



On the Initialization of Adaptive Learning in Macroeconomic Models

DOI:

[10.1016/j.jedc.2017.03.002](https://doi.org/10.1016/j.jedc.2017.03.002)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Berardi, M., & Galimberti, J. (2017). On the Initialization of Adaptive Learning in Macroeconomic Models. *Journal of Economic Dynamics and Control*, 78, 26–53. <https://doi.org/10.1016/j.jedc.2017.03.002>

Published in:

Journal of Economic Dynamics and Control

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



On the Initialization of Adaptive Learning in Macroeconomic Models*

MICHELE BERARDI
University of Manchester

JAQUESON K. GALIMBERTI[†]
ETH Zurich

Abstract

We review and evaluate methods previously adopted in the applied literature of adaptive learning in order to initialize agents' beliefs. Previous methods are classified into three broad classes: equilibrium-related, training sample-based, and estimation-based. We conduct several simulations comparing the accuracy of the initial estimates provided by these methods and how they affect the accuracy of other estimated model parameters. We find evidence against their joint estimation with standard moment conditions: as the accuracy of estimated initials tends to deteriorate with the sample size, spillover effects also deteriorate the accuracy of the estimates of the model's structural parameters. We show how this problem can be attenuated by penalizing the variance of estimation errors. Even so, the joint estimation of learning initials with other model parameters is still subject to severe distortions in small samples. We find that equilibrium-related and training sample-based initials are less prone to these issues. We also demonstrate the empirical relevance of our results by estimating a New Keynesian Phillips curve with learning, where we find that our estimation approach provides robustness to the initialization of learning. That allows us to conclude that under adaptive learning the degree of price stickiness is lower compared to inferences under rational expectations.

Keywords: expectations, adaptive learning, initialization, algorithms, hybrid New Keynesian Phillips curve.

JEL codes: C63, D84, E03, E37.

1 Introduction

Adaptive learning algorithms have been proposed to provide a procedural rationality view on agents' process of expectations formation. Reopening a long standing debate on how should

*We gratefully acknowledge the comments and evaluations provided by an Associate Editor and two anonymous reviewers. An earlier version of this paper was presented at the 2016 EEA-ESEM joint meeting in Geneva. Any remaining errors are our own.

[†]Corresponding author. E-mail: galimberti@kof.ethz.ch.

1 expectations be modeled in macroeconomic models, the heuristics provided by learning algo-
2 rithms come at the cost of introducing new degrees of freedom into the analysis. One such
3 an open node relates to how these recursive mechanisms should be initialized in order to be
4 representative of agents' learning-to-forecast behavior.

5 The main characteristic of the adaptive learning approach is its reliance on recursive al-
6 gorithms in order to represent how agents update their beliefs as new observations about the
7 economic relationship of interest become available. Such recursions naturally demand an ini-
8 tial starting point, and it is the numerical specification of these conditions that we denote as
9 the initialization problem. Clearly, the uncertainties affecting the initialization of the learning
10 process will propagate recursively into the predictions obtained with the model, and it seems
11 crucial that the researcher understands the magnitude of these distortions and how they can
12 affect structural inferences.

13 In this paper we investigate this issue with particular attention to the applied literature of
14 learning in macroeconomics. Here applied is taken to encompass both theoretical simulations
15 as well as exercises of empirical estimation and calibration. Examples can be found in Sargent
16 (1999); Marcet and Nicolini (2003), or more recently in Eusepi and Preston (2011); Milani
17 (2011), between many others cited throughout the paper. The main distinctive feature of these
18 works consists in the replacement of the rational expectations (RE) assumption of an instanta-
19 neous adjustment of agents expectations, with a characterization of agents as adaptive learners
20 of their own environment. More generally, our study will be relevant for scholars interested
21 in the methods needed to uncover the initial beliefs of economic agents in models where such
22 beliefs actually matter for economic dynamics.

23 The economic relevance of the initialization issue can be illustrated considering the long
24 debated causes of the period of Great Inflation during the 1970s in the US. One of the main
25 explanations for that episode comes from Sargent's (1999) hypothesis that the evolution of US
26 inflation rates over the period can be attributed to the evolution of the monetary authority's
27 beliefs about the trade-off between inflation and unemployment, the so-called Phillips curve.
28 Subsequent studies advanced on this issue attributing the rise in inflation rates to delayed policy
29 responses to ongoing structural changes in the economy of that period (see, e.g., Bullard and
30 Eusepi, 2005; Orphanides and Williams, 2005a). Importantly, as evidenced in Primiceri (2006);
31 Sargent et al. (2006), the point of departure in policymaker's beliefs, which is what we refer to
32 as the learning initials in this particular context, is a crucial feature in such explanation.

33 Assumptions about initial beliefs also matter for the fit of models that introduce adaptive
34 learning on the side of other market participants, such as households and firms. Examples
35 are given by Carceles-Poveda and Giannitsarou (2008) for asset pricing models, Huang et al.
36 (2009) in a standard growth model, and Slobodyan and Wouters (2012a,b) in a medium-scale
37 dynamic stochastic general equilibrium (DSGE) model. Overall, these studies present results
38 showing that whereas the introduction of learning has interesting effects on the dynamics and
39 the fit of models to the data, a great portion of the improvements may be associated to transition

1 dynamics from specific initial beliefs. Hence, it is important to have a systematic evaluation of
2 the different alternatives available as initialization methods, and we attempt to fill that gap in
3 this paper.

4 We review the literature in order to pool together the pre-existing initialization methods
5 into an archetypal classification that can be broadly defined in three major classes: equilibrium-
6 related methods, training sample-based methods, and estimation-based methods. The equilibrium-
7 related initializations are generally obtained taking rational expectations equilibrium (REE) as
8 a reference, and exploring distributional deviations from that assumption. The training sample-
9 based initializations, as the name suggests, are obtained with the application of the learning
10 algorithm (or variations) over a pre-sample of observations that is left aside from the original
11 sample of data available. Here we distinguish between two main forms of this method: (i) a
12 Weighted Least Squares (WLS) approach, which is equivalent to the application of a constant
13 gain Least Squares (LS) algorithm¹ to the training sample; and, (ii) an Ordinary Least Squares
14 (OLS) approach, which is equivalent to the decreasing gain form of the LS algorithm. Finally,
15 the more recent estimation-based methods consist of approaches involving the joint estimation
16 of the initials with other model parameters, hence allowing the use of the same data that is
17 used for inferences about structural features of the model to guide the specification of learning
18 initials.

19 We compare the initialization methods on the basis of the accuracy of their delivered initial
20 estimates and their effects over the accuracy of other estimated model parameters. To evaluate
21 accuracy we derive measures of the Mean Squared Deviation (MSD) between true parame-
22 ter values and their corresponding estimates obtained according to the different initialization
23 methods. We relate the MSD measures to two main principles to judge the quality of an initial-
24 ization. First, we look at the coherence of the initialization estimates to the dynamics implied
25 by the learning process; second, we consider how susceptible the method is to biases that push
26 up the model's explanatory power over the initial portion of observations in the sample.

27 The accuracy of initialization methods is first examined analytically within the simplified
28 framework of an example model where agents are required to learn only a constant. Under
29 these conditions we establish some important relationships between the invariant distribution
30 of the learning estimates and the MSDs associated to the initialization methods: (i) an REE-
31 based initial yields an MSD equal to the variance of the learning estimates; (ii) the accuracy
32 of training sample-based initials depends on the number of observations left aside for training,
33 and the specification of the learning gains; under a constant gain equal to that underlying the
34 data generating process (WLS-based), MSDs tend to zero as the training sample increases;
35 under decreasing gains (OLS-based), the MSDs converge to those obtained under the REE-
36 based initial; (iii) an initial estimate based on the minimization of the model's sum of squared
37 residuals tends to yield smaller MSDs as the estimation sample grows. However, we show that

¹Consistent with the majority of the adaptive learning literature, in this paper we focus on the LS algorithm as representative of how agents update their expectations (see also Berardi and Galimberti, 2014, on this point).

1 only the WLS-based method is capable of converging to the true initial estimate. We believe
2 these results should prove useful for applied researchers interested in the simulation of models
3 with adaptive learning.

4 For empirical purposes, one is confronted with a more intricate initialization problem,
5 namely, the estimation of other model parameters. To evaluate initialization methods under
6 these circumstances we conduct several simulation exercises based on a model of the New
7 Keynesian Phillips curve (NKPC) with learning. Our simulation results point to some interest-
8 ing findings. First, the analytical results based on the assumption of known model parameters
9 are confirmed for the NKPC model; particularly, the training sample-based initialization meth-
10 ods are in general favored in terms of the coherence criterion, since these methods are found
11 to provide more accurate estimates of the learning initials. Second, the performance of the ini-
12 tialization methods is sensitive to the accuracy of their associated model parameters' estimates;
13 here the equilibrium-related initials can prove useful for being less sensitive to the estimation
14 of the learning gain, which turns out to be severely affected by finite sample distortions.

15 Notwithstanding, we find that the pre-determined initials were much less vulnerable to the
16 estimation of other model parameters than the jointly estimated initials. Regarding the latter,
17 we find that the estimation approach plays a key role in determining the quality of the initials
18 and their effects over other parameters' estimates. Particularly, an estimation based solely
19 on traditional moment conditions, derived from the assumption that the model's unobserved
20 disturbance is a martingale difference sequence, leads to severe distortions to the estimated
21 initials as the sample size grows. Furthermore, these distortions can spillover to the estimates of
22 the other model parameters. Based on our analytical derivations, we proposed a solution to this
23 problem with the inclusion of a squared residual criterion to the estimation objective, denoting
24 this as the augmented approach. Although we show that our augmented approach is successful
25 in restoring the consistency of the jointly estimated initials, this solution is only effective for
26 large samples. Importantly, under the pre-determined initials the negative spillover effect from
27 the initialization errors to the model parameters' estimates is not as strong as observed under
28 the joint estimation approaches.

29 This last finding is of particular relevance for empirical analysis, where interest is usually
30 in uncovering the underlying values of structural parameters that may validate the model's con-
31 sistency with data evidence. In order to further enhance our understanding on the relevance of
32 these different initialization methods for applied macroeconomics, we also present an empiri-
33 cal application on the determination of US inflation rates under the Phillips curve framework.
34 Adopting a generalized method of moments (GMM) estimation approach, we find that the ini-
35 tials and the estimation criterion can, indeed, affect the estimates of structural parameters. Our
36 results indicate that allowing for adaptive learning in the determination of inflation rates leads
37 to a lower degree of price stickiness. This finding, nevertheless, is only robust to the alternative
38 specifications of the initial estimates under our augmented estimation approach.

1 Finally, we also contribute to the literature on the estimation of models with adaptive learn-
 2 ing by proposing the introduction of robust stability restrictions in the estimation of the model.
 3 Because adaptive learning under constant gain can generate unrealistic dynamics due to in-
 4 stabilities in the recursive estimation of the learning coefficients (see Evans and Honkapohja,
 5 2009), the estimation of models under learning can be extremely sensitive to the sample of
 6 observations and the range of values allowed for the learning gain. To deal with this issue we
 7 impose additional constraints to the estimates of the learning gains, drawing upper bounds con-
 8 ditioned on the values of other parameters through simulations of the model. One advantage
 9 of this approach, compared to an unconditional reduction of the gain upper bound, is that it
 10 provides weaker constraints robust to learning stability for the joint estimation of the model
 11 parameters with the learning gain and initials.

12 The remainder of this paper proceeds as follows. In section §2 we provide a brief introduc-
 13 tion to the use of adaptive learning in macroeconomic models, and establish the initialization
 14 problem under a simple “learn-the-average” example model. A review of the initialization
 15 methods previously adopted in the literature is presented in section §3 together with an ana-
 16 lytical evaluation of these methods under the example model. We then proceed to present our
 17 simulation analysis, in section §4, and an empirical application, in section §5, both aiming at
 18 a comparative evaluation between the different methods of initialization under more realistic
 19 circumstances, such as the case where the learning gain and other model parameters require
 20 estimation. Finally, we conclude this paper with some remarks in section §6. Some key deriva-
 21 tions and supplementary results are provided in the Appendices.

22 **2 Adaptive Learning and the Initialization Problem**

23 **2.1 A brief primer on adaptive learning**

24 Adaptive learning is introduced in macroeconomic models as an alternative to the assumption
 25 that agents hold rational expectations. One implication of the rational expectations assumption
 26 is that agents’ beliefs are always consistent with the true model of the economy. Hence, under
 27 RE the economy instantaneously adjusts itself towards an equilibrium after any kind of shock
 28 that may have realistically affected agents’ beliefs. In contrast, adaptive learning introduces
 29 some degree of persistence in the process through which agents update their beliefs, which
 30 allows such beliefs to deviate from RE in the short run, while keeping up with the idea of
 31 consistency in the long run.

32 To help fix ideas consider an univariate linear forward-looking model, where the determi-
 33 nation of the current value of a variable of interest, y_t , depends on the value expected for that
 34 same variable in the next period, y_{t+1}^e plus a mean zero random shock, u_t , i.e.,

$$y_t = \beta y_{t+1}^e + u_t. \tag{1}$$

1 Simple as it stands, this specification may represent the reduced form of the equilibrium equa-
 2 tions of an economic model which could potentially be non linear; it also corresponds, e.g., to
 3 simplified versions of two well known models: the Cagan model of inflation, letting y_t stand
 4 for the price level and u_t for a mean zero random supply of money; and, the standard model
 5 of asset pricing under risk neutrality, letting y_t stand for the asset price and u_t for a mean zero
 6 random sequence of dividends.

7 A solution to model (1) requires the specification of agents' perceived law of motion (PLM),
 8 which depicts how agents form expectations. Particularly, if agents condition their forecasts on
 9 a constant, $y_{t+1}^e = 0$ solves the model for any β . Hence, the stochastic process followed by
 10 the economy, also known as the actual law of motion (ALM), is directly determined by the
 11 specification of agents' PLM.

12 Under learning the corresponding ALM is given by

$$y_t = \beta\phi_{t-1} + u_t, \quad (2)$$

13 where ϕ_{t-1} denotes agents' estimates of the constant in their PLM based on observations avail-
 14 able up to the previous period. Different recursive algorithms have been proposed in the litera-
 15 ture to represent how agents update such estimates. Due to its widespread popularity between
 16 econometricians, one natural choice for that purpose has been the Least Squares (LS) algorithm,
 17 which can be generally defined as follows.

