

ERASMUS UNIVERSITY ROTTERDAM

Erasmus School of Economics

Bachelor Thesis Econometrie & Operationele Research  
Testing the Efficiency of Bias Predictions in Forecasting

Supervisor: Oorschot, J.A.

Second Assessor: Franses, P.H.B.F.

Ewout van Raad  
448569er@eur.nl

July 7, 2019

**Abstract**

Combining forecasts from multiple models into one forecast typically leads to better forecasting performance than individual models can achieve. However, elaborate backward looking combination strategies are often outperformed by simple strategies, such as averaging. Gibbs and Vasnev (2018) show that a forward looking approach based on predictable bias can make combinations that outperform equal weighting, individual and random walk forecasts. In this thesis, those results are generally replicated. The extension of their analysis is within the modification of bias predictions to have autoregressive terms to account for possible serial correlation. Forecast combinations have been made with these new bias predictions for both inflation and unemployment rate data and the results show that forecasting inflation generally benefits from adding autoregressive terms to the bias predictions, whereas unemployment forecasting performance does not improve.

The views stated in this thesis are those of the author and not necessarily those of Erasmus School of Economics or Erasmus University Rotterdam.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Data</b>	<b>4</b>
<b>3</b>	<b>Methodology</b>	<b>5</b>
3.1	Models . . . . .	5
3.2	Forecasting . . . . .	5
3.3	Combined Forecasting . . . . .	7
3.3.1	Conditionally Optimal Weights . . . . .	7
3.3.2	Combination Theory . . . . .	7
3.3.3	Combination Strategies . . . . .	9
3.3.4	Bias Prediction . . . . .	9
<b>4</b>	<b>Results</b>	<b>11</b>
<b>5</b>	<b>Conclusion</b>	<b>19</b>
<b>6</b>	<b>Limitations and Discussion</b>	<b>19</b>
<b>A</b>	<b>List of programs written and used</b>	<b>21</b>

# 1 Introduction

In many forecast combination applications it is customary to consider approaches that look at the historical performance of individual forecasts in order to construct optimal weights. In empirical studies it is often found that combining forecasts produces better forecasts than forecasts that originate from the superior individual model (Timmermann, 2006). Clemen (1989) finds similar results and stresses the interdisciplinary usefulness of combined forecasting. Next to econometrics, it also has potential for sciences as meteorology and psychology. The quote: 'We no longer need to justify this methodology,' referring to combined forecasting, illustrates the relevance of this topic.

Typically, simple forecast combination strategies such as equal weighting produce the most accurate forecasts when compared to combined forecasts that have a more elaborate way of weight determination. This empirical result seems to be contradictory to the theory, as the averaging strategy should only be optimal under a set of very restrictive conditions (Stock and Watson, 2004). Recently, one of the underlying reasons for this phenomenon is tackled within the literature. Gibbs and Vasnev (2018) show that an important reason for this difference between practice and theory lies within the forecast errors. Generally, forecast errors are predictable and they exhibit serial correlation. They prove that optimal weights conditional to the available information at that time yield better forecast results than unconditional weighting under a general loss function. Therefore, weight determination should be based on expected forecast performance, instead of looking at past performance. Gibbs and Vasnev (2018) also show empirically that forward looking approaches are better at forecasting inflation than backward looking ones. However, their focus was predominantly on the predictability of the forecast errors. The prediction of the forecast errors is given by a direct forecasting procedure of the model's real-time forecast error. To this end, they use an explicit regression for a given model  $i$  on the four-step ahead forecast error  $e_{i,t+4} = \pi_{t+4} - \mathbf{E}(\pi_{t+4})$ :

$$e_{i,t+4} = c + \beta_i x_t + \xi_{i,t+4},$$

in which  $\pi_{t+4}$  is the inflation rate at time  $t + 4$ ,  $\mathbf{E}(\pi_{t+4})$  is the four-quarter-ahead forecast of inflation at time  $t$  and  $x_t$  is a macroeconomic variable. Within this specification, the typically present serial correlation of the errors is not captured in any way. This could possibly be done by specifying an ARMA structure.

Moreover, only inflation forecasts are considered. Other variables may also be interesting to forecast, with possibly new model specifications based on macroeconomic relations. In particular, unemployment rates may be interesting to look at. Gibbs and Vasnev (2018) use a Phillips Curve Vector Autoregressive (VAR) model to forecast inflation, but they ignore the opportunity to forecast unemployment as well. New model specifications could be considered based on economic theory, such as a VAR model with unemployment and GDP growth (Okun's law, 1962). Incorporating new data and models is useful to create a higher degree of generality of the results found by Gibbs and Vasnev (2018), or possibly the new results find that their findings do not hold generally. The main research question of the thesis would be: Does the prediction of forecast errors improve when accounting for serial correlation using an ARMA structure and therefore lead to better forecast performance? To answer this question, it is useful to answer some subquestions, being:

1. Which ARMA structure would be most suitable?
2. Is the result replicable with another data type, specifically unemployment rates?

The following Data section describes the data used within this thesis and how it is used. The Methodology section elaborates on the forecasting models used and how the forecasting is performed. Next to that, it provides some theory about forecast combination, explains the procedure of predicting bias and describes the methods I use in this thesis to make the combinations. Logically the Results section follows, in which

I show that adding autoregressive terms is useful for inflation forecasting, but not for unemployment rate forecasting. After that section I provide the conclusion of this thesis and discuss further research possibilities and limitations of this paper.

## 2 Data

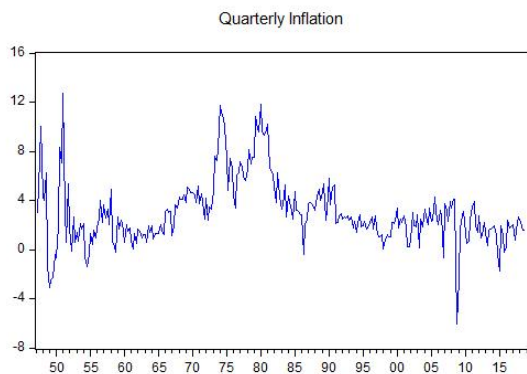
The data that I consider for this thesis is the same real-time data set that Gibbs and Vasnev (2018) use in their research. It is the Philadelphia Federal Reserve Real-Time Macroeconomic Data Set of which a couple of data types are extracted. First of all, we extract the quarterly Price Consumer Expenditure (PCE) index. This PCE is used to create the inflation variable as

$$\pi_t = 400 \ln \left( \frac{p_t}{p_{t-1}} \right),$$

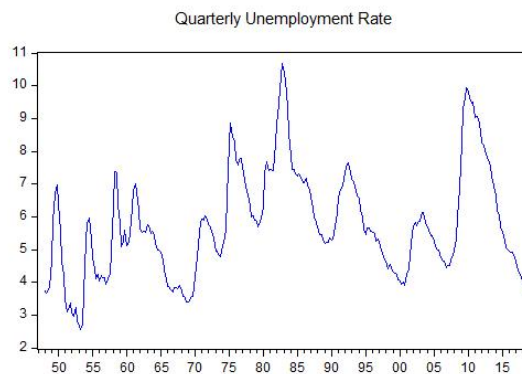
in which  $p_t$  is the PCE index. The macroeconomic variables that we consider for forecasting inflation are constructed from real GDP and unemployment rates found in the data set. We use the GDP level data to construct three new measures. The first one is the simple GDP growth, expressed in the same way as the inflation rate above. The second measure is the output gap, which we create by applying a Hodrick-Prescott filter to the GDP level data and taking the difference between the actual GDP level and the potential output, which is the result after the filtering. Finally, we consider a growth gap measure which is the difference between the current GDP growth and the highest GDP growth of the last twelve quarters.

