.988 | **1.047** | 1.041 | 0.919** | **0.985** | 0.909 | 0.870* | 0.757* | 0.986 |
| ($B_3, \alpha = \hat{\alpha}$),POOS-CV | **0.971** | 0.964 | **1.089*** | 1.076 | 0.933* | **0.887** | 0.837 | 0.908 | 0.783* | 0.904** |
| ($B_3, \alpha = \hat{\alpha}$),K-fold | **0.968** | 0.944 | **1.009** | 0.999 | 0.898*** | **0.895** | 0.872 | 0.883** | 0.744** | **0.907*** |
| ($B_3, \alpha = 1$),POOS-CV | **1.006** | 1.066 | **1.059*** | 1.039 | 0.896*** | **0.894** | 1.131 | 0.974 | 0.764* | 0.987 |
| ($B_3, \alpha = 1$),K-fold | **0.994** | 0.924 | **1.037** | 0.96 | 0.975 | **0.934** | 0.852 | 0.834** | 0.712** | 1.01 |
| ($B_3, \alpha = 0$),POOS-CV | 1.181* | 0.961 | **1.104*** | 1.056 | 0.937** | 1.215 | 0.901 | 1.013 | 0.825 | 0.919* |
| ($B_3, \alpha = 0$),K-fold | **0.999** | 0.953 | **1.036** | 0.94 | 0.97 | **0.897** | 0.845 | 0.923 | 0.735** | 0.925** |
| SVR-ARDI,Lin,POOS-CV | 1.062 | 0.967 | **1.164*** | 1.113* | 1.065 | 1.016 | **0.762*** | 1.117 | 0.714** | 1.097 |
| SVR-ARDI,Lin,K-fold | **0.990** | 0.98 | **1.011** | 0.922 | 0.909** | **0.935** | 0.885 | 0.825** | **0.667*** | 0.994 |
| SVR-ARDI,RBF,POOS-CV | **0.972** | 0.937 | **1.069** | 1.039 | 1.068 | **0.875** | 0.741 | **0.796*** | 0.707*** | 1.204* |
| SVR-ARDI,RBF,K-fold | **1.018** | 0.938 | **1.123** | 0.914* | 0.882*** | **0.931** | 0.781 | 0.858** | 0.778** | <u>0.858**</u> |

*Note: The numbers represent the relative, with respect to AR,BIC model, root MSPE. Models retained in model confidence set are in bold, the minimum values are underlined, while ***, **, * stand for 1%, 5% and 10% significance of Diebold-Mariano test.*

# Table 7: CPI Inflation: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=3 | h=9 | h=12 | h=24 | h=1 | h=3 | h=9 | h=12 | h=24 |
| **Data-poor ($H_t^-$) models** | | | | | | | | | | |
| AR,BIC (RMSPE) | 0.0312 | 0.0257 | 0.0194 | 0.0187 | 0.0188 | 0.0556 | 0.0484 | 0.032 | 0.0277 | 0.0221 |
| AR,AIC | 0.969*** | 0.984 | 0.976* | 0.988 | 0.995 | 1.000 | 0.970** | 0.999 | 0.992 | 1.005 |
| AR,POOS-CV | 0.966** | 0.988 | 0.997 | 0.992 | 1.009 | 0.961** | 0.981 | 0.995 | 0.978 | 1.003 |
| AR,K-fold | 0.972** | 0.976** | 0.975* | 0.988 | 0.987 | 1.002 | 0.965*** | 0.998 | 0.992 | 1.005 |
| RRAR,POOS-CV | 0.969** | 0.984 | 0.99 | 0.993 | 1.006 | 0.961** | 0.982 | 0.995 | 0.963* | 0.998 |
| RRAR,K-fold | 0.964*** | 0.979** | 0.970* | 0.980* | 0.989 | 0.989 | 0.973** | 0.996 | 0.992 | 0.997 |
| RFAR,POOS-CV | 0.983 | **0.944*** | **0.909*** | **0.930** | 1.022 | 1.018 | 0.998 | 1.063 | 1.047 | 0.998 |
| RFAR,K-fold | **0.975** | **0.927*** | **0.909*** | **0.956** | 0.998 | 1.032 | 0.972 | 1.065 | 1.103 | 1.019 |
| KRR-AR,POOS-CV | **0.972** | **0.905*** | **0.872*** | **0.872*** | 0.907** | 1.023 | **0.930*** | 0.927 | 0.91 | 0.852* |
| KRR,AR,K-fold | **0.931*** | **0.888*** | **0.836*** | **0.827*** | 0.942 | 0.965 | **0.920*** | 0.92 | 0.915 | 0.975 |
| SVR-AR,Lin,POOS-CV | 1.119** | 1.291** | 1.210*** | 1.438*** | 1.417*** | 1.116 | 1.196** | 1.204** | 1.055 | 1.613*** |
| SVR-AR,Lin,K-fold | 1.239*** | 1.369** | 1.518*** | 1.606*** | 1.411*** | 1.159* | 1.326* | 1.459** | 1.501* | 1.016 |
| SVR-AR,RBF,POOS-CV | 0.988 | 1.004 | 1.086* | 1.068** | 1.127** | 0.999 | 1.004 | 0.969 | 1.091** | 1.501*** |
| SVR-AR,RBF,K-fold | 0.99 | 1.025 | 1.025 | 1.003 | 1.370*** | 0.965 | 0.979 | 0.996 | 0.896** | 1.553** |
| **Data-rich ($H_t^+$) models** | | | | | | | | | | |
| ARDI,BIC | 0.96 | **0.973** | 1.024 | **0.895*** | 0.880* | 0.919* | **0.906*** | 0.779* | 0.755** | 0.713** |
| ARDI,AIC | **0.954** | **0.990** | 1.034 | **0.895** | 0.884 | 0.925 | **0.898** | 0.778* | **0.736*** | 0.676** |
| ARDI,POOS-CV | **0.950** | **0.984** | 1.017 | **0.910** | 0.916 | 0.916* | **0.913*** | 0.832** | 0.781*** | **0.669*** |
| ARDI,K-fold | **0.941*** | **0.990** | 1.028 | **0.873*** | 0.858* | 0.891** | **0.900** | 0.784* | **0.709*** | 0.635** |
| RRARDI,POOS-CV | **0.943*** | **0.975** | 1.001 | **0.917** | 0.914 | 0.905* | **0.912*** | 0.828** | **0.780*** | 0.666** |
| RRARDI,K-fold | **0.943*** | **0.983** | 1.022 | **0.875*** | 0.882 | 0.927* | **0.901** | **0.744*** | **0.664*** | **0.613*** |
| RFARDI,POOS-CV | **0.947*** | **0.908*** | **0.853*** | **0.914*** | 0.979 | 0.976 | **0.939*** | 0.988 | 1.051 | 0.964 |
| RFARDI,K-fold | **0.936*** | **0.907*** | **0.854*** | **0.868*** | 0.909* | 0.962 | **0.933*** | 0.979 | 0.93 | 1.003 |
| KRR-ARDI,POOS-CV | 1.006 | 1.043 | 0.959 | 0.972 | 1.067 | 1.046 | 1.093 | 0.952 | 0.948 | 0.946 |
| KRR,ARDI,K-fold | **0.985** | 0.999 | 0.983 | 0.977 | 0.938 | 0.998 | 0.99 | 1.023 | 1.022 | 0.986 |
| $(B_1,\alpha=\hat{\alpha})$,POOS-CV | **0.918*** | **0.916*** | 0.976 | 0.96 | 1.026 | <u>**0.803***</u> | **0.900*** | 0.8 | 0.848 | 0.974 |
| $(B_1,\alpha=\hat{\alpha})$,K-fold | <u>**0.908***</u> | **0.921*** | 1.012 | 1.056 | 1.092* | **0.823*** | <u>**0.873***</u> | **0.774** | 0.836 | 1.069 |
| $(B_1,\alpha=1)$,POOS-CV | **0.960** | **0.908*** | 1.11 | 1.03 | 1.076 | **0.813*** | **0.889*** | 0.794 | 0.825 | 0.989 |
| $(B_1,\alpha=1)$,K-fold | **0.908*** | **0.921*** | 1.012 | 1.056 | 1.092* | **0.823*** | **0.873*** | **0.774** | 0.836 | 1.069 |
| $(B_1,\alpha=0)$,POOS-CV | 0.971 | 1.035 | 1.114* | 1.048 | 1.263** | 0.848** | **0.906** | 0.935 | 0.881 | 0.99 |
| $(B_1,\alpha=0)$,K-fold | **0.945*** | 1.057 | 1.246** | 1.289** | 1.260*** | **0.850*** | **0.939** | 0.954 | 0.944 | 1.095 |
| $(B_2,\alpha=\hat{\alpha})$,POOS-CV | **0.923*** | **0.956*** | **0.940** | 0.934 | 0.945 | 0.871** | 0.959 | 0.803* | 0.802* | 0.822* |
| $(B_2,\alpha=\hat{\alpha})$,K-fold | **0.921*** | **0.963*** | 0.995 | 0.956 | 1.037 | 0.868* | 0.957* | 0.817* | 0.778** | 0.861 |
| $(B_2,\alpha=1)$,POOS-CV | **0.942** | **0.959** | 1.158* | 1.174** | 1.151** | 0.877 | 0.927 | 0.799 | 0.907 | 1.087 |
| $(B_2,\alpha=1)$,K-fold | **0.922*** | **0.970** | 1.066 | 0.995 | 1.168* | 0.879 | 0.929 | 0.853 | 0.816* | 1.009 |
| $(B_2,\alpha=0)$,POOS-CV | **0.921*** | **0.940** | 1.079 | 0.959 | 1.071 | 0.857* | **0.881** | 1.129 | 0.883 | 0.851 |
| $(B_2,\alpha=0)$,K-fold | **0.919*** | **0.929*** | 0.997 | 1.011 | 1.212** | 0.865* | **0.883** | 0.825 | 0.961 | 0.853 |
| $(B_3,\alpha=\hat{\alpha})$,POOS-CV | **0.935*** | **0.941*** | **0.961** | **0.849*** | 0.901* | 0.889* | **0.947*** | 0.791** | 0.785** | 0.808** |
| $(B_3,\alpha=\hat{\alpha})$,K-fold | **0.938*** | **0.952*** | **0.937** | **0.915** | 0.952 | 0.891* | 0.958 | 0.801* | 0.784** | 0.91 |
| $(B_3,\alpha=1)$,POOS-CV | **0.933*** | **0.960** | 1.076 | 1.000 | 1.017 | **0.856*** | **0.917*** | **0.755*** | **0.769*** | 0.86 |
| $(B_3,\alpha=1)$,K-fold | **0.943** | 0.978 | 1.006 | **0.894** | 1.002 | 0.889 | 0.946 | 0.805 | 0.806* | 0.879 |
| $(B_3,\alpha=0)$,POOS-CV | **0.946*** | **0.939*** | **0.896*** | **0.871*** | 1.022 | 0.894* | **0.931*** | 0.865 | 0.875 | 0.896 |
| $(B_3,\alpha=0)$,K-fold | **0.921*** | **0.975** | **0.926** | **0.920** | 1.106 | 0.877*** | **0.936** | 0.839 | 0.892 | 1.147 |
| SVR-ARDI,Lin,POOS-CV | 1.148*** | 1.202* | 1.251*** | 1.209*** | 1.219** | 1.068 | 1.053 | 0.969 | 0.969 | 0.943 |
| SVR-ARDI,Lin,K-fold | 1.115*** | 1.390** | 1.197** | 1.114 | 1.177* | 1.058 | 1.295* | 0.944 | 0.954 | 1.036 |
| SVR-ARDI,RBF,POOS-CV | 0.963 | 1.031 | 1.002 | 0.962 | 0.951 | 0.922 | **0.915** | 0.848 | 0.861 | 0.996 |
| SVR-ARDI,RBF,K-fold | **0.951*** | 1.002 | 0.997 | 0.945 | <u>**0.797***</u> | 0.927* | 0.964 | 0.816** | 0.826** | **0.659*** |

*Note: The numbers represent the relative, with respect to AR,BIC model, root MSPE. Models retained in model confidence set are in bold, the minimum values are underlined, while \*\*\*, \*\*, \* stand for 1%, 5% and 10% significance of Diebold-Mariano test.*