18 **Algorithm 1 (LS).** *Let agents' PLM of y_t be given by a linear regression of the form*

$$y_t = \mathbf{x}_t' \boldsymbol{\phi}_{t-1} + \varepsilon_t, \quad (3)$$

19 where $\mathbf{x}_t = (x_{1,t}, \dots, x_{K,t})'$ is a set of pre-determined variables, possibly including a constant
 20 (e.g., $x_{1,t} = 1$) and lags of y_t , $\boldsymbol{\phi}_t = (\phi_{1,t}, \dots, \phi_{K,t})'$ stands for a vector of coefficients, possibly
 21 time-varying, and ε_t denotes an unpredictable disturbance term. Under this context, the LS
 22 estimates of $\boldsymbol{\phi}_t$, conditional on observations up to time t , are given by

$$\hat{\boldsymbol{\phi}}_t = \hat{\boldsymbol{\phi}}_{t-1} + \gamma_t \mathbf{R}_t^{-1} \mathbf{x}_t (y_t - \mathbf{x}_t' \hat{\boldsymbol{\phi}}_{t-1}), \quad (4)$$

$$\mathbf{R}_t = \mathbf{R}_{t-1} + \gamma_t (\mathbf{x}_t \mathbf{x}_t' - \mathbf{R}_{t-1}), \quad (5)$$

23 where γ_t is a learning gain parameter, and \mathbf{R}_t stands for an estimate of the regressors matrix
 24 of second moments.

25 The LS algorithm is originally motivated as the result from the minimization of a weighted
 26 sum of squared errors, where the weights are determined by the learning gain parameter (see
 27 Berardi and Galimberti, 2013). Hence, the learning gain stands for a parameter determining
 28 how quickly a given information is incorporated into the algorithm's coefficients estimates.
 29 There are two particular cases of interest: (i) when $\gamma_t = 1/t$, every observation receives the

1 same weight and (4)-(5) reduces to the (recursive) Ordinary Least Squares (OLS); and, (ii) un-
 2 der a constant gain, past observations receive geometrically decaying weights and (4)-(5) can be
 3 viewed as a (recursive) Weighted Least Squares (WLS) with weights given by $(1 - \gamma)^j$, where
 4 j indexes for the number of periods between the weighted observation and the last observation
 5 in the sample. Our focus is on the constant gain specification due to its relevance for applied
 6 purposes: because it allows for a continuous operation of the algorithm's tracking capabilities,
 7 the constant gain can capture time-varying effects of different sources, such as structural breaks
 8 or the out-of-equilibrium dynamics generated by stochastic shocks.

9 For a PLM with an intercept only the constant-gain LS algorithm simplifies to

$$\phi_t = \phi_{t-1} + \gamma (y_t - \phi_{t-1}). \quad (6)$$

10 Substituting y_t from model (1)'s ALM one can find that

$$\phi_t = \delta^t \phi_0 + \gamma \sum_{i=0}^{t-1} \delta^i u_{t-i}, \quad (7)$$

where $\delta = 1 - \gamma(1 - \beta)$. Taking the unconditional expectation of this expression we find that convergence to the RE equilibrium requires that $|\delta| < 1$, which is only possible when $\beta < 1$ and $0 < \gamma < 2/(1 - \beta)$. Notice this last condition is always true for $0 < \beta < 1$ and meaningful gain values in the interval $0 < \gamma < 1$. Under this assumption the variance of the long run distribution of the learning estimates around the REE, denoted by $\bar{\sigma}_\phi^2$, is then given by

$$\begin{aligned} \bar{\sigma}_\phi^2 &= \lim_{t \rightarrow \infty} E [\phi_t^2], \\ &= \frac{\gamma \sigma_u^2}{(1 - \beta)(1 + \delta)}, \end{aligned} \quad (8)$$

11 where σ_u^2 is the variance of u_t . Thus, the dispersion of the learning estimates around the REE
 12 increases with the value of the learning gain; see Evans and Honkapohja (2001, Theorem 7.9)
 13 for a more general result on the relationship between the learning gain and the invariant distri-
 14 bution of the learning estimates.

15 **2.2 The initialization problem**

16 Recursive learning algorithms naturally demand an initial starting point, and it is the numerical
 17 specification of these conditions that we denote as the initialization problem. By the recursive
 18 nature of learning, any error in the initial estimates will propagate recursively into the predic-
 19 tions obtained with the model. Consider the case of our example model, (2)+(6); letting $\hat{\phi}_0$
 20 stand for a guess of the true value of ϕ_0 , the model prediction of y_{t+1} associated with this initial

1 is given by $\hat{y}_{t+1} = \beta \hat{\phi}_t$, where $\hat{\phi}_t$ is obtained from (4) as

$$\begin{aligned}
\phi_0 &= \hat{\phi}_0, \\
\hat{\phi}_1 &= \hat{\phi}_0 + \gamma (y_1 - \hat{\phi}_0), \\
&\vdots \\
\hat{\phi}_t &= \hat{\phi}_{t-1} + \gamma (y_t - \hat{\phi}_{t-1}).
\end{aligned} \tag{9}$$

2 Let the corresponding prediction error be denoted by $\hat{\Delta}_{t+1} = y_{t+1} - \hat{y}_{t+1}$; then, the mean
3 squared prediction error (MSPE) from this model amounts to

$$\begin{aligned}
E \left[\hat{\Delta}_{t+1}^2 \right] &= E \left[\left(\beta (\phi_t - \hat{\phi}_t) + u_{t+1} \right)^2 \right], \\
&= \beta^2 \lambda^{2t} \mathcal{D}_{\phi_0} + \sigma_u^2,
\end{aligned} \tag{10}$$

4 where $\lambda = 1 - \gamma$, and \mathcal{D}_{ϕ_0} stands for the initialization's Mean Squared Deviation (MSD), which
5 is more generally defined as:

6 **Definition 1** (MSD). The Mean Squared Deviation between the true value of a parameter,
7 θ_t , which may be a learning initial (e.g., ϕ_0 , R_0) or a time-invariant model parameter, and a
8 corresponding estimate, $\hat{\theta}_t$, is given by

$$\mathcal{D}_{\theta_t} = E \left[\left(\theta_t - \hat{\theta}_t \right)^2 \right]. \tag{11}$$

9 Clearly, assuming that $0 < \gamma < 1$, (10) shows that the effects of initialization errors tend
10 to disappear as the distance from the initial point increases. Also notice that as γ increases, the
11 prediction error associated to an initialization error decreases. Hence, the smaller the learning
12 gain, the more important are the learning initials for the accuracy of the predictions obtained
13 with the model. For empirical purposes both the model and learning parameters, say β and γ
14 in our example model, are not known *a priori* and may therefore require estimation. In that
15 case initialization errors can further aggravate the accuracy of the model predictions through
16 its effects on parameters' estimation error. In fact, as we will show in our simulation analysis
17 further below, the estimation of the model and learning parameters can be severely affected by
18 the misspecification of the learning initials.

19 2.3 Evaluation criteria

20 Our relative assessment of the initialization methods that we describe in the next section will be
21 guided by two principles that we consider of relevance for applied adaptive learning research:
22 (i) the initials COHERENCE to the learning process; and, (ii) the initials SUSCEPTIBILITY to
23 bias the model's explanatory power and the estimation of its parameters.

1 In empirical settings, a proper initialization of the learning algorithm requires to find out
2 what were agents' beliefs at the beginning of the sample of data. To achieve this goal it is
3 important to understand the statistical properties of the learning process we are trying to mimic.
4 Recursive estimation algorithms are statistically characterized by two main distinct phases: a
5 transient phase, where the estimates are so far apart from the true parameter values that the
6 upcoming sequence of updates can easily achieve substantial improvements to the accuracy of
7 the estimates; and a stationary phase, where most of the updates to the estimates are essentially
8 just tracking tiny disturbances that may affect the system under estimation. Hence, if the initial
9 beliefs should reflect the continuation of an estimation process that was already in motion prior
10 to the sample beginning, an initialization method will satisfy the COHERENCE criterion when
11 it can deliver estimates as close as possible to the algorithm's long run operation². Under our
12 evaluation measure this corresponds to a minimization of the initials' MSD.

13 Another empirical issue is how much can the learning initials affect the accuracy of the
14 estimates of other model parameters. Under standard likelihood-based estimation approaches,
15 every data point is given the same weight on the estimation of a structural parameter that is
16 assumed to be constant throughout the sample period. Under learning this weighting profile
17 can be easily manipulated by tweaking the initial learning estimates so as to induce a transient
18 phase in the portion of the sample that follows the initialization, which potentially increases
19 the model's explanatory power. In the context of our analysis, such SUSCEPTIBILITY to biases
20 is measured by looking at the MSDs of some key parameter estimates across the different
21 initialization methods.

22 These principles may inherently generate a trade-off for the assessment of the initialization
23 methods: one can always give up some degree of the COHERENCE delivered by a learning initial
24 in exchange for some SUSCEPTIBILITY to tweak that initial in order to improve the model fit to
25 the data. Thus, it is important to obtain a quantitative assessment of how the different methods
26 perform with respect to these criteria.

27 **3 Review of Initialization Methods**

28 In this section we review initialization methods adopted in the previous literature. We also
29 present expressions for their associated MSD under the example model described in the previ-
30 ous section, i.e., (2)+(6). For that purpose, throughout this section we assumed that: (i) both
31 the learning gain and the model parameter are known; and, (ii) the learning process has already
32 converged to its invariant distribution at the period where the initial is evaluated, and at the be-
33 ginning of the training sample, when applicable. Numerical simulations validating our results
34 are presented in Appendix C.1.

²Admittedly, one may also be interested in obtaining initials that reflect the transient phase that follows the occurrence of a large shock that shifted agents' beliefs away from equilibrium just before the initialization data point.

3.1 Equilibrium-related methods

One way to initialize learning algorithms is obtained by using the existing knowledge about the law of motion(s) generating the data. Particularly, conditional on the knowledge about the model specification and the parameter values, one can easily obtain the REE implied values of agents' PLM coefficients and use these equilibrium values as reference for the initial estimates. In the case of our example model, such initial would be given by $\hat{\phi}_0^{REE} = 0$, and the associated MSD equals to the variance of the learning estimates, $\bar{\sigma}_\phi^2$, as given by (8).

Although this method was naturally appealing in earlier works with theoretical simulations, such as in Bray and Savin (1986), its debut into the applied literature came in the seminal contribution of Sargent (1999). Its usage has since been prominent in studies on the effects of replacing the assumption of frictionless REE by the sticky process of expectations formation through adaptive learning (e.g., Marcet and Nicolini, 2003; Bullard and Eusepi, 2005; Orphanides and Williams, 2005b). For simulations, robust inferences can be obtained through this method by drawing the initials from a distribution centered around the REE values (see Carceles-Poveda and Giannitsarou, 2007).

Empirically, the uncertainties about the true model parameters may complicate the adoption of this method. One alternative is to approach the issue in two stages: first, model estimates are obtained under the REE assumption; second, these estimates are used to calculate the PLM coefficient values corresponding to the REE, which are then plugged back in as initial estimates for the algorithm's recursion for the analysis under learning (see Slobodyan and Wouters, 2012b; Ormeño and Molnár, 2015). One criticism to this practical solution is that it seems very likely that the REE estimates obtained in the first step will be biased for not taking the learning effects into account. Slobodyan and Wouters (2012a, p. 93) relax the first step by using model estimates obtained under a baseline learning specification instead of under RE, though still fixing the implied initials in the second step. Later, in our simulation and empirical exercises, we show that this issue can be circumvented by allowing the REE initials to be determined jointly with the model parameters under learning.

The REE-based initials do not provide ideal initial estimates for cases where there is prior information that the economy was in a transient phase at the beginning of the sample. In such a case, an alternative is provided by the *ad-hoc* initialization method, where the initials are hand-picked by the researcher. When taking the REE-based initials as a reference, this method provides a way to validate the sensitivity of results obtained under the former approach (e.g., Milani, 2007; Carceles-Poveda and Giannitsarou, 2008). In fact, one of the main uses of *ad-hoc* initials is to deal with the possibility of structural changes around the periods of the initials: when the changes affect the REE, agents may not be able to instantaneously adjust to the new equilibrium, and could therefore be forming expectations consistent with the previous equilibrium at the time of the initialization (see also Carceles-Poveda and Giannitsarou, 2007, p. 2679).

3.2 Training sample-based methods

Another common approach in learning applications is to use a pre-sample of observations in order to obtain the initial estimates. This is especially recommended for the cases where there is not enough previous knowledge about the system under estimation so as to allow an educated guess. The origins of this approach can be traced back into the engineering literature (see, e.g., Ljung and Soderstrom, 1983, pp. 299-303), where it is often suggested that the coefficients should be initialized with the value of zero and an initial training sample should be left aside to let the algorithm adjust its estimates according to the underlying calibration.

For applied purposes, it is often easier to adopt the non-recursive version of the learning algorithm to estimate the initials over the training sample. Letting P denote the number of observations set aside for the initialization, application of (4)-(5) in the training sample results in

$$\mathbf{R}_P = \gamma_P \sum_{i=1}^P w_i \mathbf{x}_i \mathbf{x}_i' + w_0 \mathbf{R}_\emptyset, \quad (12)$$

$$\hat{\boldsymbol{\phi}}_P = \gamma_P \mathbf{R}_P^{-1} \sum_{i=1}^P w_i \mathbf{x}_i y_i, \quad (13)$$

where $\{w_i\}_{i=0}^P$ stands for the sequence of weights given to each observation in the training sample, and \mathbf{R}_\emptyset may incorporate prior information regarding the uncertainty surrounding the determination of the coefficients estimates. Under the assumption of a Gaussian random walk parameter drift model for $\boldsymbol{\phi}_t$, Berardi and Galimberti (2013) have shown that \mathbf{R}_t is inversely related to the matrix of mean squared errors associated to the Kalman filter coefficients estimates, $E \left[\left(\boldsymbol{\phi}_t - \hat{\boldsymbol{\phi}}_t \right) \left(\boldsymbol{\phi}_t - \hat{\boldsymbol{\phi}}_t \right)' \right]$. Hence, in a Bayesian interpretation, as $\mathbf{R}_\emptyset \rightarrow \mathbf{0}$ the prior becomes more diffuse, since it is associated with a higher uncertainty about the coefficients estimates³.

Depending on the weighting scheme and the prior estimates, there are two main variations of this method in the learning literature: the OLS-based (e.g., Williams, 2003; Orphanides and Williams, 2005a; Sargent et al., 2006) and the WLS-based (e.g., Primiceri, 2006; Milani, 2007, 2008, 2011, 2014; Huang et al., 2009; Chevillon et al., 2010; Eusepi and Preston, 2011; Lubik and Matthes, 2014) initials.

