

# Macroeconometrics and Machine Learning Project

Louis Blanco & Jeanne Astier

31 décembre 2022

This project is based on the article "Forecasting Inflation in a Data-Rich Environment : The Benefits of Machine Learning Methods" by Marcelo C. Medeiros, Gabriel F. R. Vasconcelos, Álvaro Veiga and Eduardo Zilberman (2019). It is organized as follows : part 1 explains the problem under study and its interests. Part 2 describes our replication of the article results in forecasting US inflation. Part 3 presents an extension of these results, applying the same methodology to forecast French inflation. Part 4 provides economic interpretation of the results and methodology used.

## 1 Introduction

*Price Gains Slow More Than Expected* : this is the title of the November Inflation Report published by the New York Times on December 13<sup>th</sup> 2022. It highlights the complexity of inflation forecasting : after a year during which the Fed repeatedly underestimated inflation's hold on the U.S.<sup>1</sup>, it seems now to be overestimating it. These forecast errors, common to many Central Banks (the ECB addresses in its Economic Bulletin 2022 its recent errors in the inflation projections), are particularly problematic as inflation forecast is of paramount importance for economic decisions. Central Banks rely on inflation forecasts to inform their monetary policy. These forecasts are also used by households, private sector companies, policymakers to build their expectations and make their decisions. In some macroeconomic models, not only realized inflation but also expectations of inflation play a central role. Hence there is a discrepancy between the prime importance of inflation forecasts and the their lack of accuracy. These forecasts errors also undermine Central Banks credi-

---

1. *2022 shattered economic forecasts. Can the Fed get 2023 right ?*, The Washington Post, December 12<sup>th</sup> 2022.

bility, degrading again the expectations. Poor inflation forecasts might thus have important welfare costs.

In this context, there is a challenge to improve performance of forecasting models. According to literature, the traditional simple univariate forecasting models turn out to be difficult to improve. Even the attempts to use Machine Learning (ML) models on a small set of variables, or "Big data" summarized in a small set of factors, did not seem conclusive. However, combining the two, as it has already been done outside of the economic field, might lead to significant improvement. The objective of the article studied is not only to beat univariate benchmark models, but also to prove that the combination of ML and "Big data" is the more effective than using just ML with few variables of Big Data with factor models.

## 2 Replication of article results

We replicate the article results by estimating benchmark models and several ML ones on a large set of variables to predict US inflation. We then compare the performances of the forecasts made by these different models, to determine which ones are the most accurate.

We face several obstacles in this replication, which lead us to make choices in how to run it. The main one is the computing power required to run all the models compared. Not only the ML models are power-consuming to be trained, but also all the processes are replicated many times (forecasts are computed at 12 different time horizons at each one of the 312 out-of-sample windows, with a rolling-window framework). This factor is highly limiting since our laptops do not have very high computing power, so the predictions can take very long to be computed<sup>2</sup>. This explains our choice to stick to the replication of general results of the paper, and not to replicate exhaustively all the robustness checks made by the authors that would have asked several days of code running with our tools. For this reason, we also gave up on replicating forecasts for 4 models among the 18 tested.

An other issue is that the data available to us do not seem to correspond exactly to what the authors describe. Even sticking to their description of the source of the data (January 2016 edition of FRED-MD data) and to the treatment they apply (among others, removing the variables with missing dates in the time period considered), we end up with slightly less variables than they do,

---

2. The total computing time for the replication results presented is about 35 hours.

without being able to identify which ones and the reasons for this. This might explain why our results statistics do not always exactly correspond to the article original ones.

As in the original article, we train several models on FRED-MD data from 1960 ; we compare performance of these different models in an out-of-sample window from January 1990 to December 2015. At each date  $t$  of the out-of-sample window, the compute forecasts based on a rolling-window framework of fixed length, at time horizons  $t + 1, t + 1, \dots, t + 12$ . We also compute forecasts for the accumulated inflation over the following 3, 6, and 12 months. The models are compared according to different statistics : root mean squared error (RMSE) and mean absolute error (MAE). We also compute, for square and absolute losses, respectively, the average p-values (accross horizons) for the model confidence sets (MCSs) based on the Tmax statistic as described in Hansen, Lunde, and Nason (2011). The estimated models are the following ones :

- benchmark models : Random Walk (RW), Auto-Regressive (AR)
- shrinkage models : Ridge regression, Least Absolute Shrinkage and Selection Operator (LASSO), adaptive LASSO (adaLASSO), Elastic Net (Elnet), adaptive Elastic Net (adaElnet)
- factor models : classic factor model (Factor), Target factor model (T. Factor)
- ensemble methods : Bagging, Complete Subset Regressions (CSR)
- Breiman Random Forest (RF)
- hybrid linear-Random Forest model : Random-Forest Ordinary Least Squares (RF/OLS).