Next to GDP, we consider unemployment rates to create two measures. The first one is simply the level of the unemployment rates and the second one is an unemployment gap, computed by taking the difference between the current unemployment rate and the lowest unemployment rate of the previous 12 quarters. The growth and unemployment gap are considered to account for possible nonlinearity of the Phillips curve (Stock and Watson, 2010).

I perform the forecasting, which I describe within the Methodology section, using the data vintage that has just the right amount of information. For example, we consider a four-quarter-ahead forecast of 1970Q1. The data vintage containing 1969Q1 as its last observation will be used to construct this forecast. It occurs that some last observations are not available yet, in which case a four-quarter-ahead forecast cannot be done with that vintage since it then has to be a five-quarter-ahead forecast for the required period of time. If this is the case, the inflation of the missing observation will be set equal to the previous period.



**Figure 1:** Inflation data



**Figure 2:** Unemployment rates

### 3 Methodology

#### 3.1 Models

In order to replicate the results by Gibbs and Vasnev (2018), we consider the same individual models as they used in their analysis when considering inflation data. This comes down to three types of models. The first group is a set of univariate autoregressive models. The AR(1), AR(2), AR(4), ARMA(1,1) and ARMA(4,4) models are considered as benchmark models. The general ARMA( $p, q$ ) model is given by

$$\pi_t = c + \sum_{i=1}^p \beta_i \pi_{t-i} + \sum_{j=1}^q \phi_j \epsilon_{t-j} + \epsilon_t,$$

The general AR( $p$ ) model is a restricted ARMA( $p, q$ ) model:

$$\pi_t = c + \sum_{i=1}^p \beta_i \pi_{t-i} + \epsilon_t,$$

with  $\phi_j = 0, \forall j$ . All benchmark models can be obtained by substituting the relevant values for  $p$  and  $q$ . Next to those, we consider a random walk model introduced by Atkeson et al. (2001), which is simply the average of the four preceding quarters of inflation:

$$\hat{\pi}_{t+h}^{AO} = \frac{1}{4} \sum_{i=1}^4 \pi_{t-i}.$$

For the purpose of forecasting unemployment rates, we consider these models as well.

The Phillips Curve type models are also used in the same way as Gibbs and Vasnev (2018), namely as bivariate VAR models with two lags for inflation and two lags of a macroeconomic variable, including a VAR All model. When forecasting unemployment, these models are used as well since the typical Phillips Curve relation is that between inflation and unemployment. The bivariate VAR with unemployment and GDP growth is interesting, since their relation is known as Okun's law (1962). The VAR(2) specification is given by

$$y_t = c + A_1 y_{t-1} + A_2 y_{t-2} + \epsilon_t,$$

where  $y_t = (\pi_t, x_t)'$ . Finally direct forecasts are considered for both data types, which regress four-quarter-ahead inflation (unemployment rate) on a macroeconomic variable:

$$\pi_{t+4} = c + \beta x_t + \epsilon_{t+4},$$

with  $\pi_t$  the inflation at time  $t$  and  $x_t$  a macroeconomic variable.

#### 3.2 Forecasting

In this thesis, we consider the four-quarter-ahead forecast of quarterly inflation and unemployment following Gibbs and Vasnev (2018) with inflation following an annual rate. Forecasts are made based on the latest vintage of data available throughout time as explained thoroughly in the Data section. Due to missing the current observation in each vintage, the four-quarter-ahead forecast resembles the forecast of the current and following three quarters. We denote the forecast as  $\mathbf{E}_t^\tau \pi_{t+4}$ , where  $\tau$  is the vintage and  $t$  is the last observation available.

The four-quarter-ahead forecasts from the benchmark ARMA models and the VAR models are obtained

by iteration (dynamic forecasting), in contrast with the direct forecasting models.

The metrics and target measures to evaluate the forecasts follow those of Gibbs and Vasnev (2018), using Root Mean Squared Forecast Error (RMSFE) to measure relative forecasting performance to the benchmark AO forecasts and Mean Forecast Error (MFE) to compute the bias. These statistics are given as follows:

$$RMSFE = \sqrt{\sum_{t=1}^T \left( \frac{e_{i,t}^2}{T} \right)},$$

$$MFE = \sum_{t=1}^T \left( \frac{e_{i,t}}{T} \right),$$

with  $e_{i,t}$  the forecast error at time  $t$  for model  $i$ . We consider the test developed by Diebold and Mariano (1995) to check whether RMSFEs are significantly different. In order to compute the Diebold-Mariano (DM) statistic, it is necessary to construct loss differentials based on the squared forecast errors of a specific model and the random walk AO forecasts,

$$D_t = e_{i,t}^2 - e_{AO,t}^2.$$

The test is based on the arithmetic mean of the loss differential

$$\bar{D} = \frac{\sum_{t=1}^T D_t}{T}.$$

It is likely that the loss differential series is autocorrelated (Harvey et al., 1997). They state that the asymptotic variance of  $\bar{D}$  can be shown to be

$$V(\bar{D}) \approx \frac{\gamma_0 + 2 \sum_{k=1}^{h-1} \gamma_k}{T},$$

with  $\gamma_k$  the  $k$ -th order autocovariance of  $D_t$  and  $h = 4$  in the case of four-quarter-ahead forecasting. The autocovariance can be estimated in the following manner:

$$\hat{\gamma}_k = \sum_{t=k+1}^T (D_t - \bar{D})(D_{t-k} - \bar{D}).$$

The Diebold-Mariano statistic is given as follows:

$$DM = \frac{\bar{D}}{\sqrt{\hat{V}(\bar{D})}},$$

which follows a standard normal distribution under the null hypothesis.

Harvey et al. (1997) state that the DM test is found to be over-sized for a sample with a moderate number of observations. They modify the DM statistic for  $h$ -step-ahead forecasts to account for this possible problem:

$$DM^* = \left( \frac{T + 1 + 2h + T^{-1}h(h-1)}{T} \right) DM.$$

Harvey et al. (1997) argue that comparing the DM statistic to quantile values of the Student's  $t$ -distribution with  $(n - 1)$  degrees of freedom is more appropriate than comparing it with the standard normal extreme values. Hence, in this thesis I compare the DM statistic with the  $t(n - 1)$  extreme values at one, five and ten percent significance level.

To test the significance of the bias I consider a  $t$ -test using HAC standard errors (Newey and West, 1986). To perform this test, I regress the series of the forecast errors on only a constant with the Newey-West standard error specification. We can simply compute the test statistic as  $\frac{\text{constant}}{\text{stderror}_{NW}}$  which follows a  $t(n-2)$ -distribution.

Considering the difference between the metrics and measures used by Gibbs and Vasnev (2018) and this paper, the exceptions are:

1. The target measures of unemployment, since there is no second release data available in the real-time data set. We compare it to the most recent data available due to the lack of this data.
2. The target measure for inflation is the second-release data of the GDP deflator inflation type. The PCE and GDP deflator are similar and for the GDP deflator the second-release data is available and verified to be legitimate within the real-time data set, whereas this is not the case for PCE inflation.