3.2.1 WLS-based initials

The WLS-based method derives from the Weighted Least Squares interpretation of the learning algorithm under a constant gain specification. In the training sample the initials associated to this method are obtained by setting $\gamma_P = \gamma$, and $w_i = (1 - \gamma)^{P-i}$ in (12)-(13). Under our

³Notice that when $\mathbf{R}_\emptyset = \mathbf{0}$, a necessary condition for \mathbf{R}_P to be invertible, as required in (13), is that $P \geq K$.

example model, the MSD associated to the WLS initial is given by

$$\begin{aligned}\mathcal{D}_{\phi_0}^{WLS} &= E \left[(\lambda^P \phi_{-P})^2 \right], \\ &= \lambda^{2P} \bar{\sigma}_\phi^2,\end{aligned}\tag{14}$$

1 from which we obtain the following proposition.

2 **Proposition 1.** *Under the model given by (2)+(6), the WLS-based initialization yields, on av-*
3 *erage, smaller MSDs than the REE-based initials.*

4 *Remark 1.* This result is specific to the model under consideration, and it stems from the fact
5 that the REE is equal to the diffuse prior used for the WLS initial, i.e., $\phi_{REE} = \hat{\phi}_{-P}^{WLS} = 0$.
6 Because the WLS method incorporates information from the training sample to that prior, it
7 will always, on average, attain more precise initial estimates. In fact, if we add a constant α to
8 the model in (1), $\phi_{REE} = \alpha / (1 - \beta)$ and $\lim_{t \rightarrow \infty} E[\phi_t^2] = \bar{\sigma}_\phi^2 + \phi_{REE}^2$; then, whereas $\mathcal{D}_{\phi_0}^{REE}$
9 would still be equal to the variance of the invariant distribution of learning estimates, \mathcal{D}^{WLS}
10 would carry the effect of ϕ_{REE}^2 .

11 For models with regressors in the PLM, as we will see in the simulation exercise of next
12 section, the initialization of the matrix of second moments, (5), also becomes important. Here
13 we consider two possibilities for the prior, \mathbf{R}_\emptyset . One, based on REE reasoning, is to set it to
14 the sample estimate of the long-run covariance matrix of the regressors. Ideally, the sample
15 used for such estimation should be restricted to the training sample itself, in order to prevent
16 contamination of the initials due to the effects of changes in the statistical properties of the data
17 that were not present before the initialization period. The second alternative we consider is to
18 follow a diffuse approach and set $\mathbf{R}_\emptyset = \mathbf{0}$.

19 3.2.2 OLS-based initials

20 The OLS-based method, as the abbreviation suggests, is based on the Ordinary Least Squares
21 estimator that is widely known among econometricians for possessing some well desired prop-
22 erties, such as consistency and efficiency in the estimation of linear models. For training sample
23 initialization, it is obtained by setting $\gamma_P = 1/P$, $w_i = 1$, and $\mathbf{R}_\emptyset = \mathbf{0}$ in (12)-(13). One impor-
24 tant advantage of this method relates to its convergence speed: the fact that a relatively higher
25 gain value is used in the first iterations of the algorithm within the training sample tends to
26 accelerate its convergence to the true initials.

27 Under our example model, the MSD associated to the OLS-based initial is given by

$$\mathcal{D}_{\phi_0}^{OLS} = \left(\delta^P - \frac{\beta(\delta^P - 1)}{P(\delta - 1)} \right)^2 \bar{\sigma}_\phi^2 + \sigma_u^2 \mathcal{K},\tag{15}$$

1 where \mathcal{K} is a function of the model parameters and the length of the initialization sample.
2 Clearly, a necessary, though not sufficient, condition for the OLS-based initials to yield smaller
3 MSD than the WLS-based initial is that $\lambda^P > \delta^P - \beta (\delta^P - 1) P^{-1} (\delta - 1)^{-1}$. However, notice
4 that as $P \rightarrow \infty$, $\mathcal{D}_{\phi_0}^{OLS} \rightarrow E[\phi_t^2]$, which leads to our next proposition.

5 **Proposition 2.** *Under the model given by (2)+(6), as the length of the training sample increases*
6 *to infinity, the MSD of the OLS-based initialization converges, on average, to the MSD of the*
7 *REE-based initialization.*

8 *Remark 2.* A direct implication of Propositions 1 and 2 is that the WLS-based initialization
9 will tend to outperform the OLS-based method as the length of the training sample increases.

10 3.3 Estimation-based methods

11 Another approach to the initialization of learning coefficients is to add the initials to the set of
12 the model's parameters, and estimate them jointly. The idea can be traced back to the landmark
13 work by Sargent et al. (2006), where the estimation of the monetary authority's initial and
14 consecutive stream of beliefs provided evidence in favor of Sargent's (1999) hypothesis on the
15 "Conquest of American Inflation": namely, that the rise and fall of post-WWII inflation in
16 the US can be attributed to the evolution of the monetary authority's beliefs about the trade-
17 off between inflation and unemployment. In spite of some early criticisms (see discussion
18 in Carboni and Ellison, 2009), the approach of joint estimation of initials has been slowly
19 incorporated into broader applications of adaptive learning. After being hinted as a possibility
20 in Milani (2007, p. 2071) and Huang et al. (2009, p. 397), initial attempts have focused on the
21 effects of the joint estimation of the initial matrix of second moments, \mathbf{R}_0 (e.g., Milani, 2008,
22 2011), and more recently on the estimation of the complete set of learning initials (as in, e.g.,
23 Slobodyan and Wouters, 2012b; Gaus and Ramamurthy, 2014).

24 There are different alternatives for the estimation of dynamics macroeconomic models;
25 here, for illustrative purposes, we adopt the approach of minimizing the sum of squared resid-
26 uals (SSR). Namely, the estimation-based initials are obtained according to

$$\hat{\phi}_0^{SSR} = \arg \min_{\hat{\phi}_0} SSR = \sum_{i=0}^{T-1} \hat{u}_{T-i}^2, \quad (16)$$

27 where T stands for the number of observations used for estimation. Under our example model's
28 ALM, the initial estimate that solves this problem is given by

$$\hat{\phi}_0^{SSR} = \sum_{i=0}^{T-1} \mathcal{H}_i y_{T-i}, \quad (17)$$

29 where

$$\mathcal{H}_i = \frac{\gamma (\lambda^{-i} (1 + \lambda) - \beta \lambda (\lambda^{-i} - \lambda^i))}{\beta \lambda (\lambda^{-T} - \lambda^T)}. \quad (18)$$

The MSD associated to these estimated initials are then given by

$$\begin{aligned}\mathcal{D}_{\phi_0}^{SSR} &= E \left[\left(\phi_0 - \hat{\phi}_0^{SSR} \right)^2 \right], \\ &= \frac{\gamma(\gamma - 2)}{\beta^2 (\lambda^{2T} - 1)} \sigma_u^2.\end{aligned}\tag{19}$$

1 **Proposition 3.** *Under the model given by (2)+(6), whereas increasing the length of the estima-*
 2 *tion sample improves the accuracy of SSR-based initials, the initial estimates never converge*
 3 *to their corresponding true values.*

4 *Remark 3.* Notice that evaluating (19) in the limit as $T \rightarrow \infty$ one finds that the higher the
 5 learning gain, the higher the asymptotic error incurred by the SSR-based initial.

6 Most estimation approaches share a common idea of looking for the combination of param-
 7 eter values that maximize the fit of the model, or its implications, to available macroeconomic
 8 data. Hence, the joint estimation of learning initials can have an appealing motivation for pro-
 9 viding those estimates of the initial beliefs that are the most consistent with the data according
 10 to the chosen empirical criterion. However, as we will show in our simulation analysis in the
 11 next section, the effects of incorporating the initials into estimation routines will depend on the
 12 specification of the estimation objective.

13 3.4 Mixed approaches

14 Initialization methods can have several nuances that may not be, strictly speaking, reflected into
 15 the classes we proposed above. Particularly, there are many possibilities involving a mixture
 16 of the different approaches. For example, the REE-based initials could be computed on the
 17 basis of estimates of the model parameters obtained using data solely from the training sam-
 18 ple. A similar approach has been used in Slobodyan and Wouters (2012b), though adding the
 19 OLS-based method to the mixture: after estimates of the model under RE are obtained, using
 20 either the training or the whole sample, the initials are set to the REE-implied OLS estimates
 21 of agents' PLM. Another example is given in Milani (2011), where the mix is between the
 22 WLS-based method and the estimation-based approach: for every draw in the Bayesian esti-
 23 mation routine, a training sample of observations is used to compute the initial matrix of second
 24 moments according to (12), plugging in the corresponding estimated learning gain.

25 Another approach, proposed by Berardi and Galimberti (2017), relies on the use of Kalman
 26 smoothing within a sample of training data in order to accelerate the convergence of the WLS-
 27 based initialization method. Although this approach requires additional computations, it has
 28 been found to provide important speed improvements under alternative specifications of the
 29 learning mechanism, such as the Stochastic Gradient (SG) algorithm (Barucci and Landi, 1997;
 30 Evans and Honkapohja, 1998), which is obtained by replacing \mathbf{R}_t^{-1} by an identity matrix in the
 31 LS specification of (4). Hence, the SG does not benefit from the LS “normalization” step given

1 by the inverse of the matrix of second moments, which prevents the use of a diffuse prior on its
2 initialization.

3 Finally, although we have focused our discussion on the use of actual data on the variables
4 included in agents' PLM, another alternative is to use data from survey-based forecasts in order
5 to get information about the initial conditions. Data on survey forecasts have been broadly taken
6 as proxy for agents' actual expectations. In most of the cases, the initialization methods we
7 discussed above can be adjusted to take advantage of this information. For example, the REE-
8 based initials can be calculated using model's estimates obtained by replacing expectation terms
9 by direct measures from surveys (see, e.g., Orphanides and Williams, 2005a). Learning initials
10 consistent with surveys' information can also be obtained by adjusting the estimation-based
11 method to maximize the fit of the forecasts implied by the learning estimates to those obtained
12 from survey forecasts (Pfajfar and Santoro, 2010). Although we recognize the potential value
13 of these alternatives, we restrict the scope of our analysis to the definitions of initialization
14 methods covered by our classification.

15 4 Simulation Analysis

16 4.1 Baseline Phillips curve model

17 In order to shed further light on the comparative evaluation of the initialization methods re-
18 viewed above, we now analyze their quantitative properties with simulations. To provide a
19 meaningful economic example, we focus on a standard New Keynesian Phillips curve (NKPC)
20 model, given by

$$\pi_t = \beta\pi_{t+1}^e + \lambda x_t + u_t, \quad (20)$$

$$x_t = \rho x_{t-1} + v_t, \quad (21)$$

21 where π_t is inflation, π_{t+1}^e represents agents' expectations for next period inflation, x_t is a proxy
22 for real marginal cost, and u_t is a disturbance which can be interpreted as a measurement error
23 or as an unobserved cost-push shock. The parameters in (20) are taken as semi-structural in
24 the sense that they can be associated to deeper structural parameters of a microfounded model
25 (see, e.g., Mavroeidis et al., 2014). Particularly, β is the subjective discount factor and

$$\lambda = \frac{(1 - \theta)(1 - \theta\beta)}{\theta} \kappa, \quad (22)$$

26 where $\theta \in (0, 1)$ represents the fraction of firms that cannot change their prices in any given
27 period, i.e., an index of price rigidity under the Calvo framework, and $\kappa \leq 1$ is a function of
28 the labor elasticity of production and the price elasticity of demand.

1 The RE solution of this model is given by

$$\pi_t = \frac{\lambda}{1 - \beta\rho} x_t + u_t. \quad (23)$$

2 It can be shown that this equilibrium is E-stable if $\beta\rho < 1$, a condition that is automatically met
 3 under the usual assumptions that $0 < \beta < 1$ and $|\rho| < 1$ (see Evans and Honkapohja, 2001,
 4 pp. 198-200). Consistent with this solution, under adaptive learning agents form expectations
 5 according to a PLM given by

$$\pi_t = \phi_{t-1} x_t + z_t, \quad (24)$$

6 where ϕ_t is a parameter estimated with the univariate version of the LS algorithm given by
 7 (4)-(5), also substituting $y_t \equiv \pi_t$. Iterating (24) forward and substituting the expectations in
 8 (20) we obtain the ALM under learning

$$\pi_t = (\beta\rho\phi_{t-1} + \lambda) x_t + u_t. \quad (25)$$

9 4.2 Simulation and estimation approach

10 We generate 10,000 samples of artificial series of π_t and x_t assuming that $u_t \sim N(0, \sigma_u^2)$,
 11 $v_t \sim N(0, \sigma_v^2)$, and that $Correl(u_t, v_t) = 0$. The number of observations used for the learning
 12 initialization, in the training sample-based methods, and the estimation of the model's param-
 13 eters will be a dimension of our analysis, but in general we simulate the model for 10,000
 14 observations and assume the sample of data available for estimation starts at the 10,001 obser-
 15 vation, i.e., $t = 0$ is observation 10,000 in our artificial series. The model parameters are set
 16 to $\beta = 0.99$, $\theta = 0.65$, $\kappa = 0.25$, $\rho = 0.9$, $\sigma_u^2 = 3$, $\sigma_v^2 = 1$, whereas for the learning gain we
 17 evaluate two options, $\gamma_1 = 0.02$ and $\gamma_2 = 0.10^4$.

18 In order to estimate the model and learning parameters, we follow a generalized method of
 19 moments (GMM) approach. Following Chevillon et al. (2010), we obtain the moment condi-
 20 tions from the common assumption that the unobserved disturbance term, u_t in our model, is a
 21 martingale difference sequence, which means $E_{t-1}[u_t] = 0$. For a given set of pre-determined
 22 instruments, the model's ALM under learning is used to derive the associated GMM objective
 23 function, which is minimized using standard numerical estimation methods⁵.

24 One important dimension of our analysis of estimated initials is the choice of instruments
 25 used for the definition of the estimation criterion. According to Proposition 3, the use of a
 26 squared residuals criterion provides interesting results with respect to the accuracy of estimated
 27 initials as the estimation sample size grows. Nevertheless, the SSR criterion is not implicit in

⁴Our findings are qualitatively insensitive to these choices of parameters values, but not quantitatively. As evidenced in (10), though under a simpler model specification, the impact of initialization errors over the accuracy of the model's predictions is positively related to the magnitude of the parameter associated with the forward-looking term, β .

⁵Details about this estimation procedure are provided in Appendix B.

1 the GMM estimation approach with pre-determined instruments, and, as we will show, this can
2 turn the estimation approach susceptible to severe initialization biases for large samples. Hence,
3 we consider two estimation approaches. The first, which we denote as *standard*, is based on the
4 use of four lags of the model residuals as instrument; this approach was proposed by Chevillon
5 et al. (2010) in order to deal with the weak identification and persistence problems that are
6 induced by adaptive learning⁶. The second approach, which we call the *augmented* method, is
7 to use both four lags and the contemporaneous values of the model residuals as instruments, so
8 as to enforce the SSR criterion in the estimation of the model.