## Table 8: Housing starts: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=3 | h=9 | h=12 | h=24 | h=1 | h=3 | h=9 | h=12 | h=24 |
| Data-poor ($H_t^-$) models | | | | | | | | | | |
| AR,BIC (RMSPE) | **0.9040** | **0.4142** | **0.2499** | 0.2198 | 0.1671 | **1.2526** | 0.6658 | 0.4897 | 0.4158 | 0.2954 |
| AR,AIC | **0.998** | **1.019** | **1.000** | **1.000** | 1.000 | 1.01 | 0.965* | 1.000 | 1.000 | 1.000 |
| AR,POOS-CV | **1.001** | **1.012** | **1.019*** | 1.01 | 1.036** | 1.015 | <u>**0.936**</u>** | 1.011* | 1.013 | 1.057** |
| AR,K-fold | **0.993** | **1.017** | **1.001** | 1.000 | 1.02 | 1.01 | **0.951**** | 1.000 | 1.000 | 1.036 |
| RRAR,POOS-CV | **1.007** | **1.007** | **1.008** | 1.009 | 1.031** | 1.027* | **0.939**** | 1.001 | 1.013 | 1.050** |
| RRAR,K-fold | **0.999** | **1.014** | **0.998** | **0.998** | 1.024* | 1.013 | **0.941**** | **1.000**** | 0.999 | 1.042** |
| RFAR,POOS-CV | 1.030*** | **1.026*** | **1.028*** | 1.045** | 1.018 | 1.023 | **0.941*** | **0.992** | 1.048* | 1.013 |
| RFAR,K-fold | 1.017* | **1.022** | **1.007** | 1.031** | 1.008 | 1.02 | **0.942*** | **0.990** | 1.026 | 1.01 |
| KRR-AR,POOS-CV | **0.995** | **0.999** | **0.969*** | 1.044* | 1.037* | **0.990** | 0.972 | 0.971 | 1.050** | 0.993 |
| KRR-AR,K-fold | **0.977*** | **0.975** | **0.957**** | 0.989 | 1.001 | **0.985** | 0.976 | 1.01 | 1.006 | 1.004 |
| SVR-AR,Lin,POOS-CV | 1.032*** | **0.997** | **1.044*** | 1.064*** | 1.223** | 1.024* | **0.962*** | **0.986*** | 0.984 | **0.957**** |
| SVR-AR,Lin,K-fold | 1.036*** | 1.031 | **1.002** | 1.006 | 1.002 | 1.013 | 0.976 | 1.002 | 1.009 | 1.004 |
| SVR-AR,RBF,POOS-CV | **1.008** | 1.047** | **1.023** | 1.035*** | 1.060*** | **1.014** | 0.981 | **0.947**** | 1.015 | 1.017 |
| SVR-AR,RBF,K-fold | **1.009** | **1.011** | **1.012**** | 1.020*** | 1.034** | 1.021* | 0.969* | **1.010**** | 1.017** | 1.001 |
| Data-rich ($H_t^+$) models | | | | | | | | | | |
| ARDI,BIC | <u>**0.973**</u>* | **0.989** | **1.031** | 1.051 | 1.05 | **0.946** | 1.139 | 1.048 | 0.988 | 0.944 |
| ARDI,AIC | **0.992** | **0.995** | **1.018** | 1.06 | 1.078 | **1.000** | 1.113 | 1.025 | 1.025 | 0.96 |
| ARDI,POOS-CV | 1.01 | **1.007** | **1.080** | 1.027 | 0.998 | 1.023 | 1.128 | 1.054 | 1.015 | 1.021 |
| ARDI,K-fold | **0.992** | **0.984** | **1.026** | 1.061 | 1.094 | **1.011** | 1.093 | 1.027 | 1.027 | 0.958 |
| RRARDI,POOS-CV | **0.998** | **1.007** | **1.043** | **0.996** | 1.082 | 1.008 | 1.119 | 1.041 | 0.991 | 1.022 |
| RRARDI,K-fold | **0.998** | **0.988** | **1.051** | 1.064 | 1.089 | 1.017 | 1.118 | 1.033 | 0.998 | 0.941 |
| RFARDI,POOS-CV | **0.997** | **0.944**** | **0.930**** | **0.920*** | 0.899** | **0.982** | 0.971 | 0.965 | 0.957 | 0.972 |
| RFARDI,K-fold | **0.994** | **0.962** | **0.939*** | <u>**0.914***</u> | **0.838*** | **0.993** | 0.985 | **0.986** | 0.943 | 0.902* |
| KRR-ARDI,POOS-CV | **0.980** | <u>**0.943***</u> | <u>**0.915***</u> | **0.942*** | **0.884*** | <u>**0.941***</u> | 0.952* | 0.949 | 0.964** | 0.986 |
| KRR,ARDI,K-fold | **0.982**** | **0.949**** | **0.928** | **0.933** | **0.889**** | **0.973** | 0.973 | **1.003** | 1.022 | 0.994 |
| ($B_1.\alpha = \hat{\alpha}$),POOS-CV | **1.006** | **1.000** | **1.063** | 1.016 | **0.895*** | 1.023 | 1.099 | **0.985** | 1.026 | 1.022 |
| ($B_1.\alpha = \hat{\alpha}$),K-fold | 1.040* | 1.095** | 1.250** | 1.335** | 1.151* | 1.096* | 1.152** | 1.021 | 1.127 | **0.890** |
| ($B_1.\alpha = 1$),POOS-CV | 1.032** | 1.039 | 1.155 | 1.045 | 0.949 | 1.013 | 1.063 | **0.961** | 1.025 | 1.062 |
| ($B_1.\alpha = 1$),K-fold | 1.040* | 1.095** | 1.250** | 1.335** | 1.151* | 1.096* | 1.152** | 1.021 | 1.127 | **0.890** |
| ($B_1.\alpha = 0$),POOS-CV | **0.982** | **0.977** | 1.084 | 1.337** | 0.959 | **0.999** | 1.017 | 1.014 | 1.152** | 0.964 |
| ($B_1.\alpha = 0$),K-fold | **0.982** | **1.006** | 1.137* | 1.158** | 1.007 | **0.994** | 1.03 | **1.017** | 1.067 | <u>**0.809***</u> |
| ($B_2.\alpha = \hat{\alpha}$),POOS-CV | 1.044 | **0.992** | **0.975** | 0.988 | 0.969 | 1.177 | 1.126* | 1.034 | 0.989 | 0.972 |
| ($B_2.\alpha = \hat{\alpha}$),K-fold | **0.988** | **1.003** | **1.069** | 1.193** | 1.069 | 1.11 | 1.188* | 1.085 | 1.133* | 0.917 |
| ($B_2.\alpha = 1$),POOS-CV | **1.001** | **1.000** | **0.967** | 1.02 | **0.940*** | **0.961** | 1.047 | **0.943** | 0.985 | 1.006 |
| ($B_2.\alpha = 1$),K-fold | **0.989** | 1.095 | 1.245** | 1.203* | 1.093 | 1.007 | 1.322*** | 1.1 | **0.919** | **0.848**** |
| ($B_2.\alpha = 0$),POOS-CV | 1.091* | **0.949** | **0.987** | **0.971** | 0.939 | 1.255 | 1.027 | **0.992** | **0.956** | 0.994 |
| ($B_2.\alpha = 0$),K-fold | 1.066 | **1.068** | 1.19 | 1.044 | 1.064 | 1.248 | 1.332** | 1.057 | <u>**0.896***</u> | 0.917 |
| ($B_3, \alpha = \hat{\alpha}$),POOS-CV | 1.009 | **0.951*** | **0.935** | 0.99 | **0.891*** | 1.028 | 1.019 | **0.958** | **0.963** | 0.987 |
| ($B_3, \alpha = \hat{\alpha}$),K-fold | **0.998** | **0.977** | **1.007** | 1.055 | 1.044 | 1.019 | 1.115 | 1.017 | **0.979** | **0.882*** |
| ($B_3, \alpha = 1$),POOS-CV | **0.997** | **0.975** | **1.024** | 0.996 | 0.928* | **0.976** | 1.001 | 1.021 | **0.940** | 1.001 |
| ($B_3, \alpha = 1$),K-fold | 1.013 | 1.040 | 1.071 | 1.106 | 1.145 | 1.042 | 1.219* | 1.036 | 0.992 | 1.009 |
| ($B_3, \alpha = 0$),POOS-CV | 1.022* | **0.951*** | **0.962** | **0.944** | 0.932* | 1.022 | 0.981 | <u>**0.930**</u> | **0.915**** | 1.001 |
| ($B_3, \alpha = 0$),K-fold | 1.030** | **1.003** | **1.005** | 1.011 | 1.029 | **0.986** | 1.114 | **0.998** | 0.955 | 0.934 |
| SVR-ARDI,Lin,POOS-CV | **0.998** | 1.078* | 1.154* | 1.137* | 1.142 | 1.047 | 1.111 | **0.989** | 1.009 | 1.111 |
| SVR-ARDI,Lin,K-fold | **0.992** | **0.971** | **1.017** | 1.038 | 1.11 | 1.007 | 1.021 | **0.988** | **0.937** | 0.959 |
| SVR-ARDI,RBF,POOS-CV | **0.991** | **1.004** | **1.010** | 1.044 | 1.034 | **0.987** | 1.095 | **0.981** | 0.969 | 1.096 |
| SVR-ARDI,RBF,K-fold | **1.003** | **0.998** | **1.045** | 1.078 | 1.162* | 1.022 | 1.081 | 1.03 | 0.984 | 1.026 |

*Note: The numbers represent the relative, with respect to AR,BIC model, root MSPE. Models retained in model confidence set are in bold, the minimum values are underlined, while \*\*\*, \*\*, \* stand for 1%, 5% and 10% significance of Diebold-Mariano test.*

# B Robustness of Treatment Effects Graphs



**Figure 11:** This figure plots the distribution of $\dot\alpha_F^{(h,v)}$ from equation 11 done by $(h,v)$ subsets. The subsample under consideration here is **data-poor models**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. Variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. SEs are HAC. These are the 95% confidence bands.
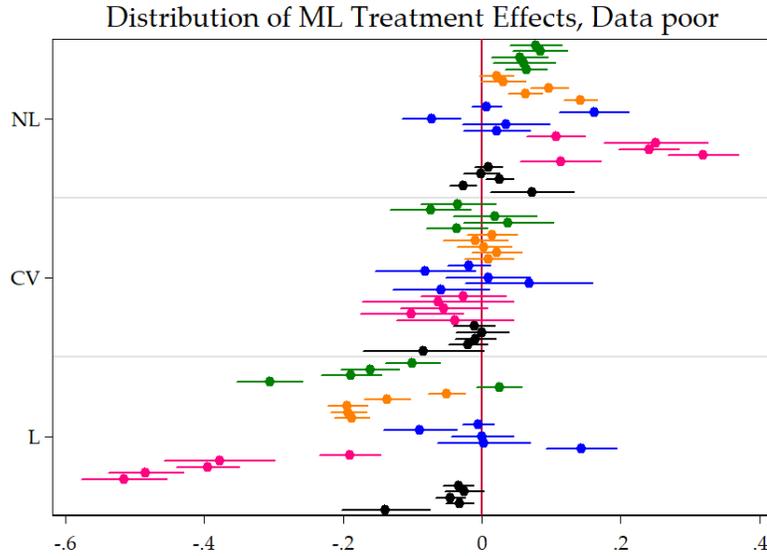


**Figure 12:** This figure plots the distribution of $\dot\alpha_F^{(h,v)}$ from equation 11 done by $(h,v)$ subsets. The subsample under consideration here is **data-rich models**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. Variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. SEs are HAC. These are the 95% confidence bands.

**Figure 13:** This figure plots the distribution of $\hat{\alpha}_F^{(h,v)}$ from equation 11 done by $(h,v)$ subsets. The subsample under consideration here are **recessions**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. Variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. SEs are HAC. These are the 95% confidence bands.
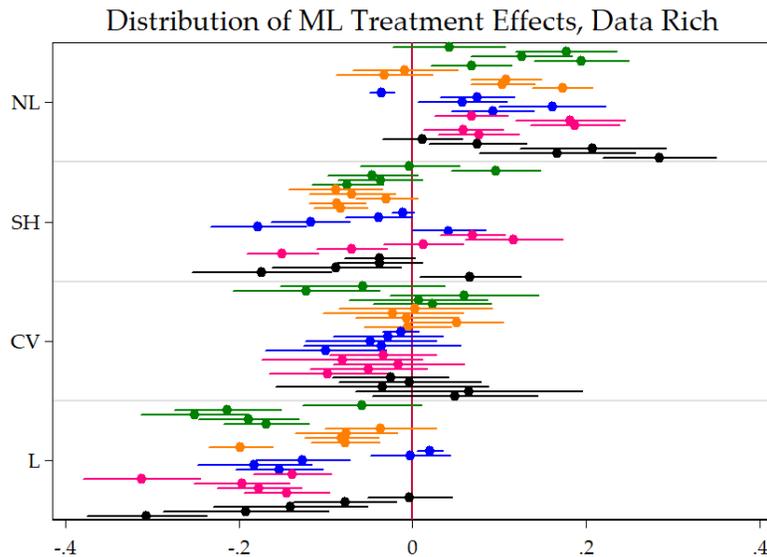


**Figure 14:** This figure plots the distribution of $\hat{\alpha}_F^{(h,v)}$ from equation 11 done by $(h,v)$ subsets. The subsample under consideration here are **expansions**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. Variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. SEs are HAC. These are the 95% confidence bands.
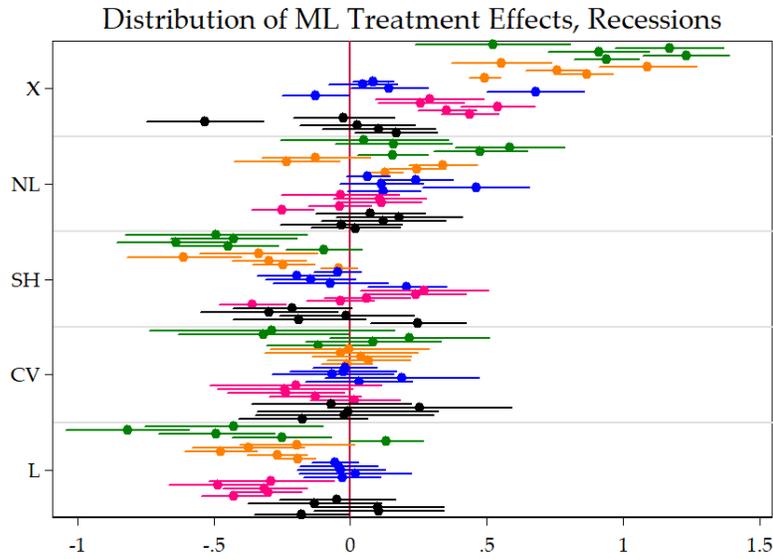
# C  Additional Results

**Figure 15:** This figure plots the distribution of $\dot{\alpha}_F^{(h,v)}$ from equation 11 done by $(h, v)$ subsets. The subsample under consideration here are **the last 20 years**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. Variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. SEs are HAC. These are the 95% confidence bands.
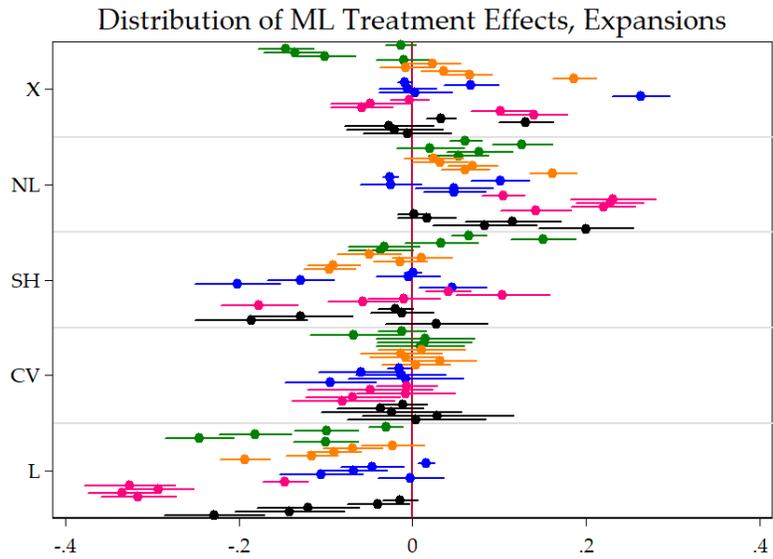


**Figure 16:** This figure plots time series of variables explaining the heterogeneity of NL treatment effects in section 6.

**Figure 17:** This figure shows the 3-year rolling window root MSPE, the cumulative root MSPE and Giacomini and Rossi (2010) fluctuation tests for linear and nonlinear data-poor and data-rich models, at 12-month horizon.

**Figure 18:** This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in comparing the data-poor and data-rich environments for linear models. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.



**Figure 19:** This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in comparing the data-poor and data-rich environments for nonlinear models. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

# D   Nonlinearites Matter – A Robustness Check

In this appendix, we trade Random Forests for Boosted Trees and KRR for Neural Networks. First, we briefly introduce the newest addition to our nonlinear arsenal. Second, we demonstrate that very similar conclusions to that of section 5.2.1 are reached using those. This further backs our claim that nonlinearities matter, whichever way they were obtained.