### 3.3 Combined Forecasting

#### 3.3.1 Conditionally Optimal Weights

When forecast combinations are considered, minimizing a loss function is typically done to determine the weights of the forecasts. Rather than just using the unconditional variance for this purpose, we consider a conditional Mean Squared Error (MSE) loss function that can be decomposed as the sum of squared bias and variance. The measure is conditional on predictable information not incorporated by classical weighting techniques. The implication that such predictable information exists follows from the notion that economic forecasting models frequently suffer from misspecification and that economic data types usually exhibit regular structural breaks. If this notion is true, Hendry and Clements (2004) show that equally weighted forecasting is effective, since the biases from the different models are leveled out against each other with this approach. Hence, an equally weighted forecast is typically hard to beat with common backward looking combination techniques. However, the structural breaks and misspecified forecasting models due to time-varying aspects imply that the forecast errors are serially correlated and those issues lead to the idea that there is some predictable information which is not considered by backward looking combination approaches. Hence, the forecast combination weights should be conditioned on that predictable information. This implies the minimization of the conditional MSE mentioned earlier, preferred over minimizing a measure like unconditional variance (Gibbs and Vasnev, 2018).

#### 3.3.2 Combination Theory

Combining forecasts from different models to minimize a certain expected loss function such as (Root) Mean Squared Error may yield better results than the individual models and simple combination strategies. Consider the four-quarter-ahead forecast  $y_{T+4}$ , information set  $I_T$  and a vector of length  $n$  consisting of the forecasts from  $n$  different models

$$\mathbf{f}_{T+4} = (f_{1,T+4}, f_{2,T+4}, \dots, f_{n,T+4})'$$

which we transform into one combined forecast using a linear transformation with weights

$$\mathbf{w} = (w_1, w_2, \dots, w_n)'$$

to construct a forecast as  $f_{c,T+4} = \mathbf{w}'\mathbf{f}_{T+4}$ , following Gibbs and Vasnev (2018). The forecast errors are given as

$$\mathbf{e}_{T+4} = y_{T+4}\mathbf{1} - \mathbf{f}_{T+4},$$

leading to a combined error  $e_{c,T+4} = \mathbf{w}'\mathbf{e}_{T+4}$ . Consider the Mean Squared Error loss function  $L(e) = e^2$ . Note that this implicitly assumes that the loss function only depends on the forecast error  $e_{c,T+4}$ . The conditionally optimal weights are the solution to the minimization problem of the MSE:

$$\mathbf{w}^*(I_T) = \underset{\mathbf{w}}{\operatorname{argmin}} \mathbf{E}(L(e_{c,T+4}|I_T)).$$

We assume that the forecasts are unbiased and minimize the MSE subject to the constraint that the weights sum up to one. We can decompose the errors of the original forecasts following Gibbs and Vasnev (2018) as

$$\mathbf{e}_{T+4} = \mathbf{b}_T + \boldsymbol{\xi}_{T+4},$$

with  $\mathbf{b}_T = \mathbf{E}(\mathbf{e}_{T+4}|I_T)$  and  $\mathbf{E}(\boldsymbol{\xi}_{T+4}|I_T) = 0$ . We can find an expression for the MSE:

$$MSE(\mathbf{w}) = \mathbf{w}'(\Sigma_\xi + \mathbf{b}_T\mathbf{b}_T')\mathbf{w},$$

where  $\Sigma_\xi = \mathbf{E}(\boldsymbol{\xi}_{T+4}\boldsymbol{\xi}_{T+4}'|I_T)$ . The MSE is minimized by the conditionally optimal weights

$$\mathbf{w}^*(I_T) = \frac{[\Sigma_\xi + \mathbf{b}_T\mathbf{b}_T']^{-1}\boldsymbol{\iota}}{\boldsymbol{\iota}'[\Sigma_\xi + \mathbf{b}_T\mathbf{b}_T']^{-1}\boldsymbol{\iota}}.$$

It is clear that not only the optimal weights, but also  $\mathbf{b}_t, \Sigma_\xi$  and  $MSE(\mathbf{w})$  depend on  $I_T$ . However, we choose to keep notation simple and reserve this dependency for the optimal weights, to emphasize the time-varying nature of the weights and to stay close to the assumptions made by Gibbs and Vasnev (2018).  $\Sigma_\xi$  is assumed to be constant in this thesis as well. The conditionally optimal weights are of course different than unconditionally optimal weights

$$\mathbf{w}^* = \frac{\Sigma_e^{-1}\boldsymbol{\iota}}{\boldsymbol{\iota}'\Sigma_e^{-1}\boldsymbol{\iota}},$$

with unconditional variance of errors  $\Sigma_e = \Sigma_\xi + \mathbf{E}(\mathbf{b}_t\mathbf{b}_t')$ . For further intuition on the dynamics of these 'new' conditionally optimal weights I refer to the observations made by Gibbs and Vasnev (2018). Finally, they propose a three-part theorem that gives a formalization of the idea that considering more information will lead to better forecast combinations.

**Theorem 1** Given the existence of the first and second (un)conditional moments, convex loss function  $L(\cdot)$  and information sets  $J_T \subset I_T$ , then:

- (a)  $\mathbf{E}(MSE(\mathbf{w}^*(I_T))) \leq \mathbf{E}(MSE(\mathbf{w}^*)),$
- (b)  $\mathbf{E}(MSE(\mathbf{w}^*(I_T)|J_T)) \leq (MSE(\mathbf{w}^*(J_T))),$
- (c)  $\mathbf{E}\left(\underset{\mathbf{w}}{\operatorname{min}} \mathbf{E}(L(e_{c,T+4}|I_T)|J_T)\right) \leq \underset{\mathbf{w}}{\operatorname{min}} \mathbf{E}(L(e_{c,T+4}|J_T)).$

Part (a) is a comparison between conditional and unconditional MSE. The theorem implies that adding more information leads to better or at worst equally good forecasting combinations. Part (b) is more or less similar to part (a) since it is a generalization using information sets, again implying that adding more information ( $I_T$  contains more information than  $J_T$ ) leads to better or equally good combinations in terms of MSE. Part (c) implies that conditioning on more information also leads to a lower or equal (unspecified but convex) loss function outcome, which is preferable of course. This concludes the theory on forecast combinations and now let us look at the different strategies that I employ in this thesis.