9 4.3 Results

10 We conduct two sets of simulation exercises comparing the accuracy of the initialization meth-
11 ods, focusing on the PLM coefficient initial⁷, ϕ_0 , and how the different initials affect the ac-
12 curacy of estimates of other model's parameters, such as the learning gain, γ , and the price
13 rigidity parameter, θ . First, we look at the case where only the learning initial is estimated, i.e.,
14 both γ and θ are assumed to be known. This corresponds to the analysis carried out analyti-
15 cally in section §3 for the simpler example model, and should therefore be informative about
16 how our previous conclusions extrapolate to more realistic models. Second, we consider the
17 more realistic case where the model parameters need to be estimated jointly with the determi-
18 nation of the learning initials. Importantly, in the latter case the pre-determined initials are also
19 conditioned on the estimates of the model parameters relevant to their determination.

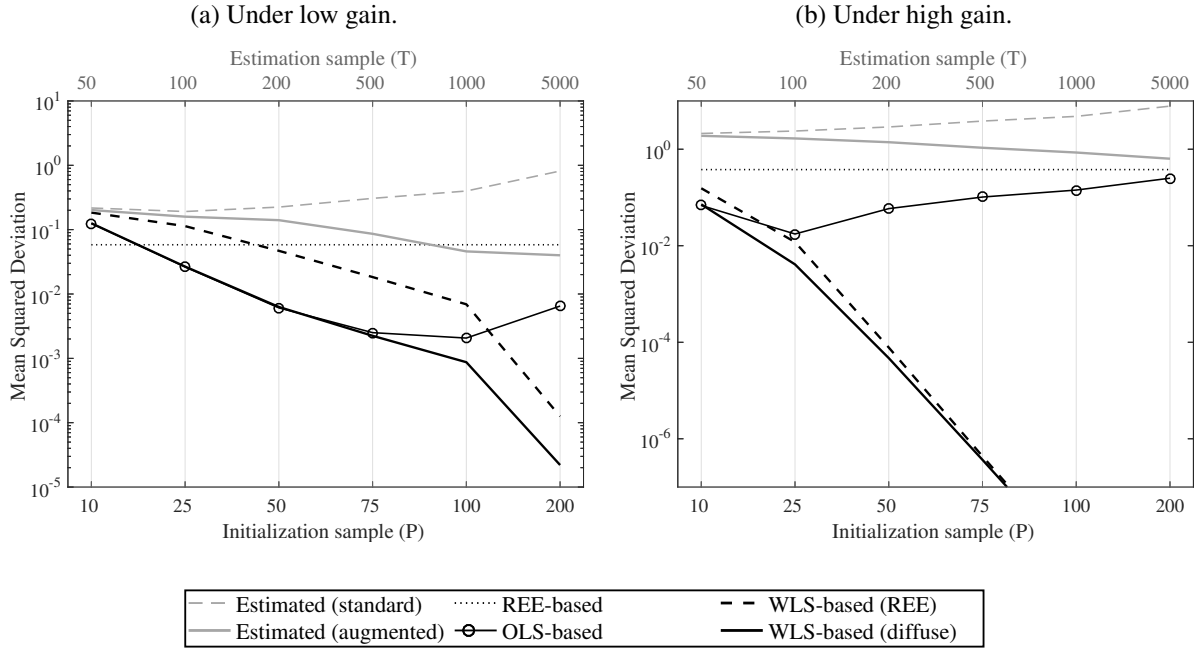
20 Under the REE-based method, based on the true model's parameter values, the initials
21 are given by $\phi_0^{REE} \simeq 0.44$, and $R_0^{REE} \simeq 5.26$. For the methods based on a training sam-
22 ple, i.e., the OLS and the WLS initials, we first set aside an initial portion of the simulated
23 series, $\{\pi_t, x_t\}_{t=1-P}^0$, and then compute the initials based on (12)-(13), varying the size of
24 the training sample over $P = \{10, 25, 50, 75, 100, 200\}$. For the WLS case we explore two
25 alternative specifications of the prior R_θ , namely, a REE prior with $R_\theta = R_0^{REE}$, and a dif-
26 fuse prior with $R_\theta = 0$. For the estimation-based initials we also consider two options: the
27 standard GMM approach, using four lags of the estimated disturbances as instrument; and
28 the augmented approach, which adds the contemporaneous disturbances to the former⁸. In
29 order to evaluate whether the sample size has any effect over the estimates we vary it over
30 $T = \{50, 100, 200, 500, 1000, 5000\}$.

⁶We have also experimented with the more traditional approach of taking lags of π_t and x_t as instruments, and found such alternatives to have no major effects over our main conclusions about initial estimates.

⁷Results on the accuracy of R_0 , which is also dependent on the initialization method, are available upon request.

⁸The initial of the matrix of second moment, R_0 , is fixed to its true value for the results reported on the estimated initials. Overall, adding R_0 to the estimation problem deteriorates the accuracy of the estimates for the PLM coefficient initial. These results are available upon request.

Figure 1: Mean squared deviations of initials with only ϕ_0 estimated.



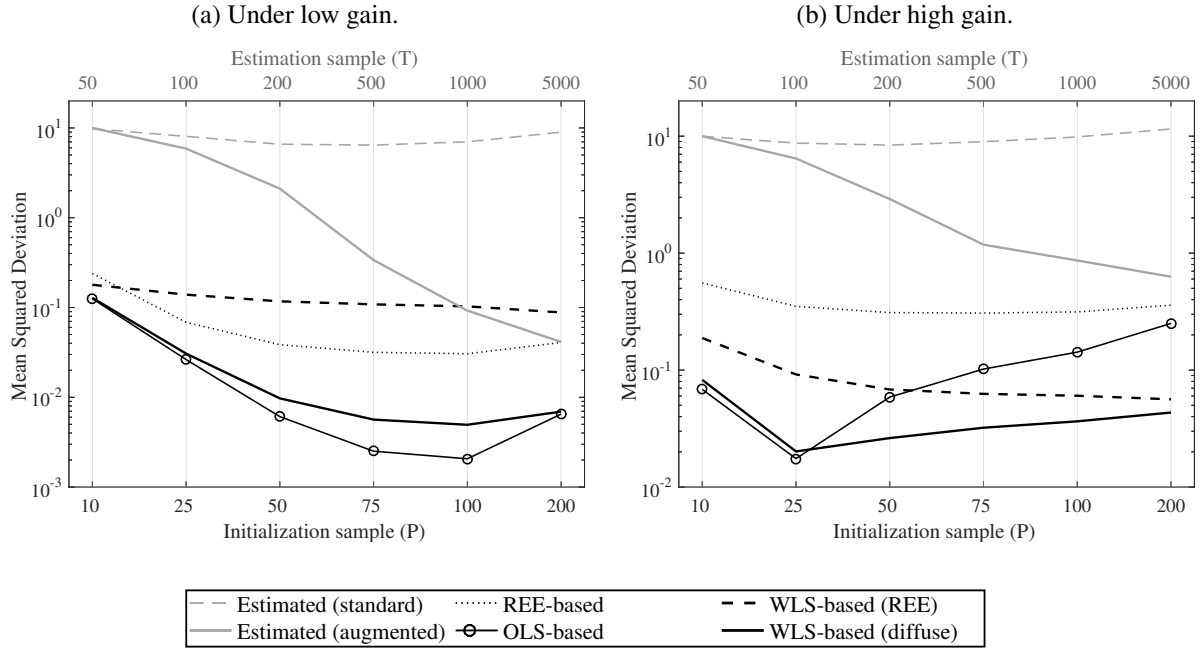
Notes: MSDs are depicted in a logarithmic scale (vertical axis) and are based on averaged statistics over 10,000 simulations of the baseline Phillips curve model. The low and high gains are set to $\gamma_1 = 0.02$ and $\gamma_2 = 0.10$, respectively. Estimated initials are obtained by GMM using the estimation sample indicated in the upper horizontal axis, whereas the training sample-based initials use the initialization samples indicated in the lower horizontal axis.

4.3.1 Initials accuracy under known model parameters

We begin looking at the case where only the learning initial is estimated. The MSDs for the initialization methods under analysis are presented in Figure 1 for the two different gain values. We draw the following observations from these results:

1. The REE-based initials are overall the less accurate among the pre-determined initials. Also, the performance of the REE-based initial deteriorates substantially for the higher gain calibration. This last result is in agreement with the result (see Eq. 8) that higher learning gains lead to noisier estimates of agents' PLM parameters, which drive out-of-equilibrium dynamics farther from the REE implied parameter values.
2. Whereas there is little difference between the OLS and the WLS initials under the lower gain calibration, the latter is clearly the best performing method under the higher gain. Also, notice that increasing the size of the training sample always improves the performance of the WLS initials, while the relationship is not monotonic for the OLS initials.
3. Between the two alternative specifications of the WLS prior on the learning coefficients uncertainty, the diffuse prior provides the most accurate initial estimates, though not in terms of the initial for the regressors' variance, where the REE prior obtains a better fit.

Figure 2: Mean squared deviations of initials with ϕ_0 , γ and θ jointly estimated.



Notes: See notes of Figure 1. The estimates of γ used for the WLS-based initials are those obtained under the augmented estimation approach with $T = 200$.

1 Hence, the use of a diffuse prior provides an interesting way to speed up the convergence
 2 of the learning estimates within the training sample.

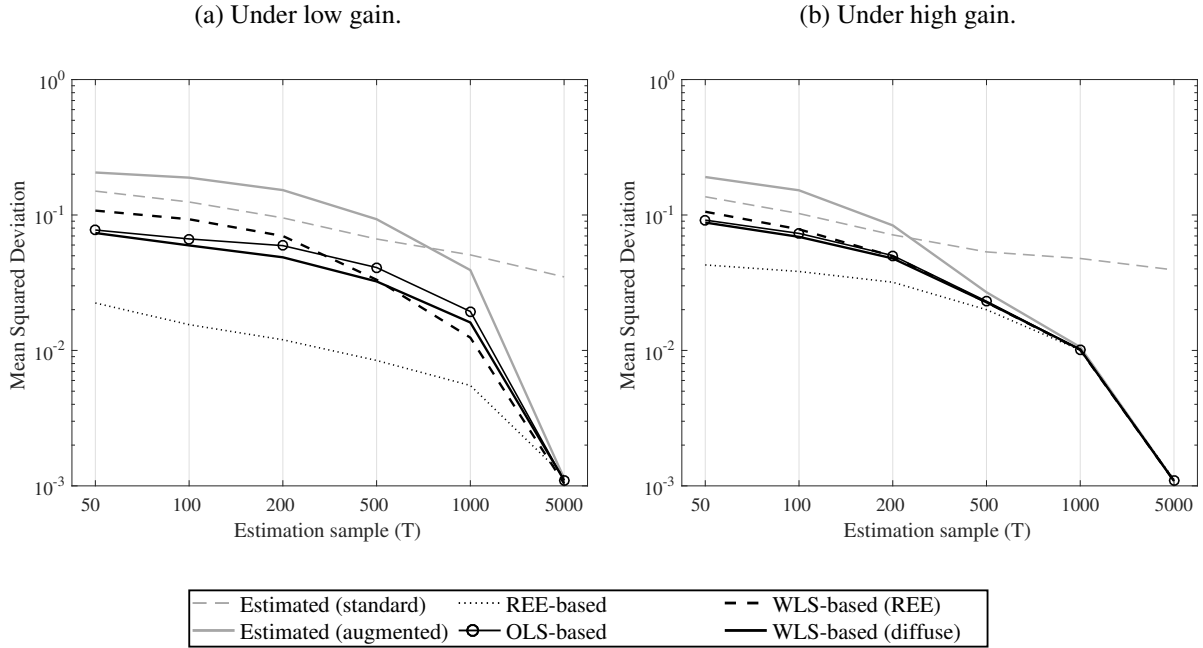
3 4. The results for the estimated initials are clearly dependent on the specification of the
 4 estimation criterion. Particularly, under the standard GMM approach larger estimation
 5 samples have a perverse effect on the accuracy of the estimated initials. Adding the
 6 squared residuals criterion to the estimation then brings significant improvements to the
 7 accuracy of this method, also avoiding the negative effects from the estimation sample
 8 size.

9 4.3.2 Initials accuracy with estimated model parameters

10 We now look at the more realistic case where model parameters need to be estimated together
 11 with the learning initial. Naturally, lack of knowledge about the true learning gain and the
 12 true index of price rigidity should affect the performance not only of the estimated initials, by
 13 adding more parameters to the estimation problem, but also the performance of those methods
 14 that are dependent on the model parameters, i.e., θ for the REE-based method, and γ for the
 15 WLS-based method. Figure 2 presents the results obtained under these circumstances⁹ and
 16 focusing on the same initialization methods and gain calibrations evaluated above.

⁹Results for the addition of each of these parameters separately to the estimation problem are presented in Appendix C.2.

Figure 3: Mean squared deviations of estimates of price rigidity index, θ .



Notes: See notes of Figure 1. The estimates refer to the exercise with the joint estimation of γ , θ and the corresponding initial. The OLS and the WLS-based initials use a training sample of $P = 75$ observations.