## D.1   Data-Poor

**Boosted Trees AR (BTAR)**. This algorithm provides an alternative means of approximating nonlinear functions by additively combining regression trees in a sequential fashion. Let $\eta \in [0,1]$ be the learning rate and $\hat{y}_{t+h}^{(n)}$ and $e_{t+h}^{(n)} := y_{t-h} - \eta \hat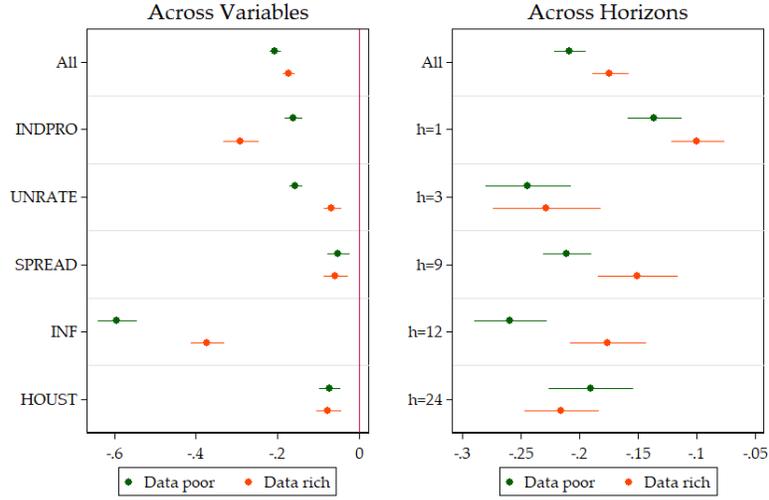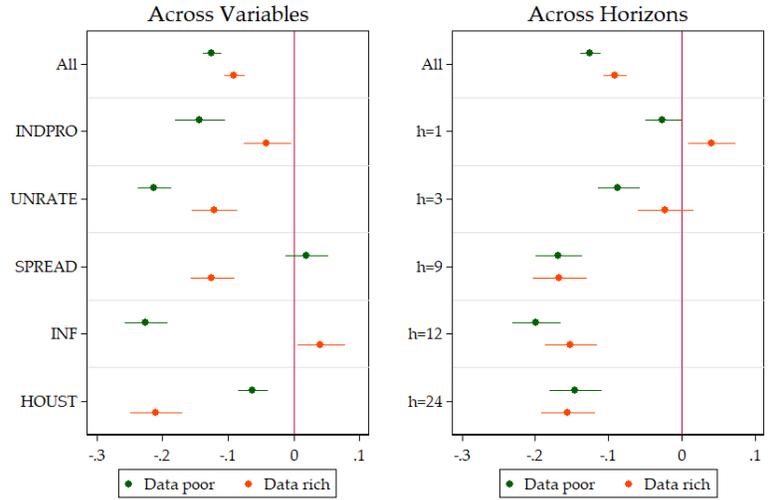{y}_{t+h}^{(n)}$ be the step $n$ predicted value and pseudo-residuals, respectively. Then, the step $n+1$ prediction is obtained as

$$\hat{y}_{t+h}^{(n+1)} = y_{t+h}^{(n)} + \rho_{n+1} f(Z_t, c_{n+1})$$

where $(c_{n+1}, \rho_{n+1}) := \underset{\rho, c}{\arg\min} \sum_{t=1}^{T} \left( e_{t+h}^{(n)} - \rho_{n+1} f(Z_t, c_{n+1}) \right)^2$ and $c_{n+1} := (c_{n+1,m})_{m=1}^{M}$ are the parameters of a regression tree. In other words, it recursively fits trees on pseudo-residuals. The maximum depth of each tree is set to 10 and all features are considered at each split. We select the number of steps and $\eta \in [0,1]$ with Bayesian optimization. We impose $p_y = 12$.

**Neural Network AR (NNAR)**. We opted for fully connected feed-forward neural networks. The value of the input vector $[Z_{it}]_{i=1}^{N_0}$ is represented by a layer of input neurons, each taking on the value of a different element in the vector. Each neuron $j$ of the first hidden layer takes on a value $h_{jt}^{(n)}$ which is determined by applying a potentially nonlinear transformation to a weighted sum of the input value. The same is true of each subsequent hidden layer until we have reached the output layer which contains a single neuron whose value is the $h$ period

ahead forecast of the model. Formally, our neural network models have the following form:

$$
h_{jt}^{(n)} = \begin{cases} f^{(1)} \left( \sum_{i=1}^{N_0} w_{ji}^{(1)} Z_{it} + w_{j0}^{(1)} \right) & n = 1 \\ f^{(n)} \left( \sum_{i=1}^{N_k} w_{ji}^{(n)} h_{it}^{(n-1)} + w_{j0}^{(n)} \right) & n > 1 \end{cases}
$$

$$
\hat{y}_{t+h} = \sum_{i=1}^{N_{N_h}} w_i^{(y)} h_{jt}^{(N_h)} + w_0^{(y)}.
$$

We restrict our attention to two fixed architectures: the first one uses a single hidden layer of 32 neurons $((N_h, N_1) = (1, 32))$ and the second one uses two hidden layers of 32 and 16 neurons, respectively $((N_h, N_1, N_2) = (2, 32, 16))$. In all cases, we use rectified linear units (ReLU) as the activation functions, i.e.

$$
f^{(n)}(z) = \max\{0, z\}, \forall n = 1, ..., N_h.
$$

The training is carried out by batch gradient descent using the Adam algorithm. This algorithm is initialized with a learning rate of 0.01 and we use an early stopping rule[28]. And, in an effort to mitigate the effects of overfitting and the impact of random initialization of weights, we train 5 neural networks with the same architecture and use their average output as our prediction value. In essence, those neural networks are simplified versions of the neural networks used in Gu et al. (2020b) where we got rid of the hyperparameter optimization and use 5 base learners instead of 10. For this algorithm, the input is a set of $p_y = 12$ lagged values of the target variable. We do not make use of cross-validation, but we do estimate model weights recursively.

## D.2 Data-Rich

**Boosted Trees ARDI (BTARDI)**. We consider a vanilla Boosted Trees where the maximum depth of each tree is set to 10 and all features are considered at each split. We select the number of steps and $\eta \in [0, 1]$ with Bayesian optimization. We impose $p_y = 12$, $p_f = 12$ and

---

[28]If improvements in the performance metric doesn't exceed a tolerance threshold for 5 consecutive epochs, we stop the training.

**Figure 20:** This figure compares the two alternative NL models averaged over all horizons. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

$K = 8$.

**Neural Network ARDI (NNARDI)**. We opted for fully connected feed-forward neural network with the same architecture as the data-poor version, but we now use $(p_y, p_f, K) = (12, 10, 12)$ for the inputs.

## D.3   Results

In line with what reported in section 5.2.1, we find that NL's treatment effect is magnified for horizons 9, 12 and 24. Additionally, it is found that both algorithms give very homogeneous improvements in the data-rich environment, another finding detailed in the main text. Results for the data-poor environment are more scattered, as they were before. Targets benefiting most from NL in the data-rich environment are INF and HOUST, which is analogous to earlier findings. However, it was found that the real activity targets benefited more from NL in our main text configuration, which is the sole noticeable difference with results reported here.

Figure 21: This figure compares the two alternative NL models averaged over all variables. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

# How is Machine Learning Useful for
# Macroeconomic Forecasting?[*]
# SUPPLEMENTARY MATERIAL

Philippe Goulet Coulombe[1][†]    Maxime Leroux[2]    Dalibor Stevanovic[2][‡]
Stéphane Surprenant[2]

[1]University of Pennsylvania
[2]Université du Québec à Montréal

This version: August 21, 2020

**Abstract**

This document contains supplementary material for the paper entitled *How Is Machine Learning Useful for Macroeconomic Forecasting?* It contains the following appendices: results for absolute loss; results with quarterly US data; results with monthly Canadian data; description of CV techniques and technical details on forecasting models.

*JEL Classification: C53, C55, E37*
*Keywords: Machine Learning, Big Data, Forecasting.*

[†]Corresponding Author: gouletc@sas.upenn.edu. Department of Economics, UPenn.
[‡]Corresponding Author: dstevanovic.econ@gmail.com. Département des sciences économiques, UQAM.

# A    Results with Absolute Loss

In this section we present results for a different out-of-sample loss function that is often used in the literature: the absolute loss. Following Koenker and Machado (1999), we generate the pseudo-$R^1$ in order to perform regressions (11) and (12): $R^1_{t,h,v,m} \equiv 1 - \frac{|e_{t,h,v,m}|}{\frac{1}{T}\sum_{t=1}^{T}|y_{v,t+h}-\bar{y}_{v,h}|}$. Hence, the figure included in this section are exact replication of those included in the main text except that the target variable of all the regressions has been changed.

The main message here is that results obtained using the squared loss are very consistent with what one would obtain using the absolute loss. The importance of each feature, figure 22, and the way it behaves according to the variable/horizon pair is the same. Indeed, most of the heterogeneity is variable specific while there are clear horizon patterns emerging when we average out variables. For instance, we clearly see by comparing figures 24 and 2 that more data and nonlinearities usefulness increase linearly in $h$. CV is flat around the 0 line. Alternative shrinkage and loss function both are negative and follow a boomerang shape (they are not as bad for short and very long horizons, but quite bad in between).

The pertinence of nonlinearities and the impertinence of alternative shrinkage follow very similar behavior to what is obtained in the main body of this paper. However, for nonlinearities, the data-poor advantages are not robust to the choice of MSPE vs MAPE. Fortunately, besides that, the figures are all very much alike.

Results for the alternative in-sample loss function also seem to be independent of the proposed choices of out-of-sample loss function. Only for hyperparameters selection we do get slightly different results: CV-KF is now sometimes worse than BIC in a statistically significant way. However, the negative effect is again much stronger for POOS CV. CV-KF still outperforms any other model selection criteria on recessions.

**Figure 22:** This figure presents predictive importance estimates. Random forest is trained to predict $R^1_{t,h,v,m}$ defined in (11) and use out-of-bags observations to assess the performance of the model and compute features' importance. NL, SH, CV and LF stand for nonlinearity, shrinkage, cross-validation and loss function features respectively. A dummy for $H^+_t$ models, $X$, is included as well.



**Figure 23:** This figure plots the distribution of $\hat{\alpha}^{(h,v)}_F$ from equation (11) done by $(h,v)$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^1$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. Finally, variables are INDPRO, UNRATE, SPREAD, INF and HOUST. Within a specific color block, the horizon increases from $h = 1$ to $h = 24$ as we are going down. As an example, we clearly see that the partial effect of $X$ on the $R^1$ of INF increases drastically with the forecasted horizon $h$. SEs are HAC. These are the 95% confidence bands.

3

**Figure 24:** This figure plots the distribution of $\dot{\alpha}_F^{(v)}$ and $\dot{\alpha}_F^{(h)}$ from equation (11) done by $h$ and $v$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^1$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. However, in this graph, $v-$specific heterogeneity and $h-$specific heterogeneity have been integrated out in turns. SEs are HAC. These are the 95% confidence bands.



**Figure 25:** This compares the two NL models averaged over all horizons. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.

## Contribution of Non-Linearities, by horizons



**Figure 26:** This compares the two NL models averaged over all variables. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.



**Figure 27:** This compares models of section 3.2 averaged over all variables and horizons. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. The base models are ARDIs specified with POOS-CV and KF-CV respectively. SEs are HAC. These are the 95% confidence bands.

Table 9: CV comparison

|  | (1) All | (2) Data-rich | (3) Data-poor | (4) Data-rich | (5) Data-poor |
|---|---|---|---|---|---|
| CV-KF | 0.0114 | -0.0233 | 0.0461 | -0.221 | -0.109 |
|  | (0.375) | (0.340) | (0.181) | (0.364) | (0.193) |
| CV-POOS | -0.765* | -0.762* | -0.768*** | -0.700 | -0.859*** |
|  | (0.375) | (0.340) | (0.181) | (0.364) | (0.193) |
| AIC | -0.396 | -0.516 | -0.275 | -0.507 | -0.522** |
|  | (0.375) | (0.340) | (0.181) | (0.364) | (0.193) |
| CV-KF * Recessions |  |  |  | 1.609 | 1.264* |
|  |  |  |  | (1.037) | (0.552) |
| CV-POOS * Recessions |  |  |  | -0.506 | 0.747 |
|  |  |  |  | (1.037) | (0.552) |
| AIC * Recessions |  |  |  | -0.0760 | 2.007*** |
|  |  |  |  | (1.037) | (0.552) |
| Observations | 91200 | 45600 | 45600 | 45600 | 45600 |

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

### CV-KF performance relative to CV-POOS



Figure 28: This compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.

**Figure 29:** This compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.



**Figure 30:** This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in **both the data-poor and data-rich environments**. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.

**Non-Linear SVR Relative Performance to KRR**

Figure 31: This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in **both recession and expansion periods**. The unit of the x-axis are improvements in OOS $R^1$ over the basis model. SEs are HAC. These are the 95% confidence bands.

# B   Results with Quarterly Data

In this section we present results for quarterly frequency using the dataset FRED-QD, publicly available at the Federal Reserve of St-Louis's web site. This is the quarterly companion to FRED-MD monthly dataset used in the main part of paper. It contains 248 US macroeconomic and financial aggregates o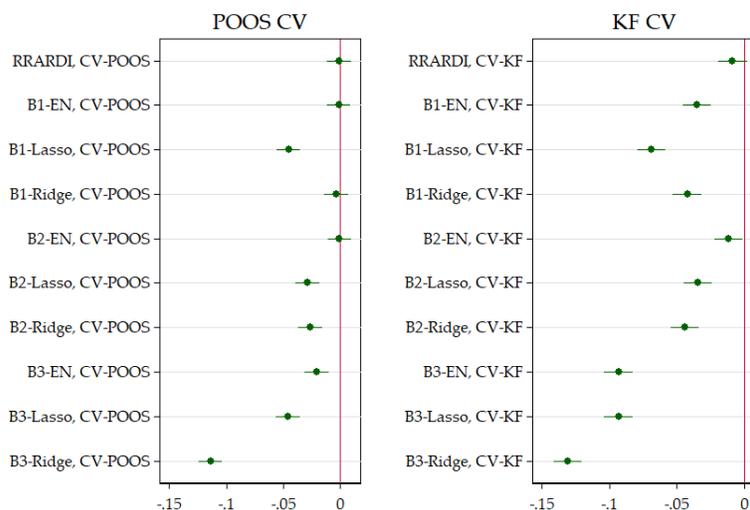bserved from 1960Q1 to 2018Q4. The series transformations to induce stationarity are the same as in Stock and Watson (2012a). The variables of interest are: real GDP, real personal consumption expenditures (CONS), real gross private investment (INV), real disposable personal income (INC) and the PCE deflator. All the targets are expressed in average growth rate over $h$ periods as in equation (4). Forecasting horizons are 1, 2, 3, 4 and 8 quarters.

The main message here is that results obtained using the quarterly data and predicting GDP components are consistent with those on monthly variables. Tables 10 - 14 summarize the overall predictive ability in terms of RMPSE relative to the reference AR,BIC model. GDP and consumption growths are best predicted at short run by the standard Stock and Watson (2002a) ARDI,BIC model, while random forests dominate at longer horizons. Nonlinear models perform well for most horizons when predicting the disposable income growth. Finally, kernel ridge regressions (both data-poor and data-rich) are the best options to predict the PCE inflation.

The ML features' importance is plotted in figure 32. Contrary to monthly data, horizons and variables fixed effects are much less important which is somehow expected because of relative smoothness of quarterly data and similar targets (4 out 5 are real activity series). Among ML treatments, shrinkage is the most important, followed by loss function and non-linearity. As in the monthly application, CV is the least relevant, while the data-rich component remains very important. From figures 33 and 34, we see that: (i) the richness of predictors' set is very helpful for most of the targets; (ii) nonlinearity treatment has positive and significant effects for investment, income and PCE deflator, while it is not significant for GDP and CONS; (iii) the impertinence of alternative shrinkage follow very similar behavior to what is obtained in the main body of this paper; (iv) CV has in general negative but small and often insignificant effect; (v) SVR loss function decreases the predictive performance as in the monthly case, especially for income growth and inflation.