### 3.3.3 Combination Strategies

We have seen that an optimal solution for the MSE exists:

$$\mathbf{w}^*(I_T) = \frac{[\Sigma_\xi + \mathbf{b}_T \mathbf{b}'_T]^{-1} \boldsymbol{\iota}}{\boldsymbol{\iota}' [\Sigma_\xi + \mathbf{b}_T \mathbf{b}'_T]^{-1} \boldsymbol{\iota}}.$$

However, estimation of  $\Sigma_\xi$  and  $\mathbf{b}_T$  inherits some issues. We consider three options to try to lessen these issues, following Gibbs and Vasnev (2018). The first one is the equal weights benchmark, since it is known to be an effective method (Hendry and Clements, 2004). The second approach is the shrinkage method. The idea of the shrinkage method is to stabilize the variance estimation, which is typically unstable. It is given explicitly by

$$\tilde{\Sigma}_\xi = \alpha \Sigma_0 + (1 - \alpha) \hat{\Sigma}_\xi,$$

with shrinkage parameter  $\alpha = 0.5$  and stabilizing matrix  $\Sigma_0$ , equal to identity matrix  $\mathbf{I}$ . The choice of the identity matrix as a stabilizer shrinks the weights towards equal weights (Gibbs and Vasnev, 2018). The optimal weights are given by

$$\hat{\mathbf{w}}^{COS}(I_T) = \frac{[\tilde{\Sigma}_\xi + \hat{\mathbf{b}}_T \hat{\mathbf{b}}'_T]^{-1} \boldsymbol{\iota}}{\boldsymbol{\iota}' [\tilde{\Sigma}_\xi + \hat{\mathbf{b}}_T \hat{\mathbf{b}}'_T]^{-1} \boldsymbol{\iota}},$$

with  $\hat{\mathbf{b}}_T$  the estimated bias and  $\boldsymbol{\iota}$  a vector of ones. Thirdly, we consider the predicted exponential weights method:

$$\hat{\mathbf{w}}^{PE}(I_T) = \frac{1}{\sum_{i=1}^n \exp(-\gamma \hat{b}_{i,T}^2)} (\exp(-\gamma \hat{b}_{1,T}^2), \dots, \exp(-\gamma \hat{b}_{n,T}^2))',$$

which accelerates the weight decrease when bias increases. In the research,  $\gamma$  is either set equal to five or sent towards infinity. This method is quite intuitive, as it 'punishes' increasing bias with decreasing weights at an exponential rate. When bias is close to or equal to zero, it yields a relatively high weight for that specific model. Looking at the limiting case where  $\gamma \rightarrow \infty$ , this model allocates a weight equal to one to the model on the model that predicts the lowest (squared) bias. This case is particularly interesting to examine, since this type of forecasting with a single model at each point in time usually performs very poorly in out-of-sample forecasting (Timmermann, 2006).

### 3.3.4 Bias Prediction

As mentioned in the introduction, we obtain the prediction of the forecast error by direct forecasting

$$e_{i,t+4} = c + \beta_i x_t + \xi_{i,t+4},$$

in which  $\pi_{t+4}$  is the inflation rate at time  $t + 4$ ,  $\mathbf{E}(\pi_{t+4})$  is the four-quarter-ahead forecast of inflation at time  $t$  and  $x_t$  is a macroeconomic variable. We use the specification for replication and extend it to have a ARX(1) and ARMAX(1,1) specification, meaning models with an ARMA structure and an explanatory variable as well. Specifically, the expressions are

$$\begin{aligned} e_{i,t+4} &= c + \beta_i x_t + \phi_i e_{i,t+4} + \xi_{i,t+4}, \\ e_{i,t+4} &= c + \beta_i x_t + \phi_i e_{i,t+3} + \psi_i \xi_{i,t+3} + \xi_{i,t+4}. \end{aligned}$$

We consider these specifications since we already use these AR(1) and ARMA(1,1) models as benchmarks for the individual forecasting performance and therefore it seems reasonable to consider these ARMA specifications for the error prediction as well. To keep the research executable within the time frame,

only these two specifications are employed.

To check which specification is best, I consider a heuristic approach based on the Akaike Information Criterion. This is necessary as I recursively perform bias prediction regressions with varying regressors for 17 different models. This procedure is done for the 'vanilla' bias prediction, the ARX(1) and AR-MAX(1,1) prediction. Since this produces a high number of AICs, I consider an indicator function  $I_t$  with the following specification:

$$I_t = 1 \text{ if } AIC_t^{AR1} > AIC_t^{ARMA(1,1)},$$
$$I_t = 0 \text{ else.}$$

Intuitively, if the sum of the indicators is sufficiently high, this would mean that the ARMA(1,1) model is preferred above the AR(1) model for a specific bias prediction. This is of course a vague criterion, so I propose to use a  $t$ -test for the indicators to test whether the mean is significantly different from 0.5. The mean of the indicator represents a probability that the ARMA model is better than the AR model, so if they are equally accurate the mean equals 0.5. The procedure is to regress the indicator on a constant with Newey-West standard errors and subsequently performing the test.

The resulting model would be a Linear Probability Model (LPM) which suffers from heteroskedasticity by definition (Horrace and Oaxaca, 2006), but the Newey-West errors fix this issue. Also, the fact that this LPM generally yields biased parameters is not a problem since only a constant is involved, which equals the mean of the indicator. The model that is significantly better on the most occasions is the model I use for the actual forecasting combinations.

By performing this heuristic approach the usage of the AIC is justified in my opinion, as the direct AIC comparison on which the value of the indicator is based is between a model and an extension of that model with an MA(1) term. In this sense, the comparison is always between nested models. This testing procedure may not be perfect for this application, since the indicator value of the present AIC is highly correlated with the indicator value of the preceding AICs. However, the Newey-West errors and parameter estimates, which are simply the means of the indicators, are not restricted by this possible dependence in the data whereas a regular  $t$ -test for testing a mean requires independence. Therefore, I choose to use this method.

We construct the forecast error series in the same manner as Gibbs and Vasnev (2018), while unemployment forecast errors are made using the latest data available due to unavailability of second-release data. For out-of-sample forecasting, we follow the procedure of Section 3.5.2 by Gibbs and Vasnev (2018), with the in-sample forecasting period starting from 1966Q4 in order to prevent the usage of information that was not available at the time a forecast was made.

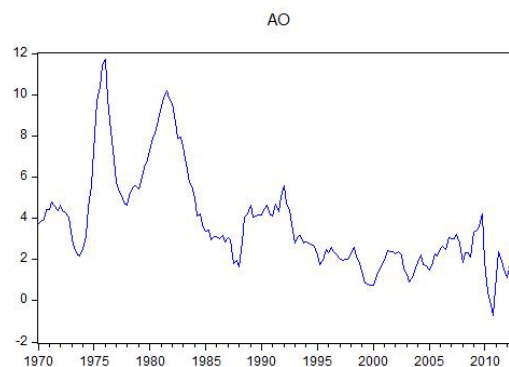
## 4 Results

**Table 1:** Individual model performance

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Benchmark</i>		
AO	2.346	-0.07 <sup>†</sup>
<i>Direct Forecasts</i>		
DF CUR	1.300	-0.25 <sup>†</sup>
DF GDP Growth	1.243	0.08 <sup>†</sup>
DF Growth Gap	1.179	0.52 <sup>†</sup>
DF Output Gap	1.303	0.20 <sup>†</sup>
DF Unemp. Gap	1.298	0.14 <sup>†</sup>
<i>VAR Forecasts</i>		
VAR CUR	1.030	-0.09 <sup>†</sup>
VAR GDP Growth	1.024	0.40 <sup>†</sup>
VAR Growth Gap	1.003	0.07 <sup>†</sup>
VAR Output Gap	1.032	0.34 <sup>†</sup>
VAR Unemp. Gap	0.988	0.21 <sup>†</sup>
VAR All	1.075	-0.16 <sup>†</sup>
<i>Benchmark Forecasts</i>		
AR(1)	1.110	0.29 <sup>†</sup>
AR(2)	0.993	0.21 <sup>†</sup>
AR(4)	1.056	0.36 <sup>†</sup>
ARMA(1,1)	0.970	0.20 <sup>†</sup>
ARMA(4,4)	1.097	0.27 <sup>†</sup>

Table 1: This table shows the performance of the individual models in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the AO benchmark. The Rel. RMSFE for AO is the real RMSFE value.