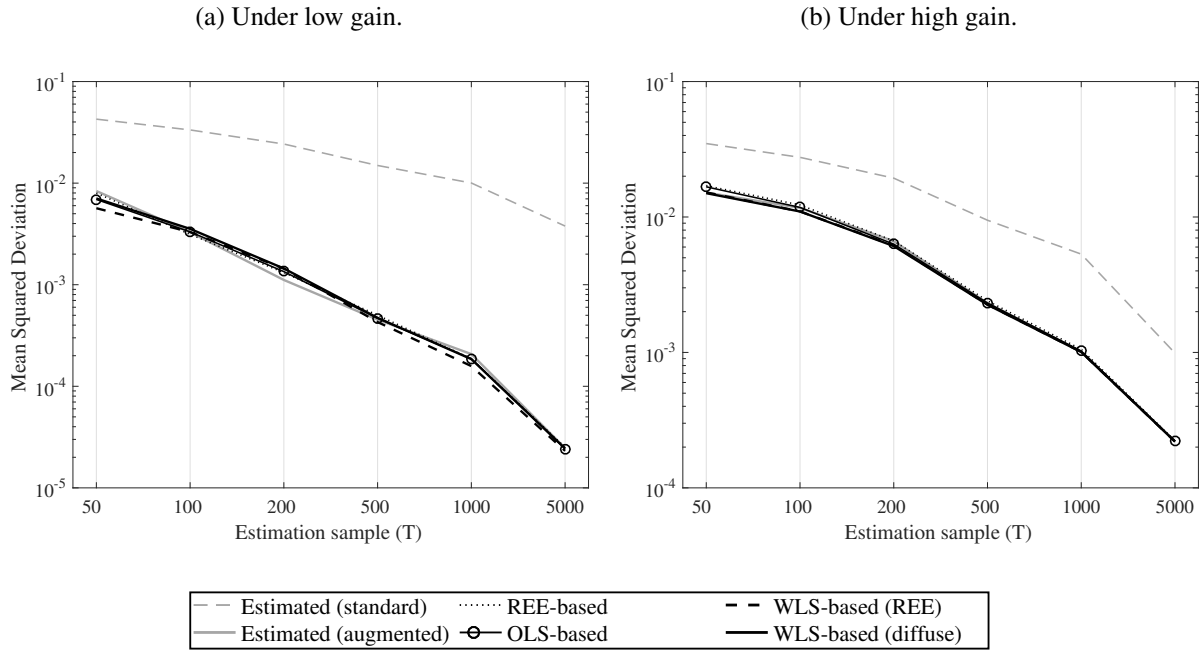
1 The main observation sprouting from these results relates to the sensitiveness of some pre-
 2 determined initials to the estimation of the model parameters. Particularly, we find that the
 3 inaccuracies in the estimation of γ can hold back the WLS convergence to the true learning
 4 coefficients within both the training sample and the estimation sample that follows the initial-
 5 ization. Besides, because the OLS-based method is not sensitive to any of the estimated model
 6 parameters, it yields performances identical to those observed without the estimation of other
 7 parameters. As in the previous exercise, the OLS-based initials can yield good convergence per-
 8 formance under small training samples, but they are not robust to variations in the underlying
 9 learning gain. With respect to the joint estimation of the initials we draw a similar conclusion
 10 to that obtained from the exercise under known model parameters: the augmented estimation
 11 approach is outperforming the standard approach.

12 4.3.3 Initials effects over accuracy of estimated model parameters

13 Perhaps a more interesting question for the applied researcher is how much the errors implied
 14 by each initialization method can affect the estimates of the other model parameters. In Fig-
 15 ures 3 and 4 we report results on the accuracy of the estimates for the parameters θ and γ ,
 16 respectively, also referring to the same simulations and estimations analyzed above.

17 Overall we find that all estimates of the model parameters tend to converge to their corre-
 18 sponding true values as the estimation sample size increases, confirming the consistency prop-
 19 erty of the adopted GMM estimators. However, we also observe different rates of convergence,

Figure 4: Mean squared deviations of estimates of learning gain, γ .



Notes: See notes of Figures 1 and 3.

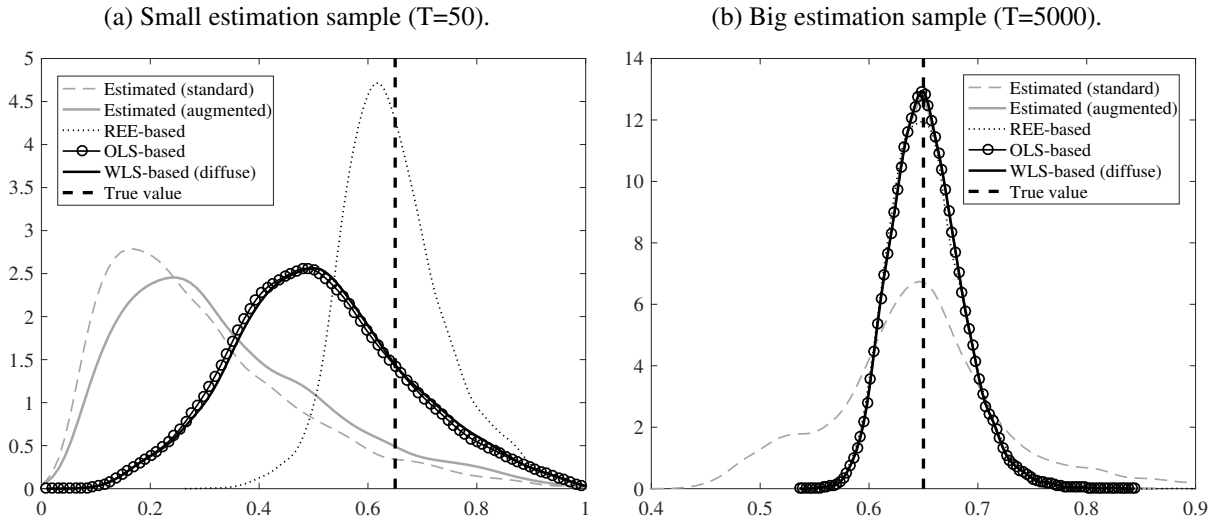
1 depending on the initialization method and the estimation approach. For the estimation of θ , the
 2 REE-based initial is performing remarkably well, whereas the training sample-based methods
 3 show a slower convergence depending on the underlying gain calibration. For the estimation of
 4 γ , the pre-determined initials are performing similarly.

5 But it is with respect to the estimation approach that we observe the greatest disparities.
 6 Whereas the augmented approach is generally preferred for large samples, its performance
 7 for the estimation of θ with smaller samples is slightly poorer relative to the standard approach.
 8 Although this result cautions against the use of non-predetermined instruments in the estimation
 9 of structural parameters in models with learning, it is also clear that our augmented approach
 10 brings substantial improvements to all estimates as the sample size grows¹⁰.

11 Nevertheless, the pre-determined methods are generally showing a better performance than
 12 that obtained under the joint estimation of the initials. This result is also evident in Figure 5,
 13 where we look at the densities of θ estimates when the gain is assumed to be known. Cor-
 14 roborating previous findings in the literature (e.g., Chevillon et al., 2010), in the left hand side
 15 (l.h.s.) panel of Figure 5 we observe that learning induces non-standard distributions in finite
 16 sample estimates of model parameters. Interestingly, we find that the joint estimation of the
 17 learning initials can lead to much stronger deviations from asymptotic distributions. Moreover,
 18 notice that even with a larger estimation sample, as in the right hand side (r.h.s.) panel of
 19 Figure 5, the θ estimates are still more dispersed under the case where the initial was jointly

¹⁰Another noteworthy caution against the augmented approach under small estimation samples is that numerical estimation instabilities were observed more frequently than in the standard approach for the estimation of the learning gain, as reported in Table 2 in the Appendix.

Figure 5: Densities of θ estimates over simulations.



Notes: Densities estimated using the normal kernel smoothing function over 10,000 simulations of the baseline Phillips curve model under the low gain setup, $\gamma_1 = 0.02$. Simulations with exact boundary estimates are discarded (see a report on these cases in the Appendix). The estimates refer to the exercise where only θ and ϕ_0 are estimated.

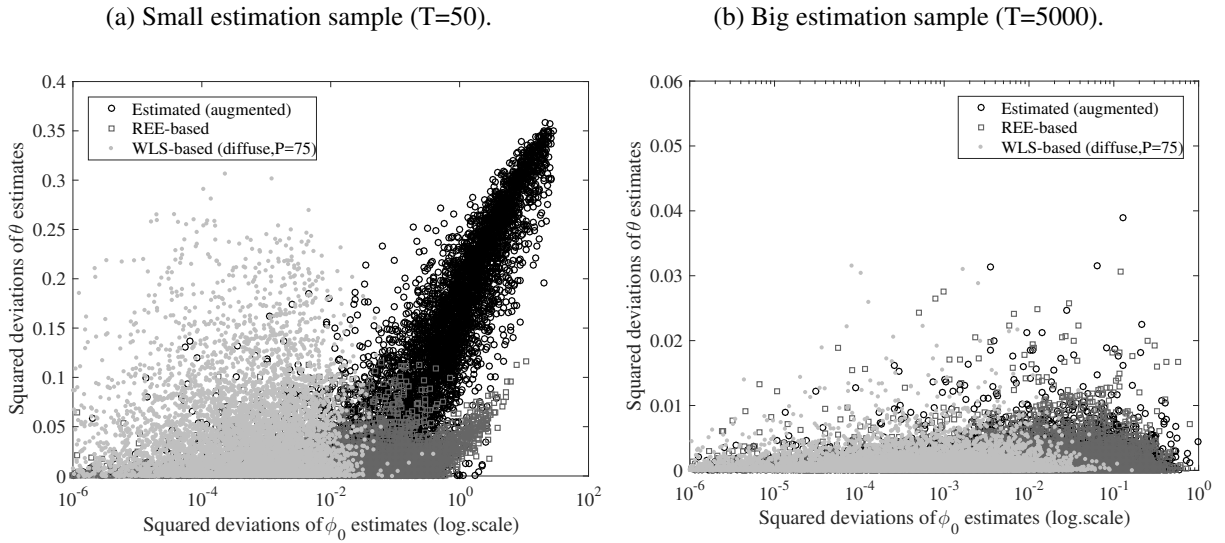
1 estimated using the standard approach than under the cases with pre-determined initials and the
 2 augmented estimation approach.

3 To better understand the impact of initialization errors over other model parameters esti-
 4 mation errors, we present in Figure 6 scatter diagrams relating the squared deviations of the
 5 initial estimates and those of θ estimates across the many simulations conducted with the joint
 6 estimation of these parameters. Importantly, we observe that the effects of the initialization
 7 errors depend on whether the initials are jointly estimated with the model parameter or pre-
 8 determined. Particularly for small estimation samples, in the l.h.s. panel of Figure 6, the joint
 9 estimation of the initial induces a positive correlation between the initials and the model pa-
 10 rameter's estimation errors. This negative spillover effect naturally disappears as the estimation
 11 sample increases, and both the model parameter and the estimated initials get more accurate.

12 4.4 Summary and discussion

13 We draw several important conclusions regarding the performance of the different initialization
 14 methods from our simulation analysis. First, our results under the assumption of known model
 15 parameters provided robustness to our analytical findings obtained under the simpler example
 16 model of the previous section. Particularly, we find that the WLS-based initials, with diffuse
 17 priors, show the best performance in terms of initials accuracy, as long as “enough” obser-
 18 vations are set aside for the training sample. The number of “enough” observations depends
 19 on the gain calibration, where smaller gains require longer training samples to guarantee the
 20 learning algorithm's convergence.

Figure 6: Relationship between initialization errors and model estimation errors.



Notes: The scatter diagrams are based on squared deviations between the indicated parameters estimates (ϕ_0 in the horizontal axis, and θ in the vertical one) and their corresponding true value over 10,000 simulations of the baseline Phillips curve model under the low gain setup, $\gamma_1 = 0.02$. Simulations with exact boundary estimates are discarded (see a report on these cases in the Appendix). The estimates refer to the exercise where only θ and ϕ_0 are estimated.

1 Second, when model parameters require estimation we found that the performance of the
 2 initialization methods is sensitive to the accuracy of the estimated parameters. However, the
 3 sensitiveness of pre-determined initials to the estimation approach was much weaker than that
 4 of the jointly estimated initials. Here, the traditional focus on pre-determined instruments was
 5 found to hinder the consistency of the initial estimates, which tended to deteriorate as the
 6 sample size increased. Our proposal to include a squared residual term in the estimation crite-
 7 rion found some success in reverting that relationship. One explanation for that result is that
 8 the SSR criterion penalizes more strongly the accumulation of learning estimation errors: the
 9 larger the estimation sample, the higher the precision attained by the learning estimates by the
 10 end of the sample, and with that, the more evident become the initialization errors (relative
 11 to end-of-sample errors). Without the SSR criterion, in contrast, the relative weights given to
 12 initialization errors only tended to decrease as the estimation sample size increased, reducing
 13 the relevance of initials and therefore turning their estimation indeterminate.

14 Notwithstanding, even under the augmented estimation approach the joint estimation of
 15 initials was found to have a negative spillover effect on the accuracy of model parameter esti-
 16 mates, particularly with small estimation samples. This spillover effect was not observed under
 17 the pre-determined initials. Hence, we conclude that their usage must be preferred; particularly,
 18 our evidence indicates that WLS-based method may be favored on the grounds of its greater
 19 robustness to uncertainties regarding the true values of the learning gain.

5 Empirical Application

5.1 Hybrid Phillips curve model

In order to evaluate the relevance of using different initialization methods, we now pursue an empirical application augmenting our baseline model and estimating it with US macroeconomic data. We follow Gali and Gertler (1999) and estimate a hybrid NKPC model given by

$$\pi_t = \psi_f \pi_{t+1}^e + \psi_b \pi_{t-1} + \delta x_t + \eta_t, \quad (26)$$

$$x_t = \rho x_{t-1} + \nu_t, \quad (27)$$

with

$$\psi_f = \beta \theta \zeta^{-1}, \quad \psi_b = \omega \zeta^{-1}, \quad \delta = (1 - \omega)(1 - \theta)(1 - \beta \theta) \zeta^{-1}, \quad (28)$$

$$\zeta = \theta + \omega(1 - \theta(1 - \beta)), \quad (29)$$

where ω is the fraction of firms that set their prices according to a backward looking rule of thumb, and the remaining parameters have the same interpretation as in the baseline specification.

Assuming the PLM is given by

$$\pi_t = a + b\pi_{t-1} + cx_t + u_t, \quad (30)$$

and that π_t is unknown when agents form expectations at period t ,

$$\pi_{t+1}^e = E_t[\pi_{t+1}] = a(1 + b) + b^2\pi_{t-1} + c(b + \rho)x_t. \quad (31)$$

The REE is given by the values of a , b , and c that solve the following equalities

$$a = \psi_f a(1 + b), \quad b = \psi_f b^2 + \psi_b, \quad c = \psi_f bc + \psi_f \rho c + \delta. \quad (32)$$

Solutions to this system are characterized in the Appendix A.5, where we also show that there is an unique E-stable REE in this model, given by $a^* = 0$, $b^* = (1 - \sqrt{1 - 4\psi_f\psi_b})/2\psi_f$, and $c^* = \delta/(1 - \psi_f(b^* + \rho))$. Under adaptive learning, agents form expectations using estimates of these parameters obtained according to the LS algorithm of (4)-(5), where $\hat{\phi}_t \equiv (\hat{a}_t, \hat{b}_t, \hat{c}_t)'$, $\mathbf{x}_t \equiv (1, \pi_{t-1}, x_t)'$ and $y_t \equiv \pi_t$.

5.2 Data and estimation approach

We use quarterly US data covering the period from 1947q1 to 2016q3. To measure inflation we use CPI inflation, whereas for the forcing variable, x_t , we use non-farm business sector labor

1 shares. To remove trends in the latter we obtain gap measures using the Hodrick-Prescott (HP)
2 filter. For simplicity, we are neglecting real-time data issues by focusing on a unique snapshot
3 of the realization of these series. All our data series are obtained from the FRED database of
4 the St. Louis Fed.