## Table 10: GDP: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=2 | h=3 | h=4 | h=8 | h=1 | h=2 | h=3 | h=4 | h=8 |
| **Data-poor ($H_t^-$) models** | | | | | | | | | | |
| AR,BIC (RMSPE) | **0.0752** | **0.0656** | **0.0619** | **0.0593** | **0.0521** | 0,1199 | 0,1347 | 0,1261 | 0,1285 | 0,1022 |
| AR,AIC | **1.004** | **0.994** | **0.999** | **1.000** | **1.000** | 1,034 | 0,995 | 1 | 1 | 1 |
| AR,POOS-CV | **0.984**** | **0.994** | **0.994** | **1.000** | **1.017** | **0.991** | 0,994 | 0,993 | 1 | 1,033 |
| AR,K-fold | **0.998** | **1.003** | **0.999** | 1,001 | **1.000** | 1,026 | 1,01 | 0,997 | 1,001 | 1 |
| RRAR,POOS-CV | **0.992** | **1.002** | **1.000** | 1,005 | **1.005** | 1,014 | 1 | 1,005 | 0,997 | 1,014 |
| RRAR,K-fold | **1.013** | **1.007** | **1.006** | 1,012 | **1.000** | 1.092* | 1.010* | 1.020*** | 1,02 | 0.999*** |
| RFAR,POOS-CV | 1.185*** | 1.104*** | 1.165*** | 1.129*** | 1.061** | 1.241** | 1.077* | 1.116** | 1.070** | 0.925*** |
| RFAR,K-fold | 1.082** | 1.124*** | 1.105** | 1.121*** | **1.064**** | 1.124* | 1,085 | 1,021 | 1.089** | 0,989 |
| KRR-AR,POOS-CV | 1,049 | 1,044 | 1,011 | 1.065* | **0.993** | 1.103** | 0,954 | 0.913* | 0.943* | 0.873*** |
| KRR,AR,K-fold | 1,044 | **1.033** | 1.051** | 1,013 | **0.995** | 1.172*** | 1,01 | 1,036 | 0,974 | 0.963*** |
| SVR-AR,Lin,POOS-CV | 1.161** | 1.136** | 1.129** | 1.143** | 1.045 | 1.233*** | 1.106** | 1.152*** | 1.061** | 1,071 |
| SVR-AR,Lin,K-fold | 1.082** | 1.092** | **1.054*** | 1.051** | **0.986** | 1.222*** | 1.110** | 1.088** | 1.054** | 0.964*** |
| SVR-AR,RBF,POOS-CV | 1,015 | **1.036*** | **1.026** | 1,051 | 1.095** | 1.038** | 1,01 | 1.037* | 0,991 | 1,016 |
| SVR-AR,RBF,K-fold | 1.043** | 1.032* | 1.029* | 1,018 | **1.011*** | 1.157*** | 1.032** | 1.041** | 0,986 | 1,002 |
| **Data-rich ($H_t^+$) models** | | | | | | | | | | |
| ARDI,BIC | <u>**0.884**</u> | <u>**0.811**</u>** | <u>**0.824**</u>** | **0.817**** | **1.002** | <u>**0.829**</u> | **0.649**** | **0.732**** | **0.704**** | 0.714*** |
| ARDI,AIC | **0.905** | **0.833*** | **0.844*** | **0.832*** | **0.989** | **0.931** | **0.652**** | **0.741**** | **0.721**** | 0.687*** |
| ARDI,POOS-CV | **0.913** | **0.861*** | **0.878** | **0.885** | **0.918** | **0.936** | **0.689**** | **0.742**** | **0.719**** | 0.735*** |
| ARDI,K-fold | **0.978** | **0.881** | **0.871** | **0.815*** | **1.070** | 1,078 | 0.709** | 0.767** | <u>**0.681**</u>*** | **0.595**** |
| RRARDI,POOS-CV | **0.938** | **0.853*** | **0.846*** | **0.924** | **0.949** | 1,034 | 0.717*** | **0.742**** | 0.740*** | 0.770*** |
| RRARDI,K-fold | **0.906** | **0.839*** | **0.842*** | <u>**0.810**</u>* | **1.021** | **0.924** | 0.720** | **0.755**** | **0.690**** | **0.587**** |
| RFARDI,POOS-CV | **0.938** | **0.929** | **0.876*** | **0.866*** | **0.887*** | 0,989 | 0.866* | 0.810** | 0.761*** | 0.739*** |
| RFARDI,K-fold | **0.941** | **0.908*** | **0.868*** | **0.856*** | <u>**0.862**</u>** | 1,022 | 0.843** | 0.813* | 0.742*** | 0.692*** |
| KRR-ARDI,POOS-CV | 1,055 | **1.048** | 1.074** | 1,049 | **1.011** | 1.135* | 0,97 | 0,979 | 0.923* | 0.921* |
| KRR,ARDI,K-fold | 1,005 | 1,038 | 1,065 | 1,074 | **0.957** | 1 | 0,969 | 0,947 | 0,95 | 0.822*** |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | 1.061* | **1.057** | **1.039** | 1.077** | **1.026** | 1.118** | 0,977 | 1,057 | 0,981 | 0.931** |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | 1,015 | **0.964** | **1.016** | 1.079** | **1.010** | 1,041 | 0,955 | 0,98 | 0,972 | 0.907*** |
| $(B_1, \alpha = 1)$,POOS-CV | 1.076** | 1.104* | **1.008** | 1.065* | **1.006** | 1.179*** | 1,007 | 1,003 | 0,954 | 0.937* |
| $(B_1, \alpha = 1)$,K-fold | **0.994** | **1.018** | 1,033 | 1.079* | **0.971** | **0.989** | 0,989 | 1,013 | 0.947* | 0.890*** |
| $(B_1, \alpha = 0)$,POOS-CV | 1.082* | **1.064** | 1.148*** | 1.145* | **0.992** | 1.242*** | 1,083 | 1.156*** | 1,033 | 0,979 |
| $(B_1, \alpha = 0)$,K-fold | 1.191** | 1.079* | 1,052 | 1.070* | **0.968** | 1.091** | 0,974 | 0,999 | 1,011 | 0.928* |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | 1,043 | **1.022** | **1.021** | 1,032 | **1.015** | 1.083* | 1,01 | 1,007 | 0,907 | 0.900** |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | **0.991** | **1.007** | **0.994** | **0.980** | **1.126** | 1,077 | 1,002 | 0,947 | **0.747**** | 0.612*** |
| $(B_2, \alpha = 1)$,POOS-CV | 1.110** | 1.072* | **1.007** | **0.991** | **0.918** | 1.217** | 1.090* | 0,998 | 0,924 | 0.782*** |
| $(B_2, \alpha = 1)$,K-fold | 1,039 | **1.027** | **1.003** | **0.961** | **1.069** | 1.136** | 1,029 | 0,957 | 0.777*** | <u>**0.563**</u>*** |
| $(B_2, \alpha = 0)$,POOS-CV | **1.000** | **1.000** | **1.001** | 0,989 | **0.978** | 1,106 | 0,959 | 0,976 | 0.852** | 0.772*** |
| $(B_2, \alpha = 0)$,K-fold | **0.986** | **0.980** | **0.980** | 1,001 | 1,132 | 1,073 | 0,958 | 0,968 | 0.819** | 0.750*** |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | 1,047 | 1,055 | 1.049* | 1,052 | **1.003** | 1,046 | 1,027 | 1.043* | 1,037 | 0.930*** |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | 1,038 | **0.975** | **1.004** | 1,021 | **0.991** | 1,056 | 0,98 | 0,988 | 0.918*** | 0.839*** |
| $(B_3, \alpha = 1)$,POOS-CV | 1.055* | 1.133** | **1.044** | 1.107** | **0.995** | 1,058 | 1.116* | 1,033 | 1,067 | 0.895** |
| $(B_3, \alpha = 1)$,K-fold | 1,045 | **1.020** | **1.009** | 1,021 | **0.982** | 1,078 | 0,994 | 1,011 | 0.942* | 0.854*** |
| $(B_3, \alpha = 0)$,POOS-CV | 1.142** | 1.153* | **0.979** | 1.217* | **0.992** | 1.124** | 1,046 | 0,976 | 1,162 | 0,973 |
| $(B_3, \alpha = 0)$,K-fold | 1.225* | **1.105** | **0.994** | 1,139 | 1.068* | 1.197** | 1,021 | 0,987 | 1,098 | 0,979 |
| SVR-ARDI,Lin,POOS-CV | **1.014** | **1.088** | 1.130* | **0.966** | **1.073** | 0,972 | 0,984 | 1,016 | 0.806*** | 0.933* |
| SVR-ARDI,Lin,K-fold | 1,027 | **1.112** | 1,064 | 1,084 | 1.237** | **0.982** | 0,998 | 0,876 | 0,957 | 0.863*** |
| SVR-ARDI,RBF,POOS-CV | 1,033 | **1.015** | **0.924** | 1,013 | **1.034** | 1,201 | 1,001 | **0.779**** | 0.871* | 0.861** |
| SVR-ARDI,RBF,K-fold | **0.896** | **0.887** | **0.930** | 0,973 | 1,089 | **0.930** | 0.781** | 0.807* | 0.823** | 0.813*** |

*Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while \*\*\*. \*\*. \* stand for 1%. 5% and 10% significance of Diebold-Mariano test.*

## Table 11: Consumption: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=2 | h=3 | h=4 | h=8 | h=1 | h=2 | h=3 | h=4 | h=8 |
| Data-poor ($H_t^-$) models | | | | | | | | | | |
| AR,BIC (RMSPE) | 0,0604 | **0.0485** | 0,0451 | 0,0476 | **0.0480** | 0,0927 | 0,0848 | 0,0851 | 0,0947 | 0,0881 |
| AR,AIC | 0.982** | **0.993** | 1,001 | 0.979** | **1.000** | 0.961*** | 0,993 | 1,004 | 0.978* | 1 |
| AR,POOS-CV | 0.961** | **0.986**** | **0.998** | **0.974**** | **0.997** | 0.920* | 0,995 | 0,999 | 0.971** | 0,998 |
| AR,K-fold | 0.987* | **1.025** | 1,015 | 0.975** | **1.035** | 0.977*** | 1,026 | 1,014 | 0.974** | 1,062 |
| RRAR,POOS-CV | **0.944**** | **0.988*** | 1 | **0.968**** | **0.998** | 0.878** | 0,989 | 1 | 0.971* | 0,99 |
| RRAR,K-fold | 0.973** | **1.013** | 1.015** | 1 | **1.011*** | 0,947 | 1,013 | 1.017* | 1.015** | 1,014 |
| RFAR,POOS-CV | 0,989 | **1.036** | 1,02 | 1,01 | 1.065** | 0,977 | 0,987 | 0.929* | 0,965 | 1,035 |
| RFAR,K-fold | 1,015 | **1.008** | 1.044* | 1.052* | 1.067** | 0,951 | 0,897 | 0,959 | 1,002 | 0,979 |
| KRR-AR,POOS-CV | 0,986 | **0.995** | 1.072* | 1.064** | **1.010** | 0,994 | 0,946 | 0,953 | 0,973 | 0,951 |
| KRR,AR,K-fold | 1,012 | **0.980** | 1,031 | 1,003 | **0.994** | 1,017 | 0,924 | 0,943 | 0,95 | 0.946** |
| SVR-AR,Lin,POOS-CV | 1,013 | 1.339*** | 1.304*** | 1.166*** | 1.012 | 0,868 | 1.225* | 1.350*** | 1.150*** | 0.935* |
| SVR-AR,Lin,K-fold | 1,085 | 1.176** | 1.222*** | 1.117*** | **1.020*** | 1,101 | 1.234* | 1.251*** | 1.133*** | 0,989 |
| SVR-AR,RBF,POOS-CV | 1.081* | **1.098**** | 1.120*** | 1.052** | **1.005** | 1,06 | 1,07 | 1,003 | 0.937*** | 0.934* |
| SVR-AR,RBF,K-fold | 0,973 | **1.026** | 1.064*** | **0.956**** | 1.083** | 0.881* | 1 | 1.054* | 0.959** | 1.109** |
| Data-rich ($H_t^+$) models | | | | | | | | | | |
| ARDI,BIC | <u>**0.897***</u> | <u>**0.879**</u> | <u>**0.903**</u> | **0.938** | **1.017** | <u>**0.782***</u> | <u>**0.729****</u> | **0.782**** | **0.829**** | 0.809*** |
| ARDI,AIC | **0.916** | **0.939** | 0,983 | 0,988 | **1.094** | 0,857 | **0.752*** | 0.800* | **0.830*** | 0.761*** |
| ARDI,POOS-CV | 1,007 | **1.002** | 1,06 | 1,069 | **0.967** | 1,071 | 0,948 | 1,05 | 1,02 | 0.860* |
| ARDI,K-fold | 1,092 | **0.948** | **0.967** | **0.959** | 1,116 | 1,31 | **0.768*** | **0.764**** | **0.819**** | 0.769*** |
| RRARDI,POOS-CV | 1,009 | **1.005** | 1,018 | 1,018 | **1.049** | 1,151 | 0,965 | 1,023 | 0,976 | 0.802** |
| RRARDI,K-fold | 1,083 | **0.924** | **0.977** | 0,995 | **1.071** | 1,339 | **0.752**** | 0,889 | 0,853 | 0.682*** |
| RFARDI,POOS-CV | 0,976 | **0.946** | **0.969** | **0.928** | **0.982** | 0,895 | **0.853*** | **0.840**** | **0.781***** | 0.808*** |
| RFARDI,K-fold | **0.937*** | **0.961** | **0.979** | <u>**0.913**</u> | **0.957** | 0.872** | **0.785**** | **0.810**** | <u>**0.775****</u> | 0.757*** |
| KRR-ARDI,POOS-CV | 1.138** | **1.112*** | 1.181** | 1.141*** | **1.021** | 1,123 | 1,059 | 1,117 | 1,028 | 0,919 |
| KRR,ARDI,K-fold | 1,054 | **1.058** | 1.118** | 1,065 | **0.994** | 1,035 | 0,909 | 0,972 | 0,955 | 0.849** |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | 1.153*** | 1.213*** | 1.168** | 1.107** | **1.038** | 1,134 | 1.238** | 1.191* | 1,009 | 0,926 |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | 1,069 | 1.193*** | 1.186*** | 1.120** | 1.079* | 1,103 | 1,155 | 1.212*** | 1.151*** | 0.901* |
| $(B_1, \alpha = 1)$,POOS-CV | 1.118** | 1.215*** | 1.184** | 1.153*** | **1.054** | 1,135 | 1,178 | 1.194* | 1,086 | 0,954 |
| $(B_1, \alpha = 1)$,K-fold | 1,056 | 1.166*** | 1.122** | 1.079** | **1.016** | 1,048 | 1,151 | 1,078 | 1.117*** | 0.878** |
| $(B_1, \alpha = 0)$,POOS-CV | 1.158*** | 1.281*** | 1.300*** | 1.171** | 1.062** | 1,119 | 1,163 | 1,172 | 1,049 | 1,012 |
| $(B_1, \alpha = 0)$,K-fold | 1.453*** | 1.219** | 1.288* | 1.103** | 1.039 | 1,325 | 0,947 | 1,069 | 1.072** | 0,966 |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | 1.092* | **1.107*** | 1.140* | 1.105* | **1.082** | 0,98 | 1,143 | 1,14 | 0,997 | 0.826** |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | 1,036 | **1.088**** | 1.167** | 1,082 | 1,129 | 1.080** | 1.139** | 1,119 | **0.814**** | <u>**0.628****</u> |
| $(B_2, \alpha = 1)$,POOS-CV | 1.158** | **1.136*** | 1.194** | 1.187*** | **1.027** | 1,051 | 1,188 | 1.223** | 1,005 | 0.839** |
| $(B_2, \alpha = 1)$,K-fold | 1,057 | 1.179*** | 1.113* | 1,072 | 1,153 | 1,107 | 1.263*** | 1,056 | 0.872* | 0.672*** |
| $(B_2, \alpha = 0)$,POOS-CV | 1.054* | **1.081*** | 1.194** | 1,049 | **1.079** | 1.084* | 1,1 | 1,056 | 0,883 | 0.865** |
| $(B_2, \alpha = 0)$,K-fold | 1.072* | **1.088** | 1.133* | 1,083 | 1.255* | 1.133** | 1,135 | 1,13 | 0.853* | 0.791*** |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | 1,061 | 1.128** | 1.165** | 1,055 | 1.052** | 1,05 | 1.164* | 1.183** | 1,027 | 1,003 |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | 1.128** | **1.057** | 1.149** | 1.125*** | **1.005** | 1,091 | 1,049 | 1.093* | 1,023 | 0.764*** |
| $(B_3, \alpha = 1)$,POOS-CV | 1.096* | 1.174** | 1.186** | 1.138** | 1.079*** | 1,095 | 1.202* | 1.192* | 1,05 | 1,006 |
| $(B_3, \alpha = 1)$,K-fold | 1,065 | **1.106**** | 1.153** | 1.188** | 1,129 | 1,052 | 1,107 | 1,149 | 1,04 | 0.825** |
| $(B_3, \alpha = 0)$,POOS-CV | 1,063 | **1.100*** | 1.118*** | 1.168** | **1.015** | 1,012 | 1,14 | 1.144** | 1.166* | 1,001 |
| $(B_3, \alpha = 0)$,K-fold | 1.441** | **1.188*** | 1.144*** | 1.152** | **1.049*** | 1.584** | 1,085 | 1.122*** | 1,104 | 0,986 |
| SVR-ARDI,Lin,POOS-CV | 1,046 | **1.201*** | 1,108 | 1,064 | 1.106* | 0,989 | 1,119 | 1,069 | 1,004 | 1,007 |
| SVR-ARDI,Lin,K-fold | 1,105 | **1.010** | 1.265** | 1,038 | 1,088 | 1,285 | 1,032 | 1,093 | 0,925 | 0.776*** |
| SVR-ARDI,RBF,POOS-CV | 1,053 | **1.021** | 1,118 | 1.080* | 1,441 | 1,077 | 1,043 | 1,069 | 0,999 | 1,754 |
| SVR-ARDI,RBF,K-fold | 0,986 | **0.987** | 1,058 | **0.981** | **1.016** | 0,932 | 0,873 | **0.755**** | **0.830*** | 0.679*** |

*Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while \*\*\*. \*\*. \* stand for 1%. 5% and 10% significance of Diebold-Mariano test.*

## Table 12: Investment: Relative Root MSPE

| | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Models | h=1 | h=2 | h=3 | h=4 | h=8 | h=1 | h=2 | h=3 | h=4 | h=8 |
| **Data-poor ($H_t^-$) models** | | | | | | | | | | |
| AR,BIC (RMSPE) | 0,4078 | 0,3385 | 0,2986 | 0,277 | 0,2036 | 0,7551 | 0,6866 | 0,5725 | 0,5482 | 0,3834 |
| AR,AIC | 1.015* | 1.011* | 1.007* | 1 | 0,996 | 1.023** | 1.015** | 1.010* | 1 | 0,991 |
| AR,POOS-CV | 0.995* | 1,004 | 1.007** | 1,004 | 1,007 | 1 | 1.008* | 1.006** | 1,008 | 1,03 |
| AR,K-fold | 1,007 | 1,004 | 1,009 | 1 | 1,021 | 1,002 | 1.018** | 1.024*** | 1,017 | 1.040* |
| RRAR,POOS-CV | 1,004 | 1,001 | 1.013*** | 1.007** | 1,001 | 1,01 | 1,002 | 1.016*** | 1.007* | 1,006 |
| RRAR,K-fold | 1.015** | 1.013* | 1.008* | 1 | 1,002 | 1.026*** | 1,012 | 1.016*** | 1.013*** | 0,998 |
| RFAR,POOS-CV | 1,055 | 1,013 | 0,979 | 0,985 | 1,046 | 1,024 | 0.905** | **0.880***** | 0,978 | 1,022 |
| RFAR,K-fold | 1,036 | 1,016 | 1,019 | 1 | 0,977 | 0,992 | 0,942 | 1,007 | 0,934 | 0,957 |
| KRR-AR,POOS-CV | 1,036 | 1 | **0.979** | 1,001 | 0,953 | 1.079* | 0,937 | 0,989 | 1,003 | 0.947** |
| KRR-AR,K-fold | 0,996 | 1,008 | 0.961* | 1 | 0.969** | 1,022 | 0,987 | 0,975 | 1,015 | 0.965*** |
| SVR-AR,Lin,POOS-CV | 1,033 | 1.097** | 1.096*** | 1.050* | 1.116** | 1,035 | 1,061 | 1.041** | 1 | 0,98 |
| SVR-AR,Lin,K-fold | 1.033* | 1.033* | 1.026** | 1.016* | 1,019 | 1.063** | 1,021 | 1.028* | 0,998 | 1,004 |
| SVR-AR,RBF,POOS-CV | 1.038*** | 1,13 | 1.062*** | 1.047** | 1.094*** | 1.050** | 1,145 | 1.069** | 1.008** | 1,006 |
| SVR-AR,RBF,K-fold | 1,03 | 1,026 | 1.039** | 1,01 | 0,986 | 1.066* | 1,018 | 1.040** | 0,994 | 0,995 |
| **Data-rich ($H_t^+$) models** | | | | | | | | | | |
| ARDI,BIC | **0.749***** | **0.774**** | **0.862*** | **0.827**** | **0.911*** | **0.603***** | **0.665***** | 0,851 | **0.827***** | 0,949 |
| ARDI,AIC | **0.757***** | **0.894*** | **0.933** | **0.831*** | 0,948 | **0.601***** | 0,847 | 0,936 | **0.773**** | 0.849** |
| ARDI,POOS-CV | **0.745***** | **0.801**** | **0.918** | 0,913 | 0,979 | **0.623***** | **0.736**** | **0.939** | **0.809***** | 0,924 |
| ARDI,K-fold | **0.765***** | **0.905** | **0.944** | **0.854** | 1,009 | **0.584***** | 0,837 | 0,993 | **0.784**** | 0.811*** |
| RRARDI,POOS-CV | **0.776***** | **0.858**** | **0.916** | 0,984 | 0,976 | **0.626***** | 0,831 | **0.937** | 0,945 | 0,969 |
| RRARDI,K-fold | **0.742***** | **0.866*** | **0.912** | **0.925** | 0,985 | **0.603***** | **0.810*** | **0.931** | 0,923 | 0.828*** |
| RFARDI,POOS-CV | 0.907** | **0.910**** | **0.884**** | **0.833**** | **0.814**** | 0,917 | 0,898 | 0,885 | **0.790***** | 0.750*** |
| RFARDI,K-fold | 0,951 | **0.927*** | **0.875**** | **0.830**** | **0.806**** | 0,966 | 0,92 | 0,922 | **0.830**** | 0.735*** |
| KRR-ARDI,POOS-CV | 0,989 | **0.945** | **0.966** | 0,942 | 0.919* | 1,01 | 0,95 | 1,028 | 0,959 | 0,933 |
| KRR,ARDI,K-fold | 0,978 | **0.952** | 0,995 | 0.937* | 0.930* | 0,974 | 0.932* | 1,049 | 0,983 | 0,987 |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | 1,036 | 0,976 | 1,014 | 1,007 | 0,939 | 0.884** | 0.916*** | 1,006 | 0.925* | 0,965 |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | 1,046 | 0,967 | **0.939** | 0.915* | 1,012 | 1.076* | 0,964 | 0,951 | 0.894*** | 0,993 |
| $(B_1, \alpha = 1)$,POOS-CV | 1,023 | 0,991 | 0,989 | 0,941 | 0,966 | 0.889* | 0,954 | 0,974 | 0.902* | 0,973 |
| $(B_1, \alpha = 1)$,K-fold | 0,953 | **0.914*** | **0.918*** | **0.887**** | 1,018 | 0.905* | 0.941** | 0,959 | 0.899*** | 0,953 |
| $(B_1, \alpha = 0)$,POOS-CV | 1,019 | 0,997 | 1.110** | 1,045 | 1,013 | 0,973 | 0,997 | 1.078*** | 1.071* | 1,008 |
| $(B_1, \alpha = 0)$,K-fold | 1.117** | 0,98 | **0.977** | 0,971 | 0,93 | 1,012 | 0.931** | 0.897 | 0,914 | 0,912 |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | 0,996 | 0,973 | 1,01 | 1,016 | **0.915** | 1,038 | 0,974 | 1,047 | 0,989 | 0.848** |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | 0,974 | 0,975 | 0,958 | 1,005 | 0,956 | 1,026 | 0,965 | 0,94 | 0.886** | 0.662** |
| $(B_2, \alpha = 1)$,POOS-CV | 0,988 | 0,961 | 1,076 | 1,069 | 1,003 | 1,008 | 0.959* | 1.150** | 1,067 | 0.874*** |
| $(B_2, \alpha = 1)$,K-fold | 0,974 | 0,965 | 0,967 | 1,014 | **0.794**** | 0,997 | 0,973 | 0,975 | 0.854* | **0.615***** |
| $(B_2, \alpha = 0)$,POOS-CV | 1,033 | 0,975 | 1,048 | 1,057 | **0.904*** | 1,056 | 0,991 | 1,102 | 1,031 | 0.871** |
| $(B_2, \alpha = 0)$,K-fold | 1,023 | 0,923 | 0,966 | 0,996 | 0,966 | 1,025 | 0.892** | 0,993 | 0,946 | 0,894 |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | 0,961 | 0,982 | 1,006 | 0,988 | **0.920*** | 0.901* | 0,991 | 1,058 | 0,996 | 0.929*** |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | 0.948* | 0,976 | **0.921** | **0.884**** | 0,941 | 0.928* | 0.967* | **0.913** | 0.845** | 0.888*** |
| $(B_3, \alpha = 1)$,POOS-CV | 0,946 | 0,985 | **0.957** | 0,977 | 0.939* | 0,916 | 0,993 | 1,037 | 0,975 | 0.941** |
| $(B_3, \alpha = 1)$,K-fold | 0,956 | 0,966 | **0.891**** | **0.894**** | 0,954 | 0,937 | 0,973 | **0.894**** | 0.881*** | 0.880*** |
| $(B_3, \alpha = 0)$,POOS-CV | 1.110* | 1.036* | 1,027 | 1,027 | 1 | 1,011 | 0,97 | 1,004 | 1,011 | 1,001 |
| $(B_3, \alpha = 0)$,K-fold | 1,151 | 0,989 | 0,982 | 1,136 | 1,023 | 0,99 | 0,965 | 0,974 | 1,089 | 0,968 |
| SVR-ARDI,Lin,POOS-CV | 0,975 | 0,995 | 1,077 | 1,013 | 1,013 | 1,042 | 0,974 | 1,086 | 0,986 | 0,938 |
| SVR-ARDI,Lin,K-fold | **0.758***** | **0.805**** | **0.908** | 1,094 | 1,098 | **0.623***** | **0.739***** | **0.808*** | 0,975 | 0,964 |
| SVR-ARDI,RBF,POOS-CV | **0.791***** | **0.909** | **0.969** | 0,956 | 0,948 | **0.711***** | 0,856 | **0.876** | 0,934 | 0.904** |
| SVR-ARDI,RBF,K-fold | **0.804***** | **0.836*** | **0.913** | 0,962 | 0,979 | 0.737*** | **0.728**** | **0.852** | 0,965 | 0.812** |

*Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while \*\*\*. \*\*. \* stand for 1%. 5% and 10% significance of Diebold-Mariano test.*