The results of this replication are summarized in Tables 1 and 2.

### 3 Extension : application to forecasting of French inflation

In the conclusion of their article, Marcelo C. Medeiros, Gabriel F. R. Vasconcelos, Álvaro Veiga and Eduardo Zilberman write : "*Although our article focuses on inflation forecasting in the US, one can easily apply ML methods to forecast other macroeconomic series in a variety of countries.*" Hence as an extension, we try to apply the benchmark and ML models comparison implemented in their article to French data, aiming at forecasting French inflation. There is a double issue at stake with this extension. First, is the domination of Random Forest over all other models (see part 4) specific to US data, or does it hold on a different dataset ? Second, are the ML models still performing better

**TABLE 1** – Replication of forecasting results : summary statistics for the out-of-sample period from 1990 to 2015

Models	Forecasting precision								Model confidence set	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
	ave. RMSE	ave. MAE	max RMSE	max MAE	min RMSE	min MAE	# min RMSE	# min MAE	ave. p-v. Tmax sq	ave. p-v. Tmax abs
RW	1.00	1.00	1.00	1.00	1.00	1.00	0	0	0.00	0.00
AR	0.61	0.81	0.88	0.87	0.56	0.70	2	0	1.00	0.32
LASSO	0.75	0.94	0.84	1.01	0.62	0.67	3	0	0.28	0.32
adaLASSO	0.85	0.96	0.88	1.04	0.64	0.66	0	0	0.07	0.13
ElNet	0.76	0.87	0.81	1.05	0.62	0.68	1	1	0.26	0.32
adaElNet	0.85	0.98	0.88	1.06	0.65	0.69	1	1	0.08	0.13
Ridge	0.79	0.89	0.86	0.92	0.73	0.71	0	1	0.05	0.13
BVAR	0.62	0.82	0.90	0.89	0.57	0.70	0	0	0.32	0.32
Bagging	1.38	1.25	1.50	1.35	0.82	0.83	0	0	0.05	0.05
CSR	0.74	0.90	0.82	0.95	0.65	0.69	0	0	0.07	0.13
Factor	0.73	0.89	0.83	1.01	0.58	0.66	1	2	0.29	0.32
T. Factor	0.65	0.80	0.84	0.85	0.58	0.66	1	0	0.32	0.32
RF	0.69	0.80	0.83	0.83	0.66	0.71	5	9	0.29	1.00
RF/OLS	0.77	0.87	0.81	0.91	0.64	0.71	1	1	0.20	0.32

NOTE : The table reports for each model different summary statistics across all the forecasting horizons (1 to 12 months and accumulated 3, 6 and 12-month). Columns (1) and (2) report the average root mean square error (RMSE) and the average mean absolute error (MAE). Columns (3) and (4) report, respectively, the maximum RMSE and MAE over the horizons considered. Columns (5) and (6) report, respectively, the minimum RMSE and MAE over the horizons considered. Assume these statistics (columns 1 to 6) are normalized for the benchmark RW model to one. Columns (7) and (8) report the number of times (across horizons) each model achieved the lowest RMSE and MAE, respectively. Columns (9) and (10) present for square and absolute losses, the average p-values for the model confidence sets (MCSs) based on the Tmax statistic as described in Hansen, Lunde, and Nason (2011).

than benchmarks with a smaller number of covariates? Indeed, we cannot not find as many monthly series available on a large time period for France as we have for the US with FRED-MD dataset. We download monthly series about France from FRED database, but way less series are available than for the US, and additionally most of them start later than 1960. We will thus see how much "Big Data" is needed for ML models to outperform univariate ones. Our dataset is made of FRED monthly time series about France. The variables cover the same fields as the US ones used in the original article : output and income ; labor market ; housing ; consumption, orders, and inventories ; money and credit ; interest and exchange rates ; prices ; stock market. The dataset we use has 64 variables from January 1985 to December 2020. Additionally to FRED time series, we use French