The results of Table 1 seem to be in line with the findings of Gibbs and Vasnev (2018). The direct forecasts are all substantially poor compared to the benchmark, with the Growth Gap specification being the best with a Relative RMSFE of 1.179. The VAR forecasts have similar performance to the AO benchmark, with the best model being the Unemployment Gap model with an RMSFE of 0.988, (insignificantly) outperforming AO. Some ARMA models perform slightly better as well, albeit not significantly. There is no individual model that can significantly outperform the random walk model which replicates the general results from Gibbs and Vasnev (2018). Note that in terms of exact numbers, the found values can be slightly different. This may be caused by the usage of most recent data for the regressors and possibly, a slightly different forecasting approach. Next to the RMSFE results, all forecasts errors are unbiased at a ten percent significance level, fitting the assumption that the expected value of the forecast errors is equal to zero.



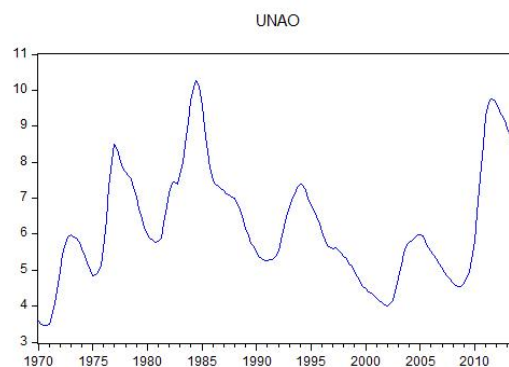
**Figure 3:** AO Random Walk forecasts inflation

**Table 2:** Individual model performance unemployment

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Benchmark</i>		
AO	1.512	0.19 <sup>†</sup>
<i>Direct Forecasts</i>		
DF Inflation	1.227	0.95
DF GDP Growth	1.284	1.03
DF Growth Gap	1.222	0.78
DF Output Gap	1.229	0.98
<i>VAR Forecasts</i>		
VAR Inflation	0.694***	0.47
VAR GDP Growth	0.682***	0.37
VAR Growth Gap	0.677***	0.48
VAR Output Gap	0.696***	0.32
VAR All	0.683***	0.21 <sup>†</sup>
<i>Benchmark Forecasts</i>		
AR(1)	0.738***	0.23 <sup>†</sup>
AR(2)	0.749***	0.50
AR(4)	0.734***	0.41
ARMA(1,1)	0.710***	0.29
ARMA(4,4)	0.736***	0.41

Table 2: This table shows the performance of the individual models in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the AO benchmark. The Rel. RMSFE for AO is the real RMSFE value.

The individual model results I obtain by using (quarterly) unemployment data have some distinct differences compared to the inflation forecast models. The AO random walk forecasts are significantly less accurate in terms of RMSFE than the autoregressive benchmark and VAR models. This result is not surprising since unemployment is known to be rigid, implying that autoregressive models are useful in explaining and forecasting unemployment. Likely due to the parameter restrictions on the AO model (all 0.25), the performance is quite poor. Although the VAR and ARMA models perform well when looking at relative RMSFE, the forecasts seem to be biased in most cases. The similarity with the inflation results lies within the direct forecasting models since the relative RMSFEs to the benchmark are approximately of the same size. The best individual model is the VAR model with the Growth Gap Measure as a regressor with a relative RMSFE of 0.677. I include this model as a benchmark next to the AO and equal weights forecast in order to check whether combined forecasts can outperform the best individual model.



**Figure 4:** AO Random Walk forecasts unemployment

**Table 3:** AR(1) vs. ARMA(1,1) contest inflation

Predictor	Regressors				
	Unemp. Gap <i>t</i> -test statistic	GDP Growth <i>t</i> -test statistic	Growth Gap <i>t</i> -test statistic	Unemp. <i>t</i> -test statistic	Output Gap <i>t</i> -test statistic
<i>Benchmark</i>					
AO	-25.13 <sup>†</sup>	-45.14 <sup>†</sup>	-27.22 <sup>†</sup>	-45.14 <sup>†</sup>	-22.14 <sup>†</sup>
<i>Direct Forecasts</i>					
DF CUR	2.86 <sup>†</sup>	1.81 <sup>†</sup>	1.71 <sup>†</sup>	1.62 <sup>†</sup>	2.15 <sup>†</sup>
DF GDP Growth	3.66 <sup>†</sup>	3.68 <sup>†</sup>	3.28 <sup>†</sup>	2.38 <sup>†</sup>	5.86 <sup>†</sup>
DF Growth Gap	2.78 <sup>†</sup>	2.71 <sup>†</sup>	2.49 <sup>†</sup>	2.16 <sup>†</sup>	4.36 <sup>†</sup>
DF Output Gap	1.80 <sup>†</sup>	1.64 <sup>†</sup>	1.47 <sup>†</sup>	1.53 <sup>†</sup>	2.46 <sup>†</sup>
DF Unemp. Gap	1.33 <sup>†</sup>	1.18	1.16	1.08	1.58 <sup>†</sup>
<i>VAR Forecasts</i>					
VAR CUR	1.25	0.39	0.84	0.39	1.84 <sup>†</sup>
VAR GDP Growth	0.47	1.02	0.89	1.17	1.69 <sup>†</sup>
VAR Growth Gap	-0.32	-0.52	-0.73	-0.25	0.33
VAR Output Gap	0.54	1.10	0.90	1.10	2.05 <sup>†</sup>
VAR Unemp. Gap	-1.48 <sup>†</sup>	-1.27	-1.72 <sup>†</sup>	-1.30 <sup>†</sup>	-0.32
VAR All	-9.67 <sup>†</sup>	-9.12 <sup>†</sup>	-13.96 <sup>†</sup>	-13.34 <sup>†</sup>	-8.47 <sup>†</sup>
<i>Benchmark Forecasts</i>					
AR(1)	1.88 <sup>†</sup>	1.80 <sup>†</sup>	1.83 <sup>†</sup>	1.62 <sup>†</sup>	2.54 <sup>†</sup>
AR(2)	1.29	1.27	1.21	1.17	1.92 <sup>†</sup>
AR(4)	1.80 <sup>†</sup>	1.81 <sup>†</sup>	1.71 <sup>†</sup>	1.61 <sup>†</sup>	2.54 <sup>†</sup>
ARMA(1,1)	1.41 <sup>†</sup>	1.27	1.21	1.17	1.92 <sup>†</sup>
ARMA(4,4)	2.65 <sup>†</sup>	2.60 <sup>†</sup>	2.50 <sup>†</sup>	2.16 <sup>†</sup>	2.61 <sup>†</sup>
<i>'Victories'</i>					
AR(1)	3	2	3	3	2
ARMA(1,1)	9	7	7	7	13

Table 3: A † represents significant difference at ten percent level. Significant negative *t*-statistics indicate a victory for AR(1) and significant positive *t*-statistics indicate an ARMA(1,1) victory. *t*-statistics are obtained by regressing the AIC indicators on a constant and using Newey-West standard errors.

Table 3 shows that the ARMA(1,1) bias predictions generally yield better AIC results than the AR(1) predictions for inflation data. The difference in number of victories is quite substantial since the AR(1) specification does not have more than three significant victories whereas the minimum number of victories from the ARMA(1,1) model equals seven for a given regressor.