5 We adopt a GMM estimation approach similar to that used in our simulation analysis. Here,
6 the set of parameters determining inflation is given by $\Theta = (\beta, \theta, \omega, \rho, \gamma, \Phi_0, \mathbf{R}_0)$. To facilitate
7 estimation we fix some of these parameters to values we consider reasonable: $\beta = 0.99$, $\rho =$
8 0.60 , the latter corresponding (approximately) to the OLS estimate of (27). We experiment
9 with both estimation approaches discussed in our simulation exercise. Under the standard
10 estimation approach the instruments are given by a constant plus four lags of the estimation
11 residuals, the HP-filtered output gap on real GDP, and a long-short interest rate spread given
12 by the difference between the 10-year Treasury Bill rate and the Federal Funds rate. Under the
13 augmented approach we then add the contemporaneous estimation errors to the former list of
14 instruments¹¹.

15 In order to conduct a comparative evaluation of the estimates associated to the different
16 initialization methods, we fix the estimation sample to start from 1972q2. That leaves up to
17 100 observations for the training sample-based initializations. For the cases of the REE-based
18 and the WLS-based methods, the initials are re-calculated within the estimation routine in order
19 to take into account the estimates of the model parameters and the learning gain, respectively.
20 We use a diffuse prior for the WLS-based method.

21 We also benchmark our estimates of the model under learning with corresponding estimates
22 under RE. For that purpose we follow the approach of replacing the expectation term in (26)
23 by actual observations of next period inflation. Imposing the same identifying assumption we
24 used under learning, which in the present context is given by $E_{t-1}[\eta_t] = 0$, leads to moment
25 conditions of the form of (52), except that the residual term now also includes an one-step-
26 ahead inflation forecast error (see Mavroeidis et al., 2014, pp. 133-4). Finally, as instruments
27 for the RE estimation we use four lags of inflation and the labor income share gap instead of
28 the lagged residuals.

29 **5.3 Robust stable learning estimates**

30 Models with adaptive learning may generate unrealistic dynamics due to instabilities in the
31 recursive estimation of the learning coefficients. That is particularly relevant for the case of
32 constant gain learning: the higher the value of the learning gain, the higher the chances that
33 a sequence of abnormal forecast errors will push agents' PLM estimates outside the region

¹¹In a previous version of this paper (Berardi and Galimberti, 2016) we also included an unrestricted constant in the model specification, which turned out to result in unreasonable and less stable parameter estimates under learning. One possible explanation for this issue relates to the small sample sensitiveness of the GMM estimator to the normalization of orthogonality conditions with nonlinearities (see Fuhrer et al., 1995). Hence, in the current version we decided to include the unrestricted constant only for the estimation of the model under RE.

1 of stability in the parameters' space. Therefore, the estimation of models under learning can
2 be extremely sensitive to the sample of observations and the range of values allowed for the
3 learning gain.

4 To deal with these potential sources of instabilities we develop a robust stability estimation
5 approach. Namely, we impose additional constraints to the estimates of the learning gains so
6 as to ensure the stability of learning¹². Importantly, due to self-referentiality, these constraints
7 can depend on the values of the model parameters. Considering that we are also estimating θ
8 and ω we draw the upper bounds to be imposed on the estimation of the learning gain through
9 simulations of the model for different combinations of those parameters¹³. The resulting upper
10 bounds are presented in Figure 7. We observe that the model becomes more sensitive to the
11 learning gain as the model parameters get closer to their boundary values. To understand this
12 result notice that from (26): (i) as $\theta, \omega \rightarrow 1$, inflation converges to a random walk process
13 because $\delta \rightarrow 0$ and $\psi_f + \psi_b \rightarrow 1$; and, (ii) as $\theta, \omega \rightarrow 0$, the model becomes indeterminate
14 because $\zeta \rightarrow 0$.

15 Another mechanism that is often coupled to the learning algorithm in order to prevent the
16 learning estimates to diverge is the so-called "projection facility". Whenever the learning es-
17 timates leave a bounded region in the parameters space, the projection facility is activated in
18 order to contain the escape. This device is particularly useful for the cases of PLMs containing
19 lagged endogenous variables, where the definition of parameters' bounds can be easily derived
20 from stationarity conditions. Here, for the estimation of the hybrid NKPC we also adopt a pro-
21 jection facility that is activated whenever $|\hat{b}_t| \geq 1.5$, in which case the learning algorithm does
22 not update the learning estimates and discards the associated observation¹⁴.

23 5.4 Results

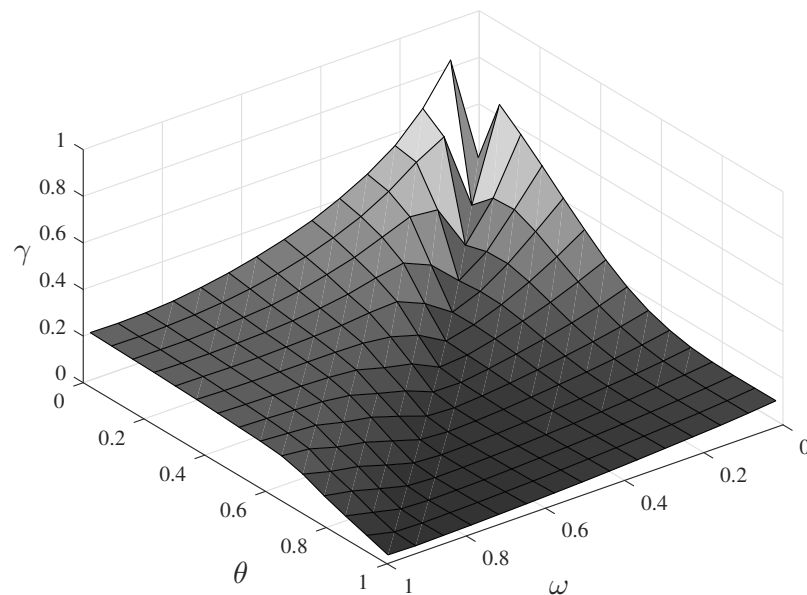
24 The results obtained with the GMM estimation of the hybrid NKPC model are presented in
25 Table 1. The RE benchmark estimates seem to be in accordance to previous estimates in the
26 literature (e.g., Gali and Gertler, 1999) where the Calvo's index of price stickiness, θ , is found
27 to be greater than the fraction of backward looking price setters, ω . This result is also robust
28 to the introduction of learning, although with some variation in the point estimates across the
29 estimation and initialization approaches. Particularly, we find that the degree of price stickiness

¹²Although we are not aware of any previous attempt to introduce such robustness conditions in the estimation of learning models in the literature, this approach is inspired by Evans and Honkapohja (2009), who showed that constant gain learning can turn an otherwise learnable equilibrium unstable, depending on the value of the gain. We thank an Associate Editor for pointing us towards this direction.

¹³See the Appendix C.3 for details about these simulations.

¹⁴Although stationarity of (30) requires that $|b| < 1$, we adopt a critical value of 1.5 so as to allow the algorithm some space to return to the stability region on its own.

Figure 7: Learning gain simulated upper bounds for hybrid New Keynesian Phillips curve model.



Notes: Upper bounds obtained from 1,000 simulations of model (26) under learning over a grid of 15 values of $\theta = \{0.01, 0.08, \dots, 0.91, 0.99\}$ and 15 values of $\omega = \{0.01, 0.08, \dots, 0.91, 0.99\}$. The reported upper bounds refer to averaged maximum gain values for which the simulated learning estimates of \hat{b}_t remained between ± 1.5 over 1,000 periods. See the Appendix C.3 for further details.

1 tended to decrease with the introduction of learning¹⁵, which is consistent with the findings of
2 Milani (2007).

3 We also find that the adaptive learning approach provides a superior fit to the data. That
4 is clearly evidenced by looking at the inflation expectations implied by the different model
5 estimates, depicted in Figure 8. These expectations are obtained from the assumed PLM, as in
6 (31), taking the learning initials as point of departure and updating these estimates according
7 to the learning algorithm and the corresponding data. As expected, the constant-gain learning
8 mechanism does a better job than the RE approach in tracking the fluctuations of inflation over
9 the sample of data.

10 Comparing the differences due to changes in the estimation and the initialization methods
11 we find that the augmented estimation approach provides more reliable estimates. Particularly,
12 we observe that there is much less variation on the estimates across the initializations when
13 the parameters are estimated under the augmented rather than the standard approach. E.g.,
14 ω varies between 0.01 (WLS- and OLS-based initial) and 0.47 (estimated initial) under the
15 standard estimation approach, while under the augmented approach ω is narrower identified
16 between 0.31 (estimated initial) and 0.44 (WLS-based initial); these variations have important
17 implications for the impact of the forcing variable on inflation, as evidenced by the implied δ
18 estimates in the last column of Table 1, which is a crucial parameter for policymakers. The same
19 applies to the other estimated parameters, except for the learning initials, where the estimates
20 of a_0 , b_0 , and c_0 still show some variation across the initialization methods.

21 Regarding the initials, empirically, it is difficult to assess which of the initialization methods
22 provided the most accurate initial estimates. In our simulation analysis we have argued in favor
23 of the WLS-based estimates for providing the estimates that are most coherent to the ongoing
24 learning process at the beginning of the estimation sample. That seems to be the case here, as
25 the OLS-based initials point to similar estimates.

26 In order to further illustrate the differences in the implied expectations across the initial-
27 ization methods, we re-estimate the model under varying sample beginnings and compare their
28 corresponding expectations in Figure 9. Particularly, we focus on augmented estimates obtained
29 with samples beginning around the Great Inflation episode in the US: (i) from 1966q4, captur-
30 ing the preceding period; (ii) from 1972q2, taken as the onset of the referred episode; and, (iii)
31 1985q1, leaving behind this period of greater volatility and focusing on the Great Moderation
32 period. Here we can clearly see the potential biases that the joint estimation of the learning ini-
33 tials can introduce in the modelled expectations: estimating the initials with data starting from
34 1972q2 leads to an expected inflation of more than 3% (quarterly rate), whereas the expecta-
35 tions implied by an ongoing process would be around 1% for that period. The expectations

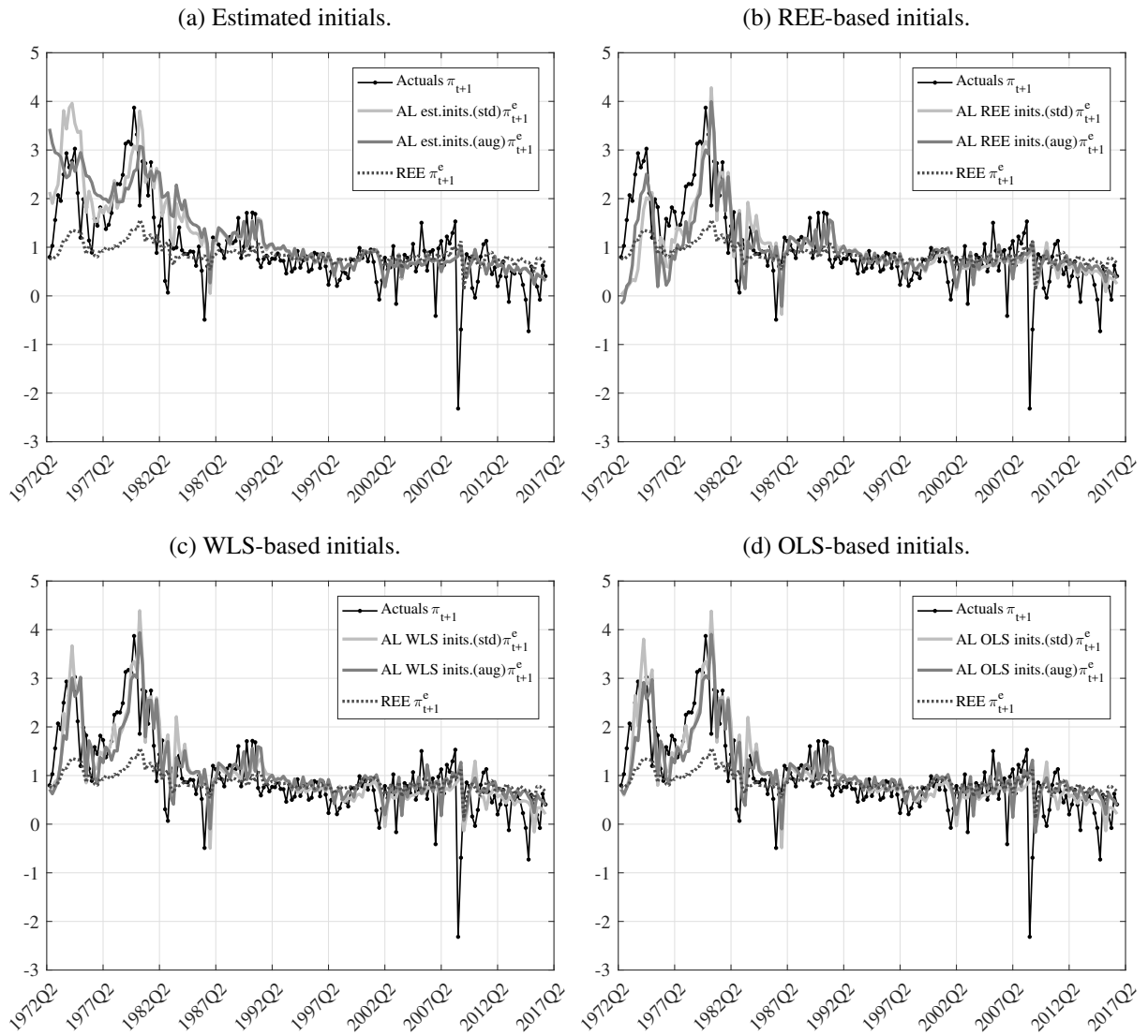
¹⁵The only exceptions to this finding occur under the standard estimation approach with jointly estimated and REE-based initials, but notice these estimates were not statistically significant. This finding is also robust to alternative estimation samples, as reported in Appendix C.4, although with substantial variation in the point estimates; this is not surprising considering that empirical estimates of the NKPC have been historically sensitive to data and estimation techniques (see Mavroudis et al., 2014).

Table 1: Empirical estimates of US NKPC - 1972q2 - 2016q3.