## Table 13: Income: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=2 | h=3 | h=4 | h=8 | h=1 | h=2 | h=3 | h=4 | h=8 |
| Data-poor ($H_t^-$) models | | | | | | | | | | |
| AR,BIC (RMSPE) | **0.1011** | **0.0669** | **0.0581** | **0.0528** | 0,0417 | 0,1336 | 0,088 | 0,0803 | 0,0772 | 0,0683 |
| AR,AIC | **0.995** | **0.991** | **0.998** | **1.000** | 1 | 1 | 0.969* | 1 | 1 | 1 |
| AR,POOS-CV | **0.985*** | **0.996** | **1.002** | **0.999** | **0.991** | **0.938**** | **0.980**** | 0.992* | 0,998 | 0,993 |
| AR,K-fold | **0.987** | **0.992** | **0.994** | **0.998** | 1,002 | **0.947**** | **0.963**** | **0.969**** | 1 | 0,999 |
| RRAR,POOS-CV | **0.987** | **0.996** | **1.002** | 1.006*** | 0,995 | **0.939**** | **0.976**** | 0,994 | 1.006*** | 0.991** |
| RRAR,K-fold | **0.988** | **0.991** | **1.000** | **1.003*** | 1 | **0.945**** | 0.972*** | 1 | 1.008*** | 0.999** |
| RFAR,POOS-CV | **1.028** | 1.068** | 1.075** | 1,016 | 1,008 | 1,072 | 1.103* | 0.939* | 0,975 | 0,975 |
| RFAR,K-fold | 1.132*** | **1.024** | 1.056* | 1,01 | 1,036 | 1.124** | 0,976 | 0,989 | 0,985 | 0,957 |
| KRR-AR,POOS-CV | **0.990** | **1.000** | 1,033 | 1.070** | **0.967** | **0.923**** | **0.905**** | 0,959 | 0,979 | 0.908* |
| KRR,AR,K-fold | **0.988** | **0.991** | **1.004** | 1.049* | 1,037 | **0.964** | **0.897**** | 0,956 | 0,978 | 0.913** |
| SVR-AR,Lin,POOS-CV | **1.000** | 1,056 | **1.009** | 1,881 | 1.165** | 0,976 | 0,954 | 0,97 | 0,993 | 1,111 |
| SVR-AR,Lin,K-fold | **0.993** | **0.995** | **0.996** | **0.988** | **0.962*** | 0,976 | 0,996 | 1,015 | 1.016** | 0.965*** |
| SVR-AR,RBF,POOS-CV | **0.975** | 1,049 | 1,022 | 1.066* | **0.969** | **0.939**** | 0,959 | 0,973 | 1,01 | 0.928*** |
| SVR-AR,RBF,K-fold | 1.012* | **0.996** | **1.009** | 1,012 | 1.018* | 1,01 | 1 | 1.026* | 1.036*** | 1.029** |
| Data-rich ($H_t^+$) models | | | | | | | | | | |
| ARDI,BIC | **1.059** | **0.981** | **0.913**** | **0.939** | **0.963** | 1,257 | **0.773**** | **<u>0.726</u>*** | **0.777**** | 0.769*** |
| ARDI,AIC | **1.016** | **0.940** | **<u>0.911</u>*** | **0.966** | **0.992** | 1,05 | **<u>0.611</u>*** | **0.757**** | 0,886 | 0.721*** |
| ARDI,POOS-CV | **1.040** | **0.975** | **0.945** | **0.933** | 1,128 | 1,149 | **0.757**** | **0.753**** | **0.757**** | 0.770** |
| ARDI,K-fold | 1,065 | **0.946** | **0.953** | **0.974** | 1,028 | 1,175 | **0.664**** | **0.796**** | 0,898 | 0.689*** |
| RRARDI,POOS-CV | **1.038** | **1.007** | **0.971** | **<u>0.917</u>** | 1,058 | 1,12 | **0.796*** | 0,869 | **0.767**** | 0.743*** |
| RRARDI,K-fold | 1,06 | **0.973** | **0.925** | **0.919** | **0.999** | 1,197 | 0,82 | 0.830* | 0,871 | **0.627*** |
| RFARDI,POOS-CV | **<u>0.954</u>*** | **<u>0.932</u>*** | **0.936*** | **0.919*** | **0.910*** | **0.916** | **0.807**** | **0.822**** | **0.762**** | 0.678*** |
| RFARDI,K-fold | **0.977** | **0.957** | **0.929**** | **0.925**** | **<u>0.886</u>*** | **0.931** | 0.821** | 0.802** | 0.795*** | **0.675*** |
| KRR-ARDI,POOS-CV | 1,026 | 1.069*** | 1,025 | 1.090* | 0,985 | **0.948** | 0,991 | 0.936** | 0,954 | 0.894** |
| KRR,ARDI,K-fold | **0.969** | 1,012 | 1,075 | 1.084* | 0,991 | **0.947** | 0.925** | 0,942 | 0.929* | 0.849*** |
| $(B_1, \alpha = \hat{\alpha})$,POOS-CV | **1.010** | 1.045* | **0.997** | 1,016 | 1,015 | **0.948*** | 0,993 | 1,018 | 1,034 | 0.922* |
| $(B_1, \alpha = \hat{\alpha})$,K-fold | **1.008** | 1,02 | **1.031** | 1,025 | 1,055 | 0,988 | 1,063 | 0.882*** | 0,972 | 0,903 |
| $(B_1, \alpha = 1)$,POOS-CV | **1.010** | 1.105** | 1.070* | 1.035* | 1,016 | 0,998 | 0,963 | 0,985 | 1.067** | 0.914** |
| $(B_1, \alpha = 1)$,K-fold | 1,017 | **1.020** | 1,014 | 1,015 | 1,091 | 1,036 | 1,066 | 0,974 | 0,958 | 0.895* |
| $(B_1, \alpha = 0)$,POOS-CV | 1.030* | 1,034 | 1.050** | 1.075*** | 1,014 | **0.942*** | 1,021 | 1,034 | 1,031 | 1.120* |
| $(B_1, \alpha = 0)$,K-fold | 1.023* | **0.996** | 1,032 | 1,010 | 0,953 | 0,972 | 0.921** | 0.904*** | 0,964 | 0,936 |
| $(B_2, \alpha = \hat{\alpha})$,POOS-CV | **1.001** | **0.976** | **0.989** | 1,027 | **0.972** | 0,994 | 0.874** | 0,998 | 1.043** | 0.772** |
| $(B_2, \alpha = \hat{\alpha})$,K-fold | **1.020** | **0.979** | **0.975** | **0.988** | 1.220** | 1.054* | 0,934 | 0,931 | 0,897 | 0.790** |
| $(B_2, \alpha = 1)$,POOS-CV | **0.992** | **0.988** | **0.991** | **1.005** | **0.947** | 0,978 | 1,003 | 0,991 | 1,002 | 0.877*** |
| $(B_2, \alpha = 1)$,K-fold | 1.080* | **0.971** | **0.958** | **0.966** | 1.262** | 1.253* | 0.872** | 0.848** | 0.838** | **0.691**** |
| $(B_2, \alpha = 0)$,POOS-CV | **1.022** | **0.978** | **0.958** | **0.993** | **0.964** | 1,061 | 0.844*** | 0,924 | 0,931 | 0.722*** |
| $(B_2, \alpha = 0)$,K-fold | 1,028 | **1.000** | **0.990** | **0.997** | 1,158 | 1,051 | 0,955 | 0,983 | 0,921 | 0.830** |
| $(B_3, \alpha = \hat{\alpha})$,POOS-CV | **1.009** | **1.010** | 1,013 | 1,032 | 1,015 | 0.953* | 0,993 | 1.047** | 1,027 | 0.935** |
| $(B_3, \alpha = \hat{\alpha})$,K-fold | **0.990** | **0.995** | **0.997** | 1,024 | 1.085* | 0,962 | 0,924 | 0,969 | 1.051* | 0.882*** |
| $(B_3, \alpha = 1)$,POOS-CV | **0.995** | **1.005** | **1.006** | 1,035 | 1.040** | 0,978 | 0,984 | 1.056** | 1,047 | 0.991* |
| $(B_3, \alpha = 1)$,K-fold | **1.003** | **1.006** | **1.005** | **0.999** | 1.171*** | 1,001 | 0.931* | 0,999 | 1,002 | 0.862*** |
| $(B_3, \alpha = 0)$,POOS-CV | **0.985** | **0.987** | **0.986** | 1,04 | **0.984** | **0.941**** | 0,954 | 0,987 | 1,145 | 0.959** |
| $(B_3, \alpha = 0)$,K-fold | **0.993** | 1,132 | **1.000** | 1,078 | 1.166** | **0.947**** | 0.906** | 0,991 | 1,134 | 1,001 |
| SVR-ARDI,Lin,POOS-CV | 1,06 | 1,081 | **1.005** | **0.982** | 1,082 | **0.958** | 1,019 | 0,906 | 0.863* | 0.888** |
| SVR-ARDI,Lin,K-fold | 1.170* | **0.968** | 1,042 | **0.984** | 1,144 | 1.512* | 0,852 | 0.821* | **<u>0.736</u>*** | 0,988 |
| SVR-ARDI,RBF,POOS-CV | 1.147** | 1,097 | **0.975** | **0.972** | 1,025 | 1.311* | 1,069 | 0,97 | 0,992 | 0,931 |
| SVR-ARDI,RBF,K-fold | **1.008** | 1,117 | **0.985** | **0.998** | 1,191 | **0.943** | 1,286 | 0.827** | 0.843** | 0.770*** |

*Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while \*\*\*. \*\*. \* stand for 1%. 5% and 10% significance of Diebold-Mariano test.*

## Table 14: PCE Deflator: Relative Root MSPE

| Models | Full Out-of-Sample | | | | | NBER Recessions Periods | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | h=1 | h=2 | h=3 | h=4 | h=8 | h=1 | h=2 | h=3 | h=4 | h=8 |
| **Data-poor ($H_t^-$) models** | | | | | | | | | | |
| AR,BIC (RMSPE) | **0.0442** | 0,0421 | 0,0395 | **0.0387** | **0.0418** | **0.0798** | 0,0827 | 0,078 | 0,069 | 0,0644 |
| AR,AIC | **1.000** | 0,999 | 0.992** | **0.991**** | **0.976*** | 1.033* | 1,018 | 0,997 | 1 | 0.976* |
| AR,POOS-CV | **0.991** | **0.969**** | **0.990*** | **0.968**** | **0.968**** | 1,025 | 0,976 | 0,998 | 0.984** | 0.974* |
| AR,K-fold | **0.992** | **0.984** | 0,998 | **0.984**** | **0.988** | 1,032 | 1,007 | 0,997 | 0,993 | 0,989 |
| RRAR,POOS-CV | **0.974**** | **0.953**** | **0.964**** | **0.967*** | **0.958**** | 1.019* | 0,965 | 0,968 | 0,981 | 0.938*** |
| RRAR,K-fold | **1.000** | **0.983** | **0.988***** | **0.992*** | **0.976*** | 1,025 | 1,005 | 0.994** | 0,993 | 0.955** |
| RFAR,POOS-CV | **0.981** | **0.917**** | **0.917*** | **0.936** | 1,053 | 1,059 | 0,937 | 0,94 | 1,022 | 0,896 |
| RFAR,K-fold | **0.969** | **0.921**** | **0.923*** | **0.917*** | 1,025 | **1.030** | 0,936 | 0,947 | 1,013 | 0.795** |
| KRR-AR,POOS-CV | **1.042** | <u>**0.894****</u> | **0.867*** | **0.891** | **0.903*** | 1,178 | **0.873*** | 0,817 | **0.760**** | 0.775** |
| KRR,AR,K-fold | **0.997** | **0.908** | <u>**0.860*****</u> | <u>**0.870*****</u> | 1,009 | 1,021 | **0.855** | **0.770*** | **0.768**** | 0.783** |
| SVR-AR,Lin,POOS-CV | **1.011** | 1.198*** | 1.075* | 1.488** | 1.410*** | 1,04 | 1.084** | 1,001 | 1,202 | 1.300* |
| SVR-AR,Lin,K-fold | 1.563*** | 1.950*** | 1.914*** | 1.805*** | 1.662*** | 1.329* | 1.622*** | 1.293** | 1,116 | 0,948 |
| SVR-AR,RBF,POOS-CV | **0.990** | 1,007 | 1,04 | **1.058** | 1.188** | 1,009 | 0,933 | 1,017 | 1,114 | 1,002 |
| SVR-AR,RBF,K-fold | **1.083**** | 1.040** | 1,059 | 1.222** | 1.189** | 1.019** | 0,992 | 0.931*** | 1,032 | 0,865 |
| **Data-rich ($H_t^+$) models** | | | | | | | | | | |
| ARDI,BIC | **1.016** | **0.978** | **0.994** | **0.990** | **0.986** | 1,048 | **0.949** | 0,939 | <u>**0.714****</u> | 0.731** |
| ARDI,AIC | **1.043** | 1,027 | 1,052 | **1.050** | 1,068 | 1,104 | 0,99 | 0,924 | **0.844** | 0.806** |
| ARDI,POOS-CV | **1.091** | 1,055 | 1,084 | **1.013** | **0.918** | 1.221** | 1,113 | 1,015 | **0.751*** | 0.686** |
| ARDI,K-fold | **1.037** | 1,027 | 1.092* | **1.069** | 1,047 | 1,107 | 1,007 | 0,926 | 0,853 | 0.816** |
| RRARDI,POOS-CV | **1.010** | **1.041** | **1.037** | **1.000** | **0.990** | 1,058 | 1,063 | 0,977 | **0.720**** | **0.639**** |
| RRARDI,K-fold | **0.988** | 1,014 | 1.117* | **1.073** | 1,167 | 1,023 | 0,972 | 0,976 | 0,857 | 0.681*** |
| RFARDI,POOS-CV | <u>**0.963**</u> | **0.900**** | **0.895*** | **0.914** | 1,088 | 1,032 | 0,944 | 0,906 | 0,956 | 0.786*** |
| RFARDI,K-fold | **0.970** | **0.904**** | **0.931** | **0.946** | 1,040 | 1,046 | 0,932 | 0,924 | 1,026 | 0.786*** |
| KRR-ARDI,POOS-CV | **1.017** | **0.914** | **0.924** | **0.958** | **0.948** | **0.996** | **0.850*** | **0.783*** | **0.835*** | 0,902 |
| KRR,ARDI,K-fold | **0.988** | **0.925** | **0.893*** | **0.904*** | <u>**0.835****</u> | 1,045 | **0.858** | 0.842* | **0.822**** | 0.668** |
| ($B_1,\alpha=\hat{\alpha}$),POOS-CV | **1.133**** | 1.200*** | 1.195** | 1.310*** | 1.267** | **0.967** | 1,018 | **0.778*** | 1,005 | 0.833** |
| ($B_1,\alpha=\hat{\alpha}$),K-fold | 1.123** | 1.221*** | 1.187* | 1.316*** | 1.179* | 1,029 | **0.871** | **0.749**** | 0,905 | 0.766*** |
| ($B_1,\alpha=1$),POOS-CV | 1.251*** | 1.276*** | 1.208** | 1.221** | 1.403*** | 1,137 | 1,01 | **0.828** | 0,973 | 1,015 |
| ($B_1,\alpha=1$),K-fold | 1.368*** | 1.340*** | 1.412*** | 1.409*** | 1.270** | 1.280** | 0,91 | 0,957 | 0,903 | 0.726** |
| ($B_1,\alpha=0$),POOS-CV | 1.488** | 1.562** | 1.269** | 1.396** | 1.431*** | 1.153* | 0,961 | 0,979 | **0.793** | 1,307 |
| ($B_1,\alpha=0$),K-fold | 1.540** | 1.493** | 1.489** | 1.429** | 1.317** | 1.125* | **0.815** | **0.706*** | **0.738** | 1,074 |
| ($B_2,\alpha=\hat{\alpha}$),POOS-CV | 1.131*** | 1.249** | 1.152** | 1.193** | 1,111 | 1,051 | 1,268 | 0.903* | 0.843** | **0.637**** |
| ($B_2,\alpha=\hat{\alpha}$),K-fold | 1.111** | 1,266 | 1.103* | **1.142*** | 1,079 | 1,115 | 1,387 | 0,925 | **0.823*** | 0,749 |
| ($B_2,\alpha=1$),POOS-CV | **1.075**** | 1.078** | 1.095* | 1.233** | 1.259** | 1,026 | 0,974 | 0.912** | 0,884 | **0.606**** |
| ($B_2,\alpha=1$),K-fold | 1.078* | 1,315 | 1.098* | **1.130*** | 1,172 | 1,11 | 1,449 | 0,933 | **0.798**** | 0.679* |
| ($B_2,\alpha=0$),POOS-CV | 1.316** | 1.332** | 1.418*** | 1.393*** | 1.169* | 1,373 | 1.345* | 1,298 | 0,948 | **0.629***** |
| ($B_2,\alpha=0$),K-fold | 1.358** | 1.291** | 1.388** | 1.313** | 1,13 | 1,487 | 1,263 | 1,339 | 1,016 | <u>**0.597*****</u> |
| ($B_3,\alpha=\hat{\alpha}$),POOS-CV | **1.033*** | 1,009 | 1,063 | 1.092** | 1,102 | 1,016 | 0.945* | 0,972 | 0.885* | 0.854** |
| ($B_3,\alpha=\hat{\alpha}$),K-fold | **1.009** | 1,033 | 1.094*** | 1,056 | 1,101 | **1.000** | 1,001 | 0.946* | 0.936* | 0.790*** |
| ($B_3,\alpha=1$),POOS-CV | **1.010** | 1.042* | 1.086* | 1.101** | 1,12 | <u>**0.955***</u> | 0.953* | 0,993 | 0,923 | 0.824** |
| ($B_3,\alpha=1$),K-fold | **0.995** | 1,032 | 1.048** | **1.042** | 1.209** | 0.965** | 1,007 | 0,997 | 0,947 | 0.907* |
| ($B_3,\alpha=0$),POOS-CV | 1.084** | **1.001** | 1,017 | **1.016** | 1.117* | 1,067 | **0.910** | 0,904 | 0,917 | 0,885 |
| ($B_3,\alpha=0$),K-fold | **1.071*** | 1.198* | 1,12 | 1.133* | 1.127* | 1.085* | 1,149 | 0,979 | 0,948 | 0,923 |
| SVR-ARDI,Lin,POOS-CV | **1.086*** | 1.271*** | 1.292*** | 1.228** | 1.220** | **1.009** | 1,13 | 1,081 | 0,945 | 0,97 |
| SVR-ARDI,Lin,K-fold | 1.136* | 1.161* | 1.351* | 1.301** | 1.169* | 1.228* | **0.881** | 1,173 | 1,145 | 1,026 |
| SVR-ARDI,RBF,POOS-CV | 1.236 | 1,019 | 1,017 | **0.958** | **0.991** | 1,47 | 0,968 | 0,939 | **0.768***** | 0.798** |
| SVR-ARDI,RBF,K-fold | **1.054** | 1,062 | 1,063 | 1.236*** | 1,075 | 1,096 | 1,048 | 0,909 | 0,985 | 0,891 |