**Table 4:** AR(1) vs. ARMA(1,1) contest unemployment

Predictor	Regressors			
	Inflation $t$ -test statistic	GDP Growth $t$ -test statistic	Growth Gap $t$ -test statistic	Output Gap $t$ -test statistic
<i>Benchmark</i>				
AO	-35.03 <sup>†</sup>	-20.85 <sup>†</sup>	-20.85 <sup>†</sup>	-13.03 <sup>†</sup>
<i>Direct Forecasts</i>				
DF Inflation	20.00 <sup>†</sup>	25.43 <sup>†</sup>	20.00 <sup>†</sup>	20.00 <sup>†</sup>
DF GDP Growth	-10.39 <sup>†</sup>	-44.88 <sup>†</sup>	-86.52 <sup>†</sup>	-45.14 <sup>†</sup>
DF Growth Gap	-6.99 <sup>†</sup>	-39.69 <sup>†</sup>	-45.14 <sup>†</sup>	-45.14 <sup>†</sup>
DF Output Gap	25.38 <sup>†</sup>	20.00 <sup>†</sup>	25.38 <sup>†</sup>	20.00 <sup>†</sup>
<i>VAR Forecasts</i>				
VAR Inflation	1.58 <sup>†</sup>	2.03 <sup>†</sup>	1.77 <sup>†</sup>	1.84 <sup>†</sup>
VAR GDP Growth	-10.74 <sup>†</sup>	-4.30 <sup>†</sup>	-4.79 <sup>†</sup>	-1.64 <sup>†</sup>
VAR Growth Gap	2.48 <sup>†</sup>	3.10 <sup>†</sup>	3.27 <sup>†</sup>	3.17 <sup>†</sup>
VAR Output Gap	-44.67 <sup>†</sup>	-15.20 <sup>†</sup>	-22.46 <sup>†</sup>	-20.00 <sup>†</sup>
VAR All	8.75 <sup>†</sup>	9.17 <sup>†</sup>	9.17 <sup>†</sup>	8.75 <sup>†</sup>
<i>Benchmark Forecasts</i>				
AR(1)	22.98 <sup>†</sup>	41.49 <sup>†</sup>	86.52 <sup>†</sup>	39.85 <sup>†</sup>
AR(2)	-24.92 <sup>†</sup>	-9.25 <sup>†</sup>	-9.25 <sup>†</sup>	-9.53 <sup>†</sup>
AR(4)	-22.98 <sup>†</sup>	-15.32 <sup>†</sup>	-15.32 <sup>†</sup>	-39.69 <sup>†</sup>
ARMA(1,1)	-1.86 <sup>†</sup>	3.41 <sup>†</sup>	3.09 <sup>†</sup>	3.25 <sup>†</sup>
ARMA(4,4)	-17.52 <sup>†</sup>	-18.14 <sup>†</sup>	-18.14 <sup>†</sup>	-22.98 <sup>†</sup>
<i>'Victories'</i>				
AR(1)	9	8	8	8
ARMA(1,1)	6	7	7	7

Table 4: A <sup>†</sup> represents significant difference at ten percent level. Significant negative  $t$ -statistics indicate a victory for AR(1) and significant positive  $t$ -statistics indicate an ARMA(1,1) victory.  $t$ -statistics are obtained by regressing the AIC indicators on a constant and using Newey-West standard errors.

The results from Table 4 show that the results for unemployment data differ from the inflation data results. The number of victories are closer to each other, but the AR(1) bias prediction is strictly superior for each regressor in terms of these victories. Hence, I use the ARMA(1,1) specification for bias prediction for inflation data whereas I use the AR(1) specification for unemployment data.

**Table 5:** Combined forecasting performance inflation

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Shrinkage Weights</i> ( $\alpha = 0.5$ )		
CUR	0.784***	-0.002 <sup>†</sup>
GDP Growth	0.964	0.41
Growth Gap	0.927	0.40
Output Gap	0.873***	0.21 <sup>†</sup>
Unemp. Gap	0.903**	0.40
<i>Predicted Exponential Weights</i> ( $\gamma = 5$ )		
CUR	0.879***	0.09 <sup>†</sup>
GDP Growth	0.987	0.28 <sup>†</sup>
Growth Gap	0.971	0.42
Output Gap	0.959	0.33 <sup>†</sup>
Unemp. Gap	0.909**	-0.31 <sup>†</sup>
<i>Predicted Exponential Weights</i> ( $\gamma \rightarrow \infty$ )		
CUR	0.906**	0.06 <sup>†</sup>
GDP Growth	1.031	0.20 <sup>†</sup>
Growth Gap	1.006	0.34 <sup>†</sup>
Output Gap	1.019	0.40
Unemp. Gap	0.928*	0.34 <sup>†</sup>
<i>Benchmark Forecasts</i>		
Equal weights	2.440	0.19 <sup>†</sup>
AO	0.962	-0.07 <sup>†</sup>

Table 5: This table shows the performance of the combined forecasts in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \* $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the Equal Weights benchmark. The Rel. RMSFE for equal weights is the real RMSFE value.

The results from Table 5 indicate that using a forward looking approach when it comes to combining forecasts is useful. The shrinkage method yields significant improvements over the equal weights benchmark for three regressors, with the best combined forecasts coming from the bias prediction using unemployment rate. Predicted exponential weighting for both  $\gamma$  values also leads to two significant improvements. However, sending  $\gamma$  towards infinity leads to three forecast series that are worse than the equal weights benchmark considering RMSFE and bias. Timmermann (2006) notes that using a single model for each period of time generally leads to poor out-of-sample forecasting results, so in that sense these results are not out of line with the theory. The outcome is also still in line with Gibbs and Vasnev (2018) since bias predictions can still lead to RMSFE improvement for two forecast combinations with this method, even when this type of forecasts would typically yield inferior results.

The bias predictions using the unemployment variables are important for all methods, since they (more or less) yield the best results considering relative RMSFE. The forecasts are generally unbiased with a few exceptions. Since all forecasts were unbiased at a ten percent level when looking at the individual model forecasts, combining forecasts does not necessarily improve bias results.

These results again generally follow Gibbs and Vasnev (2018), since the idea that a forward looking approach should work well in combined forecasting is confirmed by this empirical analysis for inflation data. There are some differences considering which regressors work best for each method and which regressors for bias prediction yield the significant improvements. However, for reasons mentioned earlier at the individual model results, this is not surprising. Moreover, it seems more intuitive that unemployment variables are important for explaining and forecasting inflation, since it is established economic theory that inflation and unemployment have a distinct relationship (Phillips, 1958).

**Table 6:** Combined forecasting performance ARMA(1,1)

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Shrinkage Weights</i> ( $\alpha = 0.5$ )		
CUR	0.875***	-0.14 <sup>†</sup>
GDP Growth	0.885**	-0.02 <sup>†</sup>
Growth Gap	0.886**	-0.01 <sup>†</sup>
Output Gap	0.876***	0.03 <sup>†</sup>
Unemp. Gap	0.817***	-0.10 <sup>†</sup>
<i>Predicted Exponential Weights</i> ( $\gamma = 5$ )		
CUR	0.925***	-0.14 <sup>†</sup>
GDP Growth	0.894**	-0.04 <sup>†</sup>
Growth Gap	0.895**	-0.06 <sup>†</sup>
Output Gap	0.921***	-0.05 <sup>†</sup>
Unemp. Gap	0.913***	-0.05 <sup>†</sup>
<i>Predicted Exponential Weights</i> ( $\gamma \rightarrow \infty$ )		
CUR	0.948	-0.17 <sup>†</sup>
GDP Growth	0.904**	-0.07 <sup>†</sup>
Growth Gap	0.926*	-0.10 <sup>†</sup>
Output Gap	0.966	-0.08 <sup>†</sup>
Unemp. Gap	0.941	-0.09 <sup>†</sup>
<i>Benchmark Forecasts</i>		
Equal weights	2.440	0.19 <sup>†</sup>
AO	0.962	-0.07 <sup>†</sup>

Table 6: This table shows the performance of the combined forecasts in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \* $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the Equal Weights benchmark. The Rel. RMSFE for equal weights is the real RMSFE value.