Exercise	Parameter estimates						Implied reduced form parameters		
	Structural		Learning				ψ_f	ψ_b	δ
	θ	ω	γ	a_0	b_0	c_0			
Rational expectations	0.85 (0.09)	0.40 (0.06)					0.68	0.32	0.01
Learning with jointly estimated initials									
- Standard	0.93 (1.07)	0.47 (0.71)	0.04 (0.03)	0.79 (0.54)	0.99 (0.32)	0.13 (0.26)	0.66	0.34	0.00
- Augmented	0.62 (0.07)	0.31 (0.11)	0.05 (0.02)	2.31 (2.02)	0.29 (0.89)	-0.37 (0.53)	0.66	0.33	0.11
Learning with REE-based initials									
- Standard	0.95 (0.58)	0.17 (0.10)	0.08 (0.02)	0.00	0.18	0.01	0.84	0.15	0.00
- Augmented	0.62 (0.11)	0.39 (0.06)	0.03 (0.02)	0.00	0.63	0.34	0.61	0.39	0.09
Under learning with WLS-based initials									
- Standard	0.74 (0.06)	0.01 (0.09)	0.11 (0.02)	0.53	0.45	0.12	0.98	0.01	0.09
- Augmented	0.56 (0.06)	0.44 (0.05)	0.02 (0.01)	0.17	0.76	0.04	0.56	0.44	0.11
Under learning with OLS-based initials									
- Standard	0.75 (0.07)	0.01 (0.09)	0.11 (0.02)	0.17	0.76	0.04	0.98	0.01	0.08
- Augmented	0.58 (0.06)	0.43 (0.06)	0.02 (0.01)	0.17	0.76	0.04	0.57	0.43	0.10

Notes: Parameters estimated by GMM, as explained in the text, using US data from 1972q2 to 2016q3. The implied reduced form parameters are obtained according to (28)-(29). For the cases with pre-determined initials, the initial learning coefficients are obtained from training samples of 75 observations over the period from 1953q3 to 1972q1. Values in parentheses are standard errors of the estimates and are computed based on numerical approximations of the objective function first derivatives. The standard errors under learning should be interpreted with caution since the estimators distribution, and corresponding test statistics, can become non-standard (see Chevillon et al., 2010). Estimates in bold are those found to be statistically different than zero at a 5% level of significance.

Figure 8: Expectations implied by model estimates from 1972q2 to 2016q3.



Notes: The inflation expectations are computed according to (31) departing from the estimates of the learning initials and the gain reported in Table 1. For the REE expectations we use the equilibrium values implied by the structural estimates, according to (32).

1 implied by the REE-based initials also present some deviation, which is consistent with our
2 previous finding that the REE fails to capture the dynamics of inflation over this period.

3 Summing up, we interpret these empirical results as providing further evidence in favor of
4 the use of training samples to initialize learning estimates together with the use of the aug-
5 mented estimation approach in order to estimate models with learning. The inclusion of the
6 SSR criterion seems to play a major role in facilitating the identification of the model parame-
7 ters, particularly for turning these estimates less sensitive to the learning initials. Irrespective of
8 the initialization method, all the augmented estimates of the hybrid NKPC tell a similar story:
9 when agents need to learn how to form their inflation expectations, our inference points to a
10 lower degree of price stickiness than suggested by estimates under the RE assumption. Further-
11 more, the augmented estimates do not reject the hybrid form of the NKPC, although the effect
12 of learning on the fraction of backward looking firms seems to depend on the initialization
13 adopted.

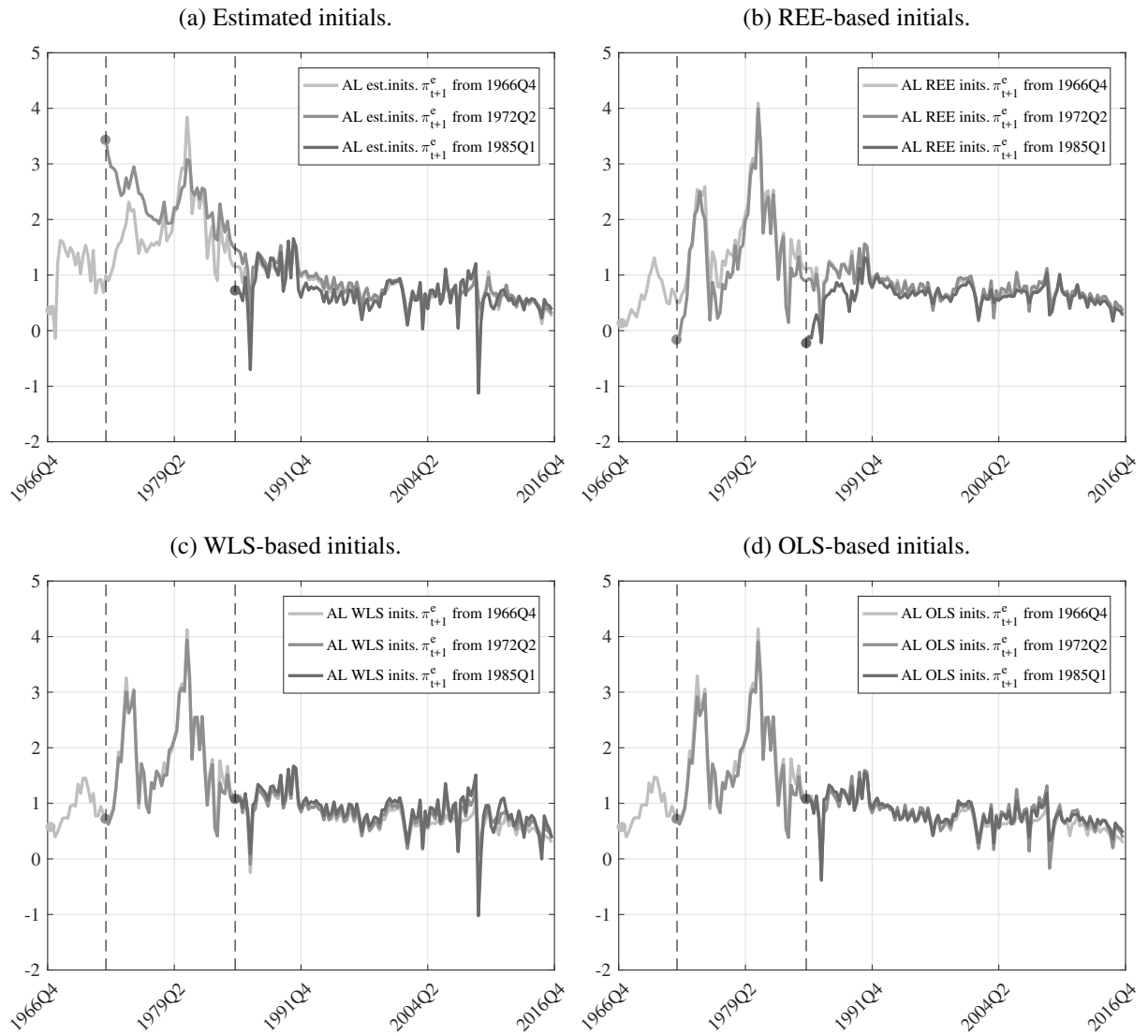
14 **6 Concluding remarks**

15 In this paper we provided a critical review on several methods previously proposed in the lit-
16 erature of learning and expectations in macroeconomics in order to initialize its learning algo-
17 rithms. We proposed a taxonomy of initialization methods that can be broadly defined in three
18 major classes: equilibrium-related methods, training sample-based methods, and estimation-
19 based methods. We conducted extensive simulation exercises comparing different initialization
20 methods that can be conceived within this classification.

21 Our analysis led us to draw the following recommendations. First, though equilibrium-
22 related initialization methods seem to provide rather conservative initials, they are often inco-
23 herent with the dynamics implied by learning, particularly under high learning gains. Second,
24 among the training sample-based methods, the use of standard OLS estimates can also turn out
25 to provide incoherent estimates since it does not take into account the particular specification
26 of the learning gain. Direct application of the learning algorithm into the training sample, the
27 WLS-based method in our terminology, was overall favored by our evaluation criteria. Particu-
28 larly, we found that a diffuse specification of this method leads to an accelerated convergence,
29 facilitating the feasibility of the method in macroeconomic contexts.

30 We also uncovered some problems with the approach of joint estimation of the learning
31 initials with other model parameters. Strikingly, the accuracy of these estimates tended to
32 deteriorate with the sample size, and, perhaps more importantly, we found evidence of spillover
33 effects from the biases introduced by estimation of the initials into the estimates of the model's
34 structural parameters. We proposed a solution to the first problem by including a squared
35 residual criterion to the estimation objective; whereas this approach restored consistency to the
36 joint estimation of the learning initials, attenuating the negative spillover effect asymptotically,
37 the estimation of initials was still problematic under small samples.

Figure 9: Expectations implied by augmented model estimates varying the sample beginnings.



Notes: The inflation expectations are computed as described in the notes of Figure 8, although based on parameter estimates obtained under the different sample beginnings as indicated in the legends. The estimates for these alternative samples are reported in Appendix C.4.

1 Finally, the sensitiveness of model estimates to the learning initialization method and the
 2 estimation approach was illustrated with an empirical application on the determination of US
 3 inflation rates under the New Keynesian Phillips curve framework. Our estimation approach
 4 was found to be the most robust to alternative specifications of the initials, with results in-
 5 dicating that the introduction of learning leads to a decrease in the estimated degree of price
 6 stickiness relative to what would otherwise be inferred under rational expectations. This find-
 7 ing is consistent with the idea that macroeconomic persistency is not only due to firms inability
 8 to adjust prices instantaneously, but also to the sluggish adjustment of their expectations in
 9 real-time.

10 **A Derivations**

11 **A.1 MSPE implied by initialization error**

12 From (2) we have that $y_{t+1} = \beta\phi_t + u_{t+1}$ and $\hat{y}_{t+1} = \beta\hat{\phi}_t$, so that the prediction error is given
 13 by

$$\begin{aligned}\hat{\Delta}_{t+1} &= y_{t+1} - \hat{y}_{t+1}, \\ &= \beta(\phi_t - \hat{\phi}_t) + u_{t+1}.\end{aligned}\tag{33}$$

14 Defining $\hat{\epsilon}_t = \phi_t - \hat{\phi}_t$, from (9) we find that

$$\begin{aligned}\hat{\epsilon}_t &= \phi_{t-1} + \gamma(y_t - \phi_{t-1}) - \hat{\phi}_{t-1} - \gamma(y_t - \hat{\phi}_{t-1}), \\ &= (1 - \gamma)(\phi_{t-1} - \hat{\phi}_{t-1}), \\ &= (1 - \gamma)\hat{\epsilon}_{t-1},\end{aligned}$$

15 which can be solved recursively to result in $\hat{\epsilon}_t = (1 - \gamma)^t \hat{\epsilon}_0$. Substituting this back into (33)
 16 and taking the expectation of the squared value results in

$$\begin{aligned}\hat{\Delta}_{t+1} &= \beta(1 - \gamma)^t \hat{\epsilon}_0 + u_{t+1}, \\ E[\hat{\Delta}_{t+1}^2] &= \beta^2(1 - \gamma)^{2t} \hat{\epsilon}_0^2 + \sigma_u^2.\end{aligned}$$

17 For the case where β is unknown, assuming $E[\hat{\beta}u_t] = 0$, the MSPE is then given by

$$\begin{aligned}E[\hat{\Delta}_{t+1}^2] &= E\left[\left(\beta\phi_t - \hat{\beta}\hat{\phi}_t + u_{t+1}\right)^2\right], \\ &= E[\hat{\epsilon}_t^2] + \sigma_u^2,\end{aligned}\tag{34}$$

1 where $\hat{\epsilon}_t = \beta\phi_t - \hat{\beta}\hat{\phi}_t$, leading to

$$\hat{\epsilon}_t = (1 - \gamma)\hat{\epsilon}_{t-1} + \gamma(\beta - \hat{\beta})y_t.$$

2 Solving this recursively yields

$$\hat{\epsilon}_t = (1 - \gamma)^t \hat{\epsilon}_0 + \gamma(\beta - \hat{\beta}) \sum_{i=0}^{t-1} (1 - \gamma)^i y_{t-i}.$$

3 **A.2 Long run variance of learning estimates**

4 Substituting (2) into (4) we obtain

$$\phi_t = \delta\phi_{t-1} + \gamma u_t,$$

5 where $\delta = 1 - \gamma(1 - \beta)$. This recursion is equivalent to

$$\phi_t = \delta^t \phi_0 + \gamma \sum_{i=0}^{t-1} \delta^i u_{t-i}, \quad (35)$$

6 In the limit, as $t \rightarrow \infty$, $E[\phi_\infty] = 0$ as long as $|\delta| < 1$. Hence, the long run variance of ϕ_t ,
7 denoted by $\bar{\sigma}_\phi^2$, is given by

$$\begin{aligned} \bar{\sigma}_\phi^2 &= \lim_{t \rightarrow \infty} E[\phi_t^2], \\ &= \lim_{t \rightarrow \infty} E \left[\delta^{2t} \phi_0^2 + 2\delta^t \phi_0 \gamma \sum_{i=0}^{t-1} \delta^i u_{t-i} + \left(\gamma \sum_{i=0}^{t-1} \delta^i u_{t-i} \right)^2 \right], \end{aligned}$$

8 which, because u_t is assumed to be serially independent, simplifies to

$$\begin{aligned} \bar{\sigma}_\phi^2 &= \lim_{t \rightarrow \infty} \delta^{2t} \phi_0^2 + \gamma^2 \sum_{i=0}^{t-1} \delta^{2i} \sigma_u^2, \\ &= \lim_{t \rightarrow \infty} \delta^{2t} \phi_0^2 + \gamma^2 \sigma_u^2 \frac{1 - \delta^{2t}}{1 - \delta^2}, \\ &= \frac{\gamma^2 \sigma_u^2}{1 - \delta^2}, \\ &= \frac{\gamma \sigma_u^2}{(1 - \beta)(1 + \delta)} \end{aligned}$$

9 where the limit is solved under the assumption that $|\delta| < 1$. Notice that $\partial \bar{\sigma}_\phi^2 / \partial \gamma > 0$, i.e., as
10 the gain increases the dispersion of the learning estimates tends to increase as well.