Note: The numbers represent the relative. with respect to AR,BIC model. root MSPE. Models retained in model confidence set are in bold. the minimum values are underlined. while ***. **. * stand for 1%. 5% and 10% significance of Diebold-Mariano test.
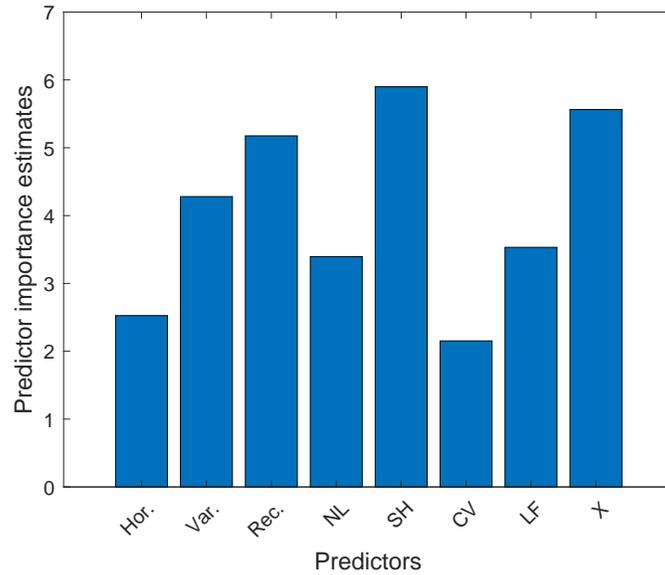
**Figure 32:** This figure presents predictive importance estimates. Random forest is trained to predict $R^2_{t,h,v,m}$ defined in (11) and use out-of-bags observations to assess the performance of the model and compute features' importance. NL, SH, CV and LF stand for nonlinearity, shrinkage, cross-validation and loss function features respectively. A dummy for $H_t^+$ models, $X$, is included as well.
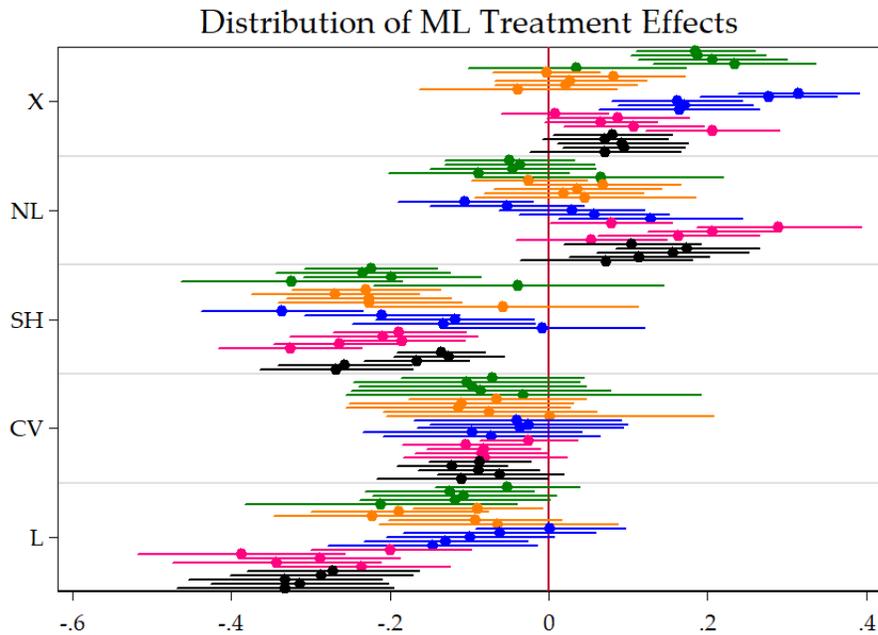


**Figure 33:** This figure plots the distribution of $\hat{\alpha}_F^{(h,v)}$ from equation (11) done by $(h,v)$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^2$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. Finally, variables are GDP, CONS, INV, INC and PCE. Within a specific color block, the horizon increases from $h = 1$ to $h = 8$ as we are going down. SEs are HAC. These are the 95% confidence bands.

15

## Distribution of averaged ML Treatment Effects



Figure 34: This figure plots the distribution of $\hat{\alpha}_F^{(v)}$ and $\hat{\alpha}_F^{(h)}$ from equation (11) done by $h$ and $v$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^2$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. However, in this graph, $v-$specific heterogeneity and $h-$specific heterogeneity have been integrated out in turns. SEs are HAC. These are the 95% confidence bands.

## Contribution of Non-Linearities, by variables



Figure 35: This compares the two NL models averaged over all horizons. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

16

## Contribution of Non-Linearities, by horizons



Figure 36: This compares the two NL models averaged over all variables. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.
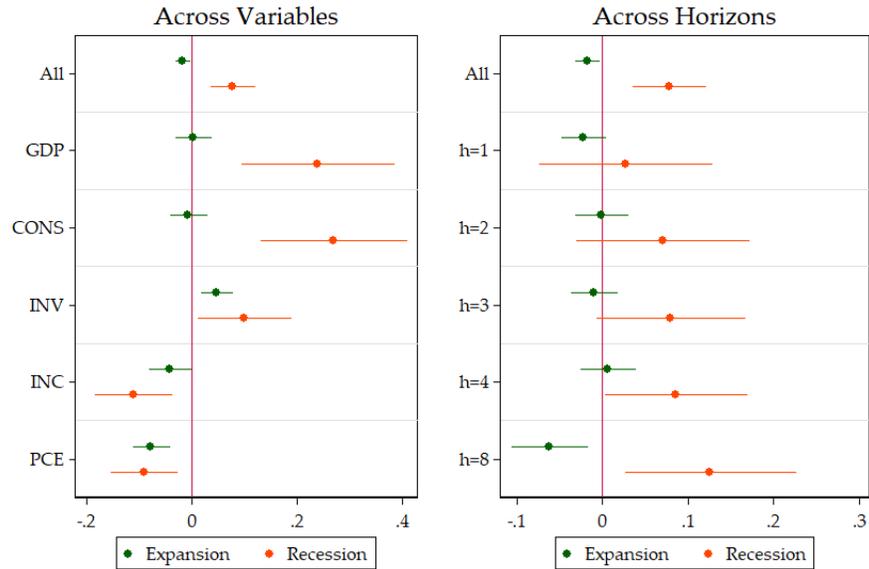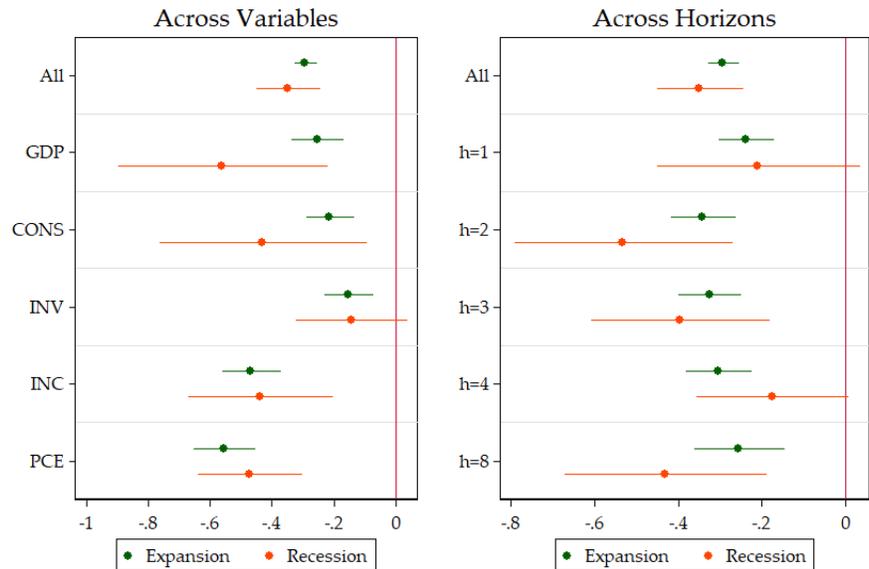
## Alternative shrinkage wrt ARDI



Figure 37: This compares models of section 3.2 averaged over all variables and horizons. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. The base models are ARDIs specified with POOS-CV and KF-CV respectively. SEs are HAC. These are the 95% confidence bands.

## Table 15: CV comparison

|  | (1) All | (2) Data-rich | (3) Data-poor | (4) Data-rich | (5) Data-poor |
|---|---|---|---|---|---|
| CV-KF | -4.248* | -8.304*** | -0.192 | -9.651*** | 0.114 |
|  | (1.940) | (1.787) | (0.424) | (1.886) | (0.386) |
| CV-POOS | -2.852 | -6.690** | 0.985** | -6.772** | 0.917* |
|  | (1.887) | (2.163) | (0.382) | (2.270) | (0.386) |
| AIC | -2.182 | -4.722** | 0.358 | -5.557** | 0.373 |
|  | (1.816) | (1.598) | (0.320) | (1.694) | (0.303) |
| CV-KF * Recessions |  |  |  | 13.21** | -2.956* |
|  |  |  |  | (4.893) | (1.500) |
| CV-POOS * Recessions |  |  |  | 1.002 | 0.683 |
|  |  |  |  | (5.345) | (1.125) |
| AIC * Recessions |  |  |  | 8.421 | -0.127 |
|  |  |  |  | (4.643) | (1.101) |
| Observations | 36960 | 18480 | 18480 | 18360 | 18360 |

Standard errors in parentheses

$^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$



Figure 38: This compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

**Figure 39:** This compares the two CVs procedure averaged over all the models that use them. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.
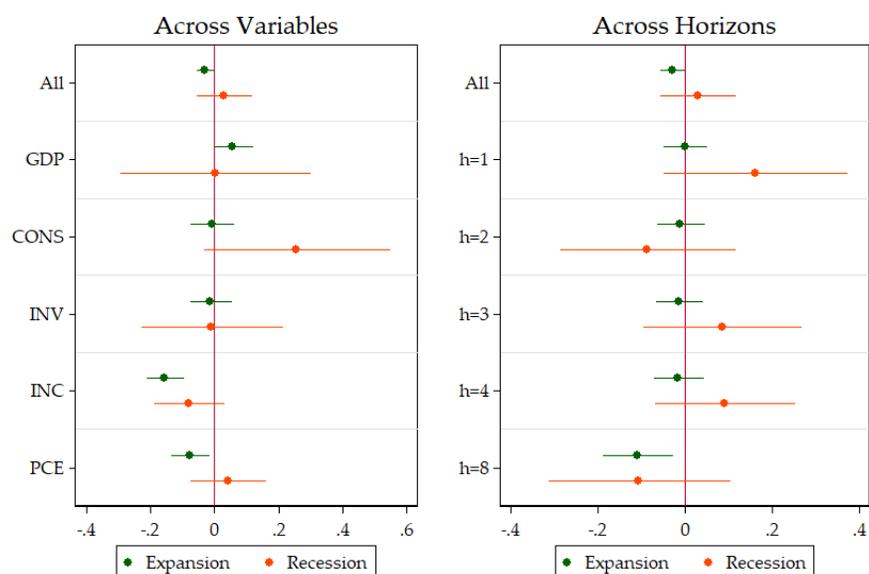


**Figure 40:** This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in **both the data-poor and data-rich environments**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

Figure 41: This graph display the marginal (un)improvements by variables and horizons to opt for the SVR in-sample loss function in **both recession and expansion periods**. The unit of the x-axis are improvements in OOS $R^2$ over the basis model. SEs are HAC. These are the 95% confidence bands.

# C   Results with Canadian data

In this section we present results obtained with Canadian data from Fortin-Gagnon et al. (2020). It is a monthly dataset of 139 macroeconomic and financial variables, with categories similar to those from McCracken and Ng (2016), except that it contains much more international trade indicators to take into account the openness of Canadian economy. Data starts on 1981M01 and ends on 2017M12. The out-of-sample starts on 2000M01. The variables of interest are the same as in US application: industrial growth, unemployment rate change, term spread, CPI inflation and housing starts growth. Forecasting horizons are 1, 3, 9, 12 and 24 months. We do not compute results for recession periods separately since Canada has experienced only one downturn in the evaluation period.

The results with Canadian data are overall similar to those in the paper. The main difference is a smaller NL treatment effect. That can be potentially explained through lenses of the analysis in section 6. The pseudo-out-of-sample covers 2000-2017 period during which Canadian financial system did not experience a dramatic nonfinancial cycle as in the US., and the housing bubble did not burst. The main reason for this discrepancy being more concentrated and strictly regulated (since 80's) Canadian financial system (Bordo et al., 2015). Hence, the nonlinearities associated to financial frictions found in the US case were probably less important and nonlinear methods did not have a significant effect on predicting real activity series on average. However, NL treatment is very important for inflation and housing. Shrinkage is still not a good idea for industrial production and unemployment rate, but can be very helpful other variables at some specific horizons. Cross-validation does not have a big impact and the SVR loss function is still harmful.
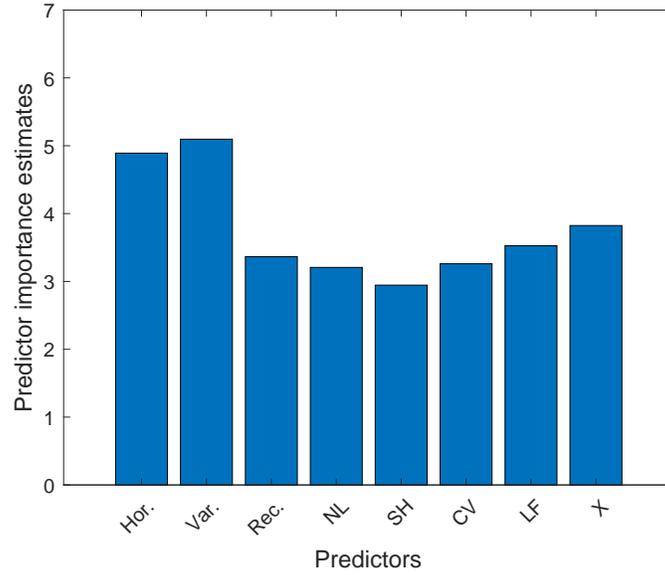
**Figure 42:** This figure presents predictive importance estimates. Random forest is trained to predict $R^2_{t,h,v,m}$ defined in (11) and use out-of-bags observations to assess the performance of the model and compute features' importance. NL, SH, CV and LF stand for nonlinearity, shrinkage, cross-validation and loss function features respectively. A dummy for $H_t^+$ models, $X$, is included as well.
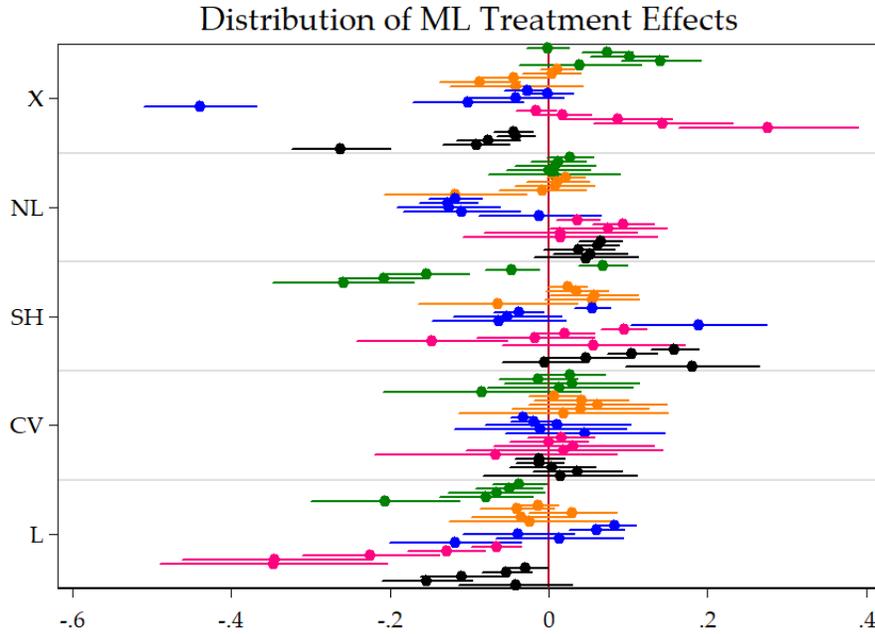


**Figure 43:** This figure plots the distribution of $\dot{\alpha}_F^{(h,v)}$ from equation (11) done by $(h,v)$ subsets. That is, we are looking at the average partial effect on the pseudo-OOS $R^2$ from augmenting the model with ML features, keeping everything else fixed. $X$ is making the switch from data-poor to data-rich. Finally, variables are INDPRO, UNRATE, SPREAD, INF and HOUS. Within a specific color block, the horizon increases from $h=1$ to $h=24$ as we are going down. SEs are HAC. These are the 95% confidence bands.