The contents of Table 6 show that adding an ARMA(1,1) specification to the bias predictions leads to superior results in comparison to the 'regular' bias predictions which only use a macro-economic variable and a constant. The first interesting fact to note is that generally the relative RMSFE values move together more closely. A logical result, since the bias predictions now all share two explanatory variables.

The shrinkage and exponential weighting ( $\gamma = 5$ ) methods now yield significant improvements over equal weights for all bias predictors. Moreover, all forecasts are unbiased, which also implies an improvement in bias performance from the combined forecast results from Table 5. The limiting case of the exponential weights method ( $\gamma = \infty$ ) also leads to better results in the sense that now, no forecast series is less accurate than the equal weights benchmark considering both RMSFE and bias. However, the modified bias predictions do not lead to more significant improvements for this method.

Since forecast errors are likely serially correlated (Gibbs and Vasnev, 2018), the result that explicitly modelling autoregressive terms to account for this possible serial correlation actually improves forecasting performance is somewhat expected.

Despite the general improvement, the best forecast series for shrinkage and exponential weighting ( $\gamma = 5$ ) from the regular combined forecasts now yield less accurate results. For these methods the unemployment rate predictor leads to the most accurate results in Table 5, which is no longer the case in Table 6. The fact that unemployment rate is no longer the best predictor is no problem of course, but the best results in terms of RMSFE are lost. Only the limiting case of exponential weighting yields a slightly better 'best' forecast series, being the GDP Growth specification with a relative RMSFE of 0.904, whereas the 'best' forecast series from 5 for this method (Unemployment Rate) leads to a relative RMSFE of 0.906. In general the forecasting performance improves, but the 'best' performing forecast series may be lost by modifying the bias predictions.



**Table 7:** Combined forecasting performance unemployment data

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Shrinkage Weights</i> ( $\alpha = 0.5$ )		
Inflation	0.878**	-0.28 <sup>†</sup>
GDP Growth	0.874**	-0.25 <sup>†</sup>
Growth Gap	0.871**	-0.25 <sup>†</sup>
Output Gap	0.843**	-0.25 <sup>†</sup>
<i>Predicted Exponential Weights</i> ( $\gamma = 5$ )		
Inflation	0.940**	-0.45
GDP Growth	0.959*	-0.39
Growth Gap	0.924**	-0.39
Output Gap	0.977	-0.40
<i>Predicted Exponential Weights</i> ( $\gamma \rightarrow \infty$ )		
Inflation	1.058	-0.46
GDP Growth	1.031	-0.23 <sup>†</sup>
Growth Gap	1.026	-0.29 <sup>†</sup>
Output Gap	1.164	-0.30 <sup>†</sup>
<i>Benchmark Forecasts</i>		
Equal weights	1.186	-0.51
AO	1.275	0.15 <sup>†</sup>
VAR Growth Gap	0.863**	0.48

Table 7: This table shows the performance of the combined forecasts in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \* $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the Equal Weights benchmark. The Rel. RMSFE for equal weights is the real RMSFE value.

The results from combining unemployment rate forecasts displayed in Table 7 are not exactly following the inflation data results.

The shrinkage method leads to four significant improvements over the equal weights benchmark with unbiased forecasts, making it the best method in this table. However, only the Output Gap predictor achieves a lower relative RMSFE as well as lower bias than the best individual model, the VAR Output Gap model. The exponential weighting ( $\gamma = 5$ ) lead to three significant improvements over the benchmark. However, the forecasts are still biased. The limiting case of exponential weighting exhibits poor forecasting results, with no single improvement over the benchmark. However, most forecasts are unbiased.

Combining forecasts leads to better performance than a simple averaging approach, but in this case most individual models already lead to better forecasting results than the equal weighting forecasts. As mentioned, only one combined forecast series performs better than the best individual model. Because of this it is difficult to say whether combining forecasts results in substantial improvement. To answer this question, it would be a good idea to remove the direct forecasting specifications from the combinations, since these offer very poor individual results in comparison to the (V)AR models. Their results may drag the performance of combined forecasts down quite considerably, which possibly makes the (V)AR individual models look better than they actually are and consequently letting the combinations look worse. Sadly this is not included within the thesis due to time constraints, but will be mentioned as a limitation/further research possibility.

**Table 8:** Combined forecasting performance unemployment data AR(1)

	1970Q1- 2014Q1	
Predictor	Rel. RMSFE	Bias
<i>Shrinkage Weights</i>		
<i>(<math>\alpha = 0.5</math>)</i>		
Inflation	0.985	-0.34 <sup>†</sup>
GDP Growth	0.986	-0.33 <sup>†</sup>
Growth Gap	0.986	-0.33 <sup>†</sup>
Output Gap	0.980	-0.33 <sup>†</sup>
<i>Predicted</i>		
<i>Exponential</i>		
<i>Weights (<math>\gamma = 5</math>)</i>		
Inflation	0.935**	-0.40
GDP Growth	0.961	-0.40
Growth Gap	0.961	-0.40
Output Gap	0.942**	-0.40
<i>Predicted</i>		
<i>Exponential</i>		
<i>Weights (<math>\gamma \rightarrow T\infty</math>)</i>		
Inflation	0.951*	-0.40
GDP Growth	0.956	-0.40
Growth Gap	0.959	-0.40
Output Gap	0.953	-0.39
<i>Benchmark</i>		
<i>Forecasts</i>		
Equal weights	1.186	-0.51
AO	1.275	0.15 <sup>†</sup>
VAR Growth Gap	0.863**	0.48

Table 8: This table shows the performance of the combined forecasts in terms of Relative Root Mean Squared Error and bias. Explanation of symbols: \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \* $p < 0.1$ , a <sup>†</sup> indicates unbiasedness at ten percent level. The RMSFEs are shown relative to the Equal Weights benchmark. The Rel. RMSFE for equal weights is the real RMSFE value.

Modifying the bias predictions with ARMA terms generally yields better results in the case of inflation data, but the results in Table 8 show a different outcome. The performance of the shrinkage method has declined substantially, leading to zero significant improvements and worse bias performance, albeit still having unbiased forecasts.

The exponential weighting ( $\gamma = 5$ ) leads to more or less similar results in comparison to Table 7. Although there are now two instead of three significant improvements, the GDP Growth and Growth Gap predictors lead to relative RMSFE results that are very close to (10 percent) significance. Some predictors perform better, some perform worse. Bias results are quite similar, all forecasts are biased around 0.40.

The performance of the limiting case of exponential weighting now in fact does improve by the addition of the AR(1) term, leading to one significant improvement over the benchmark, with the other forecast series being quite close to significance as well. All forecasts from this method are biased though.