**TABLE 2** – Replication of forecasting results : RMSE and MAE ratios (1990–2015)

<b>Panel (a) : RMSE ratio</b>															
Forecasting horizon															
Model	1	2	3	4	5	6	7	8	9	10	11	12	3m	6m	12m
AR	0,88	0,80	0,78	0,79	0,77	0,77	0,76	0,74	0,75	0,81	0,81	0,74	<b>0,64</b>	<b>0,66</b>	<b>0,56</b>
RF	<b>0,84</b>	<b>0,75</b>	<b>0,73</b>	<b>0,76</b>	<b>0,73</b>	<b>0,74</b>	<b>0,75</b>	<b>0,74</b>	<b>0,74</b>	<b>0,78</b>	<b>0,80</b>	<b>0,72</b>	0,66	0,73	0,67

  

<b>Panel (b) : MAE ratio</b>															
Forecasting horizon															
Model	1	2	3	4	5	6	7	8	9	10	11	12	3m	6m	12m
AR	0,87	0,78	0,77	0,78	0,79	0,78	0,75	0,74	0,78	0,83	0,85	0,75	<b>0,70</b>	<b>0,80</b>	0,82
RF	<b>0,81</b>	<b>0,72</b>	<b>0,71</b>	<b>0,74</b>	<b>0,74</b>	<b>0,75</b>	<b>0,73</b>	<b>0,72</b>	<b>0,75</b>	<b>0,77</b>	<b>0,81</b>	<b>0,73</b>	0,71	0,84	<b>0,81</b>

NOTE : The table reports, for each forecasting horizon (1 to 12 months and accumulated 3, 6 and 12 months), the root mean squared error (RMSE) and mean absolute error (MAE) ratios with respect to the random walk model for the full out-of-sample period (1990–2015). The statistics for the best-performing model are highlighted in bold.

inflation data series from OECD database, as a larger time span is available.

We train and test on these French data the same models we have trained and tested on the US data; but taking into account the different time span available, their performances are compared on a reduced out-of-sample window, from 2005 to 2020. Hence, compared to the replication part, the lengths of both train and test datasets are reduced. We compute, at each date  $t$  of the out-of-sample window, the same forecasts based on a rolling-window framework of fixed length, at time horizons  $t+1$ ,  $t+1$ , ...,  $t+12$ , as well as accumulated inflation over the following 3, 6, and 12 months. RMSE and MAE statistics and average p-values for MCSs based on the Tmax are used to compare the different models.

The results of this replication are available in Tables 3 and 4.

## 4 Interpretation of the results

Replication of forecasting US inflation leads to results summarized in Tables 1 and 2. Table 1 reports summary statistics across all forecasting horizons, on the period 1990 to 2015. It shows that ML models (except Bagging) and the use of a large set of predictors systematically improve the quality of inflation forecasts over RW benchmark. The AR benchmark however seems to perform better than several ML models. The RF model clearly outperforms all the other ML alternatives; it

**TABLE 3** – Extension : results for forecasting of French inflation - summary statistics for the out-of-sample period from 2000 to 2015

Models	Forecasting precision								Model confidence set	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
	ave. RMSE	ave. MAE	max RMSE	max MAE	min RMSE	min MAE	# min RMSE	# min MAE	ave. p-v. Tmax sq	ave. p-v. Tmax abs
RW	1.00	1.00	1.00	1.00	1.00	1.00	0	0	0.00	0.18
AR	0.59	0.61	1.05	1.08	0.32	0.36	0	0	0.01	0.15
LASSO	0.78	0.64	1.73	1.17	0.11	0.12	1	1	0.00	0.18
adaLASSO	0.90	0.62	2.41	1.22	0.33	0.25	0	2	0.00	0.06
ElNet	0.70	0.68	1.48	1.21	0.12	0.12	0	0	0.00	0.07
adaElNet	0.75	0.63	1.54	1.22	0.39	0.27	1	1	0.00	0.07
Ridge	0.87	0.94	2.84	2.99	0.64	0.70	2	3	0.17	0.00
BVAR	0.57	0.59	1.01	1.06	0.30	0.34	0	0	0.00	0.00
Bagging	5.16	1.40	20.32	4.01	0.89	0.44	0	0	0.00	0.00
CSR	0.72	0.61	1.41	1.19	0.44	0.39	0	0	0.00	0.01
Factor	0.49	0.52	0.96	0.99	0.23	0.28	11	8	1.00	1.00
RF	0.61	0.62	1.27	1.19	0.36	0.36	0	0	0.01	0.01
RF/OLS	1.13	0.71	3.12	1.49	0.30	0.30	0	0	0.00	0.00