# D  Detailed Implementation of Cross-validations

All of our models involve some kind of hyperparameter selection prior to estimation. To curb the overfitting problem, we use two distinct methods that we refer to loosely as cross-validation methods. To make it feasible, we optimize hyperparameters every 24 months as the expanding window grows our in-sample set. The resulting optimization points are the same across all models, variables and horizons considered. In all other periods, hyperparameter values are frozen to the previous values and models are estimated using the expanded in-sample set to generate forecasts.



POOS                                                                      K folds

- Training (INS)
- Test (INS)
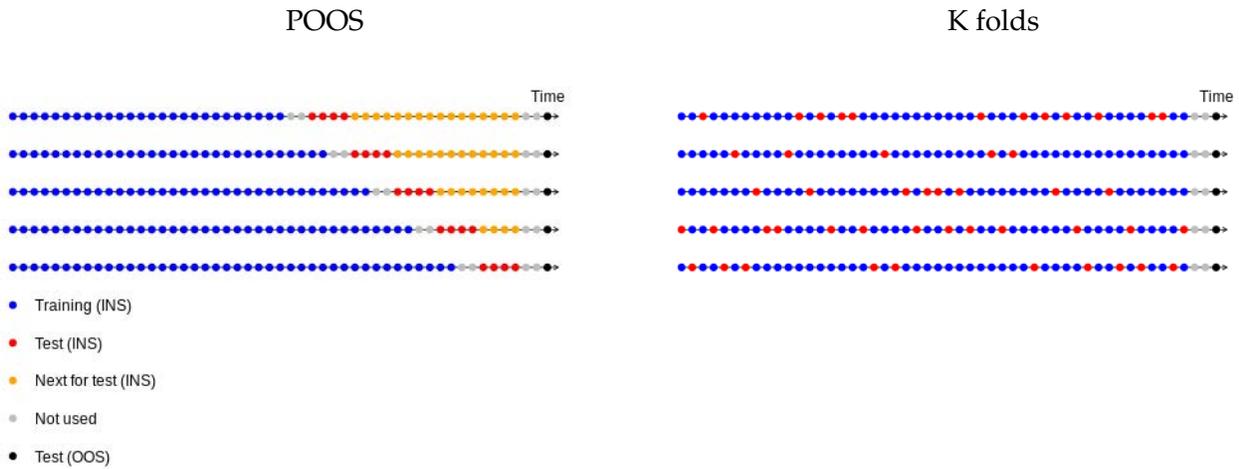- Next for test (INS)
- Not used
- Test (OOS)

Figure 44: Illustration of cross-validation methods

Notes: Figures are drawn for 3 months forecasting horizon and depict the splits performed in the in-sample set. The pseudo-out-of-sample observation to be forecasted here is shown in black.

The first cross-validation method we consider mimics in-sample the pseudo-out-of-sample comparison we perform across models. For each set of hyperparameters considered, we keep the last 25% of the in-sample set as a comparison window. Models are estimated every 12 months, but the training set is gradually expanded to keep the forecasting horizon intact. This exercise is thus repeated 5 times. Figure 44 shows a toy example with smaller jumps, a smaller comparison window and a forecasting horizon of 3 months, hence the gaps. Once hyperparameters have been selected, the model is estimated using the whole in-sample set and used to make a forecast in the pseudo-out-of-sample window that we use to compare all models (the black dot in the figure). This approach is a compromise between two methods used to evaluate time series models detailed in Tashman (2000), rolling-origin recalibration

and rolling-origin updating.[29] For a simulation study of various cross-validation methods in a time series context, including the rolling-origin recalibration method, the reader is referred to Bergmeir and Benítez (2012). We stress again that the compromise is made to bring down computation time.

The second cross-validation method, K-fold cross-validation, is based on a re-sampling scheme (Bergmeir et al., 2018). We chose to use 5 folds, meaning the in-sample set is randomly split into five disjoint subsets, each accounting on average for 20 % of the in-sample observations. For each one of the 5 subsets and each set of hyperparameters considered, 4 subsets are used for estimation and the remaining corresponding observations of the in-sample set used as a test subset to generate forecasting errors. This is illustrated in figure 44 where each subset is illustrated by red dots on different arrows.

Note that the average mean squared error in the test subset is used as the performance metric for both cross-validation methods to perform hyperparameter selection.

# E    Forecasting models in detail

## E.1    Data-poor ($H_t^-$) models

In this section we describe forecasting models that contain only lagged values of the dependent variable, and hence use a small amount of predictors, $H_t^-$.

**Autoregressive Direct (AR)**    The first univariate model is the so-called *autoregressive direct* (AR) model, which is specified as:

$$y_{t+h}^{(h)} = c + \rho(L)y_t + e_{t+h}, \quad t = 1, \ldots, T,$$

where $h \geq 1$ is the forecasting horizon. The only hyperparameter in this model is $p_y$, the order of the lag polynomial $\rho(L)$. The optimal $p$ is selected in four ways: (i) Bayesian Information Criterion (AR,BIC); (ii) Akaike Information Criterion (AR,AIC); (iii) Pseudo-out-of-sample cross validation (AR,POOS-CV); and (iv) K-fold cross validation (AR,K-fold). The lag order is selected from the following subset $p_y \in \{1, 3, 6, 12\}$. Hence, this model enters the following categories: linear $g$ function, no regularization, in-sample and cross-validation selection of hyperparameters and quadratic loss function.

**Ridge Regression AR (RRAR)**    The second specification is a penalized version of the previous AR model that allows potentially more lagged predictors by using Ridge regression. The model is written as in (6), and the parameters are estimated using Ridge penalty. The Ridge hyperparameter is selected with two cross validation strategies, which gives two mod-

---

[29]In both cases, the last observation (the origin of the forecast) of the training set is rolled forward. However, in the first case, hyperparameters are recalibrated and, in the second, only the information set is updated.

els: RRAR,POOS-CV and RRAR,K-fold. The lag order is selected from the following subset $p_y \in \{1, 3, 6, 12\}$ and for each of these value we choose the Ridge hyperparameter. This model creates variation on following axes: linear $g$, Ridge regularization, cross-validation for tuning parameters and quadratic loss function.

**Random Forests AR (RFAR)** A popular way to introduce nonlinearities in the predictive function $g$ is to use a tree method that splits the predictors space in a collection of dummy variables and their interactions. Since a standard tree regression is prompt to the overfit, we use instead the random forest approach described in Section 3.1.2. We adopt the default value in the literature of one third for 'mtry', the share of randomly selected predictors that are candidates for splits in each tree. Observations in each set are sampled with replacement to get as many observations in the trees as in the full sample. The number of lags of $y_t$, is chosen from the subset $p_y \in \{1, 3, 6, 12\}$ with cross-validation while the number of trees is selected internally with out-of-bag observations. This model generates nonlinear approximation of the optimal forecast, without regularization, using both CV techniques with the quadratic loss function: RFAR,K-fold and RFAR,POOS-CV.

**Kernel Ridge Regression AR (KRRAR)** This specification adds a nonlinear approximation of the function $g$ by using the Kernel trick as in Section 3.1.1. The model is written as in (13) and (14) but with the autoregressive part only

$$y_{t+h} = c + g(Z_t) + \varepsilon_{t+h},$$
$$Z_t = \left[ \{y_{t-0}\}_{j=0}^{p_y} \right],$$

and the forecast is obtained using the equation (16). The hyperparameters of Ridge and of its kernel are selected by two cross-validation procedures, which gives two forecasting specifications: (i) KRRAR,POOS-CV, (ii) KRRAR,K-fold. $Z_t$ consists of $y_t$ and its $p_y$ lags, $p_y \in \{1, 3, 6, 12\}$. This model is representative of a nonlinear $g$ function, Ridge regularization, cross-validation to select $\tau$ and quadratic $\hat{L}$.

**Support Vector Regression AR (SVR-AR)** We use the SVR model to create variation along the loss function dimension. In the data-poor version the predictors set $Z_t$ contains $y_t$ and a number of lags chosen from $p_y \in \{1, 3, 6, 12\}$. The hyperparameters are selected with both cross-validation techniques, and we consider 2 kernels to approximate basis functions, linear and RBF. Hence, there are 4 versions: (i) SVR-AR,Lin,POOS-CV, (ii) SVR-AR,Lin,K-fold, (iii) SVR-AR,RBF,POOS-CV and (iv) SVR-AR,RBF,K-fold. The forecasts are generated using (19).

# E.2 Data-rich ($H_t^+$) models

We now describe forecasting models that use a large dataset of predictors, including the autoregressive components, $H_t^+$.

**Diffusion Indices (ARDI)**  The reference model in the case of large predictor set is the autoregression augmented with diffusion indices from Stock and Watson (2002b):

$$y_{t+h}^{(h)} = c + \rho(L)y_t + \beta(L)F_t + e_{t+h}, \quad t = 1, \dots, T \tag{20}$$

$$X_t = \Lambda F_t + u_t \tag{21}$$

where $F_t$ are $K$ *consecutive static* factors, and $\rho(L)$ and $\beta(L)$ are lag polynomials of orders $p_y$ and $p_f$ respectively. The feasible procedure requires an estimate of $F_t$ that is usually done by PCA. The optimal values of hyperparamters $p$, $K$ and $m$ are selected in four ways: (i) Bayesian Information Criterion (ARDI,BIC); (ii) Akaike Information Criterion (ARDI,AIC); (iii) Pseudo-out-of-sample cross validation (ARDI,POOS-CV); and (iv) K-fold cross validation (ARDI,K-fold). These are selected from following subsets: $p_y \in \{1, 3, 6, 12\}$, $K \in \{3, 6, 10\}$, $p_f \in \{1, 3, 6, 12\}$. Hence, this model following features: linear $g$ function, PCA regularization, in-sample and cross-validation selection of hyperparameters and $L^2$.

**Ridge Regression Diffusion Indices (RRARDI)**  As for the small data case, we explore how a regularization affects the predictive performance of the reference model ARDI above. The predictive regression is written as in (7) and $p_y$, $p_f$ and $K$ are selected from the same subsets of values as for the ARDI case above. The parameters are estimated using Ridge penalty. All the hyperparameters are selected with two cross validation strategies, giving two models: RRARDI,POOS-CV and RRARDI,K-fold. This model creates variation on following axes: linear $g$, Ridge regularization, CV for tuning parameters and $L^2$.

**Random Forest Diffusion Indices (RFARDI)**  We also explore how nonlinearities affect the predictive performance of the ARDI model. The model is as in (7) but a Random Forest of regression trees is used. The ARDI hyperparameters are chosen from the grid as in the linear case, while the number of trees is selected with out-of-bag observations. Both POOS and K-fold CV are used to generate two forecasting models: RFARDI,POOS-CV and RFARDI,K-fold. This model generates nonlinear treatment, with PCA regularization, using both CV techniques with the quadratic loss function.

**Kernel Ridge Regression Diffusion Indices (KRRARDI)**  As for the autoregressive case, we can use the KT to generate nonlinear predictive functions $g$. The model is represented by equations (13) - (15) and the forecast is obtained using the equation (16). The hyperparameters of Ridge and of its kernel, as well as $p_y$, $K$ and $p_f$ are selected by two cross-validation procedures, which gives two forecasting specifications: (i) KRRARDI,POOS-CV, (ii) KRRARDI,K-fold. We use the same grid as in ARDI case for discrete hyperparameters. This model is representative of a nonlinear $g$ function, Ridge regularization with PCA, cross-validation to select $\tau$ and quadratic $\hat{L}$.

**Support Vector Regression ARDI (SVR-ARDI)**   We use four versions of the SVR model: (i) SVR-ARDI,Lin,POOS-CV, (ii) SVR-ARDI,Lin,K-fold, (iii) SVR-ARDI,RBF,POOS-CV and (iv) SVR-ARDI,RBF,K-fold. The SVR hyperparameters are chosen by cross-validation and the ARDI hyperparameters are chosen using a grid that search in the same subsets as the ARDI model. The forecasts are generated from equation (19). This model creates variations in all categories: nonlinear $g$, PCA regularization, CV and $\bar{\epsilon}$-insensitive loss function.

### E.2.1   Generating shrinkage schemes

The rest of the forecasting models relies on using different $B$ operators to generate variations across shrinkage schemes, as depicted in section 3.2.

$B_1$: **taking all observables** $H_t^+$   When $B$ is identity mapping, we consider $Z_t = H_t^+$ in the Elastic Net problem (18), where $H_t^+$ is defined by (5). The following lag structures for $y_t$ and $X_t$ are considered, $p_y \in \{1, 3, 6, 12\}$ $p_f \in \{1, 3, 6, 12\}$, and the exact number is cross-validated. The hyperparameter $\lambda$ is always selected by two cross validation procedures, while we consider three cases for $\alpha$: $\hat{\alpha}$, $\alpha = 1$ and $\alpha = 0$, which correspond to EN, Ridge and Lasso specifications respectively. In case of EN, $\alpha$ is also cross-validated. This gives six combinations: $(B_1, \alpha = \hat{\alpha})$,POOS-CV; $(B_1, \alpha = \hat{\alpha})$,K-fold; $(B_1, \alpha = 1)$,POOS-CV; $(B_1, \alpha = 1)$,K-fold; $(B_1, \alpha = 0)$,POOS-CV and $(B_1, \alpha = 0)$,K-fold. They create variations within regularization and hyperparameters' optimization.

$B_2$: **taking all principal components of** $X_t$   Here $B_2()$ rotates $X_t$ into $N$ factors, $F_t$, estimated by principal components, which then constitute $Z_t$ to be used in (18). Same lag structures and hyperparameters' optimization from the $B_1$ case are used to generate the following six specifications: $(B_2, \alpha = \hat{\alpha})$,POOS-CV; $(B_2, \alpha = \hat{\alpha})$,K-fold; $(B_2, \alpha = 1)$,POOS-CV; $(B_2, \alpha = 1)$,K-fold; $(B_2, \alpha = 0)$,POOS-CV and $(B_2, \alpha = 0)$,K-fold.

$B_3$: **taking all principal components of** $H_t^+$   Finally, $B_3()$ rotates $H_t^+$ by taking all principal components, where $H_t^+$ lag structure is to be selected as in the $B_1$ case. Same variations and hyperparameters' selection are used to generate the following six specifications: $(B_3, \alpha = \hat{\alpha})$,POOS-CV; $(B_3, \alpha = \hat{\alpha})$,K-fold; $(B_3, \alpha = 1)$,POOS-CV; $(B_3, \alpha = 1)$,K-fold; $(B_3, \alpha = 0)$,POOS-CV and $(B_3, \alpha = 0)$,K-fold.