When comparing the best forecast series from Table 8 to the best individual forecasting model, we see that the individual model has a strictly lower relative RMSFE than any other combined forecast series. Also, the 'best' forecast series from Table 7 are lost similarly to the inflation data results. It is difficult to draw a conclusion from these results, since: one method performs worse, one method performs similarly and one method performs better compared to the original forecast combinations. I would say that it is easier to defend not adding AR terms to the bias predictions, due to the loss of quite some significant improvements from the shrinkage method, which is not compensated for with considerable significant improvements from other methods.

However, the results may be more credible if the direct forecasting specifications were removed. Also, this approach of forecasting unemployment rates is somewhat unorthodox and is typically done in other ways, such as using labor force flows data and multiple-state models (Barnichon et al., 2012).

## 5 Conclusion

In this thesis, the research goal was to replicate the results from Gibbs and Vasnev (2018) and extend their analysis by considering the direct modelling of ARMA terms within bias predictions to see whether this captures the typically present serial correlation in forecast errors and whether these results would hold for another data type as well.

For inflation data, the results are generally replicated; combining forecasts is useful and outperforms individual models and a forward looking combination approach leads to improvement compared to a simply strategy, an averaging benchmark, which often outperforms backward looking approaches. An ARMA(1,1) specification is shown to be the most suitable for modifying the bias predictions for this data type. Adding the ARMA(1,1) terms to the bias predictions leads to general improvement in comparison to the original forecast combinations, but the forecast series that had the best performance was one of the original series. This implies that adding autoregressive terms may overparametrize the bias prediction of the better models, decreasing their forecasting performance. However, the general conclusion should be that the modification substantially increases forecasting performance.

For unemployment data, the results are different. The forecast combinations do generally not outperform the best individual model with one exception. The forward looking approach does result in improvement over an equal weights benchmark, but most individual models are already better in terms of forecasting performance. The AR(1) specification is shown to be more suitable than the ARMA(1,1) specification for unemployment data. Although one method yields better results with this specification, the modification leads to the loss of four significant improvements over the benchmark, implying that the modification does not generally lead to better results.

## 6 Limitations and Discussion

There are some issues which may have had an impact on the outcome and credibility of this thesis. First of all, it is necessary to state that the replication part was not based on the most recent version of the paper by Gibbs and Vasnev (2018), but an older version. However, I did choose to use the most recent data available as regressors, although I only computed forecasts up to and including 2014Q1. In terms of forecasting this should not be a huge problem since I did use the real-time inflation and unemployment data and over the last few years the data has not been revised substantially.

The forecast procedure used by Gibbs and Vasnev (2018) seems to use some information from the 'future'. Since they estimate a model from 1947Q2 to 1965Q4 initially and then compute the first four-quarter-ahead forecast in 1966Q1 it feels like the parameters used to make a forecast contain information from three additional (future) data points. I chose to compute the first forecast in 1966Q4, leading to three less observations in the bias prediction regressions. However, the results still hold in a general sense and this approach does not include any future information, which is definitely preferable in my opinion.

Choosing to forecast two data types in a real-time manner, cost me quite some time to program, execute and process into the thesis report. For unemployment data, it would be nice to examine the results of forecast combinations excluding the poor performing direct forecasting specifications as mentioned within the results. Due to time constraints this is not done within the thesis but it may be a good idea for further research.

The contest about which ARMA specification is best for modification of the bias predictions consisted only of two specifications. An idea for further research would be to examine more possible modifications and perhaps deriving another strategy in determining which modifications are most suitable, as mine was somewhat heuristic and hence not necessarily backed by existing literature.

The unemployment data extension yields contradicting results to the inflation data analysis in some cases. However, the inflation data analysis is probably more credible since the Phillips Curve style of forecasting is widely used for this purpose. Also, the macro-economic variables considered are often used as explanatory variables for inflation, whereas unemployment data is usually forecasted with other models, that can have two or three states, and variables like labor force flows (Barnichon et al., 2012). I chose to keep the same macro-economic variables (as far as possible) for both inflation and unemployment to keep the results close, but further research could consist of unemployment forecasts using different models and explanatory variables. Moreover, the direct forecasts are likely to make results worse than necessary in case of unemployment forecasting. It also remains to be seen whether the comparison between inflation and unemployment forecasting is completely justified since both data types have different characteristics and in the literature, forecasting methods for both variables are usually different.

## References

- A. Atkeson, L. E. Ohanian, et al. Are phillips curves useful for forecasting inflation? *Federal Reserve bank of Minneapolis quarterly review*, 25(1):2–11, 2001.
- R. Barnichon, C. J. Nekarda, J. HATZIUS, S. J. STEHN, and B. PETRONGOLO. The ins and outs of forecasting unemployment: Using labor force flows to forecast the labor market [with comments and discussion]. *Brookings Papers on Economic Activity*, pages 83–131, 2012.
- R. T. Clemen. Combining forecasts: A review and annotated bibliography. *International journal of forecasting*, 5(4):559–583, 1989.
- F. X. Diebold and R. S. Mariano. Comparing predictive accuracy. *Journal of Business & Economic Statistics*, 13(3), 1995.
- C. Gibbs and A. L. Vasnev. Conditionally optimal weights and forward-looking approaches to combining forecasts. *Available at SSRN 2919117*, 2018.
- D. Harvey, S. Leybourne, and P. Newbold. Testing the equality of prediction mean squared errors. *International Journal of forecasting*, 13(2):281–291, 1997.
- D. F. Hendry and M. P. Clements. Pooling of forecasts. *The Econometrics Journal*, 7(1):1–31, 2004.
- W. C. Horrace and R. L. Oaxaca. Results on the bias and inconsistency of ordinary least squares for the linear probability model. *Economics Letters*, 90(3):321–327, 2006.
- W. K. Newey and K. D. West. A simple, positive semi-definite, heteroskedasticity and autocorrelation-consistent covariance matrix, 1986.
- A. W. Phillips. The relation between unemployment and the rate of change of money wage rates in the united kingdom, 1861–1957 1. *economica*, 25(100):283–299, 1958.
- J. H. Stock and M. W. Watson. Combination forecasts of output growth in a seven-country data set. *Journal of forecasting*, 23(6):405–430, 2004.
- J. H. Stock and M. W. Watson. Modeling inflation after the crisis. Technical report, National Bureau of Economic Research, 2010.
- A. Timmermann. Forecast combinations. *Handbook of economic forecasting*, 1:135–196, 2006.

## A List of programs written and used

estimateallgdpgrowth.prg

EViews program that makes all forecasts, computes the relevant metrics and makes the bias predictions for inflation data.

unempeestimateall.prg

EViews program that makes all forecasts, computes the relevant metrics and makes the bias predictions for unemployment data.

main.m

Main Matlab code for acquiring the forecast combinations and their evaluation metrics of regular inflation data bias predictions.

mainunemp.m

Main Matlab code for acquiring the forecast combinations and their evaluation metrics of unemployment data bias predictions, including ARMA results.

mainarma.m

Main Matlab code for acquiring the forecast combinations and their evaluation metrics of inflation ARMA bias predictions.

shrinkage.m

Shrinkage method for combinations.

AOcreator.m

Creates an random walk sequence.

DMstat.m

Computes Diebold Mariano statistic with Harvey correction.

exponentialweighting.m

Employs exponential weighting method.

obtainData.m

Obtains some of the regressors such as the growth gap measure.

equalweights.m

Creates equal weights.