NOTE : The table reports for each model different summary statistics across all the forecasting horizons (1 to 12 months and accumulated 3, 6 and 12-month). Columns (1) and (2) report the average root mean square error (RMSE) and the average mean absolute error (MAE). Columns (3) and (4) report, respectively, the maximum RMSE and MAE over the horizons considered. Columns (5) and (6) report, respectively, the minimum RMSE and MAE over the horizons considered. As these statistics (columns 1 to 6) are normalized for the benchmark RW model to one. Columns (7) and (8) report the number of times (across horizons) each model achieved the lowest RMSE and MAE, respectively. Columns (9) and (10) present for square and absolute losses, the average p-values for the model confidence sets (MCSs) based on the Tmax statistic as described in Hansen, Lunde, and Nason (2011).

also seems to beat the AR benchmark. RF has the lowest MAEs across the horizons and its RMSEs are of same orders of magnitude as AR ones. The RF model also has the highest p-value in the MCS with absolute errors (but AR model has the highest p-value with squared errors). Finally, RF model presents the highest number of lowest RMSEs and MAEs accross horizons.

Our replication thus leads to the same conclusion as the original article : superiority of RF. However, our results show a slightly more qualified supremacy of RF than the original results. In particular, in the replication the AR benchmark seems to perform better than in the original article, and RF model slightly worse, leading to close forecast qualities for those two models. This could be explained by the issue mentioned in Part 2, that the data we use to contains less variables than what is described in the article. The "Big Data" would then be less "Big", leading to lower accuracy in

**TABLE 4** – Extension : results for forecasting of French inflation - RMSE and MAE ratios (1990–2015)

<b>Panel (a) : RMSE ratio</b>															
Forecasting horizon															
Model	1	2	3	4	5	6	7	8	9	10	11	12	3m	6m	12m
AR	<b>1,01</b>	<b>1,04</b>	<b>1,05</b>	<b>1,04</b>	1,03	1,01	0,99	0,97	0,96	0,95	0,94	0,94	<b>0,44</b>	<b>0,32</b>	0,39
RF	1,27	1,21	1,16	1,10	<b>1,02</b>	<b>0,96</b>	<b>0,93</b>	<b>0,90</b>	<b>0,90</b>	<b>0,90</b>	<b>0,90</b>	<b>0,90</b>	0,74	0,48	<b>0,36</b>

  

<b>Panel (b) : MAE ratio</b>															
Forecasting horizon															
Model	1	2	3	4	5	6	7	8	9	10	11	12	3m	6m	12m
AR	<b>1,02</b>	<b>1,03</b>	<b>1,05</b>	<b>1,06</b>	1,08	1,07	1,06	1,03	1,02	1,01	0,99	0,98	<b>0,45</b>	<b>0,36</b>	0,40
RF	1,19	1,17	1,15	1,11	<b>1,03</b>	<b>0,98</b>	<b>0,94</b>	<b>0,95</b>	<b>0,96</b>	<b>0,99</b>	<b>0,98</b>	<b>0,98</b>	0,73	0,51	<b>0,36</b>

NOTE : The table reports, for each forecasting horizon (1 to 12 months and accumulated 3, 6 and 12 months), the root mean squared error (RMSE) and mean absolute error (MAE) ratios with respect to the random walk model for the full out-of-sample period (2000–2015). The statistics for the best-performing model are highlighted in bold.

forecasting of RF model compared to benchmark.

The replication results also show, as the original article does, that among shrinkage models, sparsity-inducing methods are slightly worse than nonsparsity-inducing ones. Since RF does not impose sparsity, this might suggest that sparsity is not a desired feature to improve forecasting accuracy. On the contrary, quite different from article results, factor models do not seem to be strongly outperformed by other methods. The adoption of target factors improves the quality of the forecasts and make the results quite close to the RF ones, even if less accurate.

Table 2 shows more in detail the results of the comparison between RF and the AR benchmark, according to RMSE and MAE ratios with respect to the RW alternative for all forecasting horizons. We see that RF is beating AR and RW benchmarks on all time horizons from  $t + 1$  to  $t + 12$ , according to both RMSE and MAE; but the AR benchmark is performing better on accumulated inflation forecasts. This better performance of AR on accumulated inflation explains why on average, AR accuracy is close to (but lower than) RF one.

Extension to forecasting French inflation leads to results summarized in Tables 3 and 4. Table 3 reports summary statistics across all forecasting horizons, on the period 2005 to 2020. It shows that, if ML models allow to improve the quality of inflation forecasts over RW benchmark on average,

their performance is not systematically better as it was the case with US data. Particularly, the accuracy improvement does not seem to be constant over horizons : all ML models perform worse than RW benchmark at least once (they have a maximum RMSE and maximum MAE greater than one). We can also notice that the AR benchmark seems to perform better than several ML models. Contrary to what we found with US data, the supremacy of RF model over other ML alternatives is not very clear : it has the lower average and maximum RMSE and MAE, but is outperformed by LASSO and ElNet concerning minimum RMSE and MAE, and has a very low p-value in MCSs. Table 4 shows more in detail the results of the comparison between RF and the AR benchmark : as opposed to extension results, we see that AR is more accurate than RF not only for the accumulated inflation forecasts, but also for the horizons  $t + 1$  to  $t + 4$ . It is thus clear that AR benchmark is not outperformed by RF model as much as it is with US data. Most importantly, the factor model is performing way better than it was with US data, and is outperforming all benchmarks and ML alternatives : it beats the RW benchmark in terms on RMSE and MAE at all horizons, it has the lowest average and max RMSE and MAE, it is at most time horizons the model with lower RMSE and MAE, and it has the highest p-value in MCSs.

Our extension interestingly develops the results of original article and replication. It shows that, even though ML models can beat RW benchmark, factor model performs the best when the number of variables available is reduced. It thus confirms that improving usual forecasts cannot be done just by replacing traditional models with ML ones ; an other key feature is so-called "Big Data". Applying ML models to a reduced number of variables, as we did with French data, does not allow them to forecast inflation more accurately than AR benchmark. In this context of no "Big Data", the Factor model seems to perform the best.

This conclusion is however a little vague, and we do not precisely give a definition for "Big Data". In the article, 122 variables were used to forecast US inflation, leading to supremacy of RF model. In our extension, with use 62 variables and Factor model remains better than ML models. Studying what number of variables is needed for ML models to bring an improvement would be interesting. It would of course also depend on what variables are used.

The extension results, as the original article and replication results, also show that, among shrinkage models, sparsity-inducing methods (LASSO, adaLASSO, ElNet, adaELNet) are slightly worse than non-sparsity-inducing ones (RR). This confirms that the higher the number of variables, the



better the results of ML models.

The high performance of factor model in forecasting of French inflation is consistent with its high popularity for macroeconomic forecasting in a context of no Big Data : when there are not a lot a regressors available (as it is the case in our extension and more generally in traditionnal macroeconomic forecasting), the high level of aggregation of factor models seems to be adequate, but is not anymore when the number of regressors increase (as it is the case in the original article and, to a lesser extent, in our replication). The RF model becomes more performant only when the number of regressors increase. This is consistent with a conclusion of the original article, that the superiority of RF is due both to the variable selection mechanism induced by the method and to the presence of nonlinearities in the relation between inflation and its predictors. When more variables are available, the variable selection mechanism can only be more relevant. Moreover, the higher the number of variables, the higher the probability to have non-linear relationships between regressors and inflation.

These results lead us to re-evaluate the authors' stance that "*one can easily apply ML methods to forecast other macroeconomic series in a variety of countries.*" For ML methods to be as relevant for forecasting other countries macroeconomic series as they are for the US, these countries need to have a high data quantity available over a large time span, which is not necessarily the case, as we have seen with France.