1 A.3 MSDs of training sample-based initials

2 For the training sample-based initials we consider a sample of P observations, and in this case
 3 (7) translates to

$$\phi_0 = \delta^P \phi_{-P} + \gamma \sum_{i=0}^{P-1} \delta^i u_{-i}. \quad (36)$$

4 A.3.1 WLS-based initials

The WLS-based initial is equivalent to the application of the constant gain LS algorithm to the training sample, departing from $\hat{\phi}_{-P} = 0$ (so-called diffuse prior¹⁶). Yet, because the actuals are determined by the true learning estimates, (36) can not be directly employed to obtain the WLS estimates. Taking that into account we find that the WLS-based initial is given by

$$\begin{aligned} \hat{\phi}_0^{WLS} &= \gamma\beta\lambda^{P-1}\phi_{-P} + \gamma\beta\lambda^{-1} \sum_{i=1}^{P-1} \lambda^i \phi_{-i} + \gamma \sum_{i=0}^{P-1} \lambda^i u_{-i}, \\ &= (\delta^P - \lambda^P) \phi_{-P} + \gamma \sum_{i=0}^{P-1} \delta^i u_{-i}, \end{aligned} \quad (37)$$

where $\lambda = 1 - \gamma$. The corresponding MSD is then given by

$$\begin{aligned} \mathcal{D}_{\phi_0}^{WLS} &= E \left[(\lambda^P \phi_{-P})^2 \right], \\ &= \lambda^{2P} E [\phi_t^2]. \end{aligned} \quad (38)$$

5 A.3.2 OLS-based initials

The OLS-based initial is obtained by application of the Ordinary Least Squares estimator to the training sample of observations, which is equivalent to the use of (6) with a decreasing gain, $1/t$, instead of the constant, γ , and also departing from the diffuse prior, $\hat{\phi}_{-P} = 0$. The OLS-based initial is then given by

$$\begin{aligned} \hat{\phi}_0^{OLS} &= \frac{\beta}{P} \phi_{-P} + \frac{\beta}{P} \sum_{i=1}^{P-1} \phi_{-i} + \frac{1}{P} \sum_{i=0}^{P-1} u_{-i}, \\ &= \frac{\beta (\delta^P - 1)}{P (\delta - 1)} \phi_{-P} + \frac{1}{P (1 - \beta)} \sum_{i=1}^{P-1} (1 - \beta \delta^i) u_{-i} + \frac{1}{P} u_0, \end{aligned} \quad (39)$$

6 from which the corresponding MSD follows as

$$\mathcal{D}_{\phi_0}^{OLS} = \left(\delta^P - \frac{\beta (\delta^P - 1)}{P (\delta - 1)} \right)^2 E [\phi_t^2] + \sigma_u^2 \mathcal{K}, \quad (40)$$

¹⁶These derivations can be easily extended to the case with any other arbitrary prior, such as the REE-based prior

1 where

$$\mathcal{K} = \left(\gamma - \frac{1}{P} \right)^2 + \sum_{i=1}^{P-1} \left(\gamma \delta^i + \frac{\beta \delta^i - 1}{P(1-\beta)} \right)^2. \quad (41)$$

2 **A.4 SSR estimation-based initials and associated MSD**

3 The SSR estimation-based initials are obtained according to

$$\hat{\phi}_0^{SSR} = \arg \min_{\hat{\phi}_0} SSR = \sum_{i=0}^{T-1} \hat{u}_{T-i}^2. \quad (42)$$

4 Using (2) we find that the FOC to this minimization problem is given by

$$\beta \sum_{i=0}^{T-1} \frac{\partial \hat{\phi}_{T-i-1}}{\partial \hat{\phi}_0} \hat{\phi}_{T-i-1} = \sum_{i=0}^{T-1} \frac{\partial \hat{\phi}_{T-i-1}}{\partial \hat{\phi}_0} y_{T-i}. \quad (43)$$

5 Using (6) we find that, departing from an estimated initial and using the data available in the
6 estimation sample,

$$\hat{\phi}_{T-i-1} = \lambda^{T-i-1} \hat{\phi}_0 + \gamma \sum_{j=0}^{T-i-2} \lambda^j y_{T-i-j-1}, \quad (44)$$

7 so that

$$\frac{\partial \hat{\phi}_{T-i-1}}{\partial \hat{\phi}_0} = \lambda^{T-i-1}, \quad (45)$$

8 where $\lambda = 1 - \gamma$. Substituting (44) and (45) into (43), we find that the initial estimate is given
9 by

$$\hat{\phi}_0^{SSR} = \sum_{i=0}^{T-1} \mathcal{H}_i y_{T-i}, \quad (46)$$

10 where

$$\mathcal{H}_i = \frac{\gamma (\lambda^{-i} (1 + \lambda) - \beta \lambda (\lambda^{-i} - \lambda^i))}{\beta \lambda (\lambda^{-T} - \lambda^T)}. \quad (47)$$

11 The MSD associated to the estimated initials can be derived by using the ALM to obtain

$$y_{T-i} = \beta \delta^{T-i-1} \phi_0 + \beta \gamma \sum_{j=0}^{T-i-2} \delta^j u_{T-i-j-1} + u_{T-i}, \quad (48)$$

which introduced into (46) leads to

$$\begin{aligned} \hat{\phi}_0^{SSR} &= \phi_0 \beta \delta^{T-1} \sum_{i=0}^{T-1} \mathcal{H}_i \delta^{-i} + \beta \gamma \delta^{-1} \sum_{i=1}^{T-1} u_{T-i} \delta^i \sum_{j=0}^{i-1} \mathcal{H}_j \delta^{-j} + \sum_{i=0}^{T-1} \mathcal{H}_i u_{T-i}, \\ &= \phi_0 + \frac{\gamma (\gamma - 2)}{\beta \lambda (\lambda^T - \lambda^{-T})} \sum_{i=0}^{T-1} \lambda^{-i} u_{T-i}. \end{aligned} \quad (49)$$

This expression can be used to evaluate the MSD according to

$$\begin{aligned}
\mathcal{D}_\phi^{SSR} &= E \left[\left(\phi_0 - \hat{\phi}_0^{SSR} \right)^2 \right], \\
&= E \left[\phi_0^2 \right] - 2E \left[\phi_0 \hat{\phi}_0 \right] + E \left[\hat{\phi}_0^2 \right], \\
&= E \left[\hat{\phi}_0^2 \right] - E \left[\phi_0^2 \right], \\
&= \frac{\gamma(\gamma - 2)}{\beta^2 (\lambda^{2T} - 1)} \sigma_u^2.
\end{aligned} \tag{50}$$

1 A.5 Hybrid Phillips curve REEs

2 There are six solutions to the RE conditions in (32). Starting with b , there are two possible
3 solutions given by $b^\pm = \frac{1 \pm \sqrt{1 - 4\psi_f \psi_b}}{2\psi_f}$. For a there are three possibilities: $a = 0$ or a is indeter-
4 minate with $\beta = 1$ and $\omega \neq 1$, or $\omega = 1$ and $\beta \neq 1$. Finally, c is uniquely determined by b .
5 Putting these combinations together we have the following RE solutions:

$$RE.1 = \{a = 0, b^+, c^+\}; RE.2 = \{a = 0, b^-, c^-\};$$

$$6 RE.3 = \{a = any, b^+, c^+, \beta = 1, \omega \neq 1\}; RE.4 = \{a = any, b^+, c^+, \omega = 1, \beta \neq 1\};$$

$$7 RE.5 = \{a = any, b^-, c^-, \beta = 1, \omega \neq 1\}; RE.6 = \{a = any, b^-, c^-, \omega = 1, \beta \neq 1\}.$$

8 Naturally, the solutions with indeterminate a can be ruled out on economic grounds: the sub-
9 jective discount factor is generally assumed to be smaller than unity, whereas $\omega = 1$ would
10 break the relationship between π_t and x_t underlying the Phillips curve.

11 To check for E-stability of these solutions we first define the T -mapping associated to this
12 model:

$$T \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{bmatrix} \psi_f a (1 + b) \\ \psi_f b^2 + \psi_b \\ \psi_f b c + \psi_f c \rho + \delta \end{bmatrix}.$$

13 E-stability requires that the eigenvalues of the Jacobian matrix of T , evaluated at the given RE
14 solution, are smaller than unity. These eigenvalues depend only on the value of b and are given
15 by

$$\{(1 + b) \psi_f, 2b\psi_f, (b + \rho) \psi_f\}.$$

16 Focusing on the range of reasonable parameter values, $0 < \beta < 1$, $0 < \theta < 1$, $0 \leq \omega < 1$, and
17 $-1 < \rho < 1$, we find that only the RE solution with b^- can be E-stable.

B GMM estimation

Moment conditions are obtained using the model's ALM under learning, (25), from which we can define the residual function according to

$$h_t(\Theta) = \pi_t - \beta\rho\phi_{t-1}x_t - \lambda x_t, \quad (51)$$

where Θ denotes the set of parameters requiring estimation. For a given set of instruments, Z_t , the corresponding moment conditions are given by

$$E[Z_t h_t(\Theta)] = 0. \quad (52)$$

The model parameters are then estimated by minimization of the associated GMM objective function

$$g_T(\hat{\Theta}) = \left[T^{-1} \sum_{t=1}^T Z_t h_t(\hat{\Theta}) \right]' W_T \left[T^{-1} \sum_{t=1}^T Z_t h_t(\hat{\Theta}) \right], \quad (53)$$

which is constructed from the sample counterpart of the moment conditions in (52) and a weighting matrix, W_T . This weighting matrix is optimally defined as a consistent estimator of the inverse of the long-run variance of the moment conditions. Because the variance of (52) depends on the values of Θ , we adopt an iterative GMM estimator (see, e.g., Hall, 2005): we first set $W_T^{(0)} = I$ to obtain the preliminary estimates $\hat{\Theta}^{(0)}$ that minimize (53); we then use the Newey and West (1987) heteroskedasticity and autocorrelation consistent estimator of the variance of the moment conditions evaluated at $\hat{\Theta}^{(0)}$ to obtain a new estimate of $W_T^{(1)}$; we repeat this process until a convergence criterion is achieved.

A numerical optimization routine is used to find the values of $\hat{\Theta}$ that minimize the GMM objective function, (53). For that purpose we adopt a sequential quadratic programming algorithm, namely the 'sqp' option in the fmincon function in Matlab optimization toolbox. The convergence criterion for the iterative estimation of the weighting matrix is based on the Euclidean distance between the successive parameters estimates, i.e., $\|\hat{\Theta}^{(i)} - \hat{\Theta}^{(i-1)}\| < \epsilon$. In our simulations we set $\epsilon = 10^{-4}$, for which convergence is achieved in about 4 (5) iterations, on average, under $\gamma = 0.02$ ($\gamma = 0.10$).

Whereas the model parameters are reasonably constrained by theory implied boundaries, the parameters associated to the learning algorithm require artificial constrains to avoid numerical instabilities during estimation. Our experimental analysis led us to adopt the following constrains: $0 \leq \gamma \leq 0.5$, $\phi^{REE} - 5 \leq \hat{\phi}_0 \leq \phi^{REE} + 5$, and $0 < \hat{R}_0 \leq 50$. Although these constraints were never violated in the artificial data, the numerical estimation of $\hat{\phi}_0$ and \hat{R}_0 often resulted in boundary solutions. These cases are summarized in Table 2, where we observe that increasing the number of estimated parameters tends to increase the number of boundary solutions for the initial estimates. These effects are also amplified when the data true learning gain increases, except for the case of the estimation of the gain itself, where a higher gain facil-

1 itates estimation. The number of boundary solutions also tends to increase with the sample size
2 under the standard estimation approach, whereas it decreases under the augmented approach.
3 Perhaps the main drawback associated to the augmented approach relates to the estimation of
4 the learning gain, where the estimates turned out unstable for more than 80% of the simulations
5 under the small learning gain and using the small estimations samples.

6 **C Supplementary results**

7 **C.1 Validation of analytical MSD expressions**

8 In order to validate the MSD expressions derived in section 3, we ran 100,000 simulations of
9 the example model of section 2 for different combinations of parameters. For each simulation
10 we draw $10,000 + T$ (pseudo-)random values for u_t from a normal distribution, and use these
11 disturbances to generate artificial observations of y_t and ϕ_t according to the ALM and the
12 learning algorithm specification given by (2) and (6), respectively. We then obtain estimates
13 for $\phi_{t=10,000}$ according to each initialization method described in section 3. The results are
14 presented in table 3 and show a very good fit between our analytical derivations and numerical
15 estimates.

16 **C.2 Initials MSDs with jointly estimated model parameters**

17 Figures 10 and 11 report the initials MSD results for the addition of γ and θ , respectively, to
18 the estimation problem. These results are consistent with our findings discussed on the main
19 text on the basis of Figure 2 for the exercise where both γ and θ are jointly estimated with the
20 learning initial.

21 **C.3 Simulated upper bounds on learning gain for hybrid Phillips curve 22 model**

23 The upper bounds on the learning gain adopted for the empirical exercise are obtained by
24 running 1,000 simulations of the hybrid NKPC model, using the same parametrizations that we
25 are adopting in the estimation, i.e., fixing $\beta = 0.99$ and $\rho = 0.60$, and using RE estimated
26 residuals in order to fix the variances of the residuals needed to simulate the model. The
27 simulations are conducted over a grid of 15 values of $\theta = \{0.01, 0.08, \dots, 0.91, 0.99\}$ and 15
28 values of $\omega = \{0.01, 0.08, \dots, 0.91, 0.99\}$, and increasing the learning gain starting from $\gamma = 0.005$
29 with a step equal to 0.005 until the learning estimates diverge: 1,000 observations are simulated
30 with the first gain, and then if the estimates did not diverge, we move to the next γ , and so on.
31 We check for divergence by looking at the value of the \hat{b}_t learning coefficient estimate, and stop
32 the simulation whenever it goes above the value of 1.5. In general, when the learning estimate

Table 2: Number of simulations where estimates hit lower/upper bounds.

Estimation sample	Standard Estimation Experiments				Augmented Estimation Experiments											
	1	2	3	4	1	2	3	4								
(T)	$\hat{\phi}_0$	$\hat{\gamma}$	$\hat{\theta}$	$\hat{\phi}_0$	$\hat{\gamma}$	$\hat{\theta}$	$\hat{\phi}_0$	$\hat{\gamma}$	$\hat{\theta}$							
Data generated with $\gamma = 0.02$:																
50	0	1030	4151	224	4554	2801	2743	4196	0	462	8137	218	4562	3382	5341	4678
100	0	1060	3028	3	4279	2036	2488	3947	2	156	8326	15	4016	1536	6895	4627
200	0	1080	2007	0	3567	1483	1930	3586	3	46	7487	2	2918	271	7343	4507
500	0	1203	886	0	2467	1465	1004	2883	4	7	3987	5	1216	7	4934	3456
1000	0	1323	284	0	1724	1657	331	2368	0	2	1207	0	340	2	2032	1676
5000	0	1941	0	0	807	2318	0	1352	0	0	1	0	0	0	2	1
Data generated with $\gamma = 0.10$:																
50	149	1436	3287	135	3968	2842	2428	4108	216	903	6211	268	3903	3412	4113	4564
100	157	1599	1548	136	3429	2270	1489	3814	197	609	3881	189	3276	1924	3384	4106
200	197	1787	323	175	2879	2099	503	3384	173	356	1314	160	2331	679	1569	3080
500	335	2163	5	297	2062	2339	47	2858	102	142	135	103	997	144	200	1275
1000	479	2540	0	454	1372	2695	2	2562	57	71	14	62	288	70	21	354
5000	1282	3644	0	1092	922	3443	0	2387	16	15	0	17	0	16	0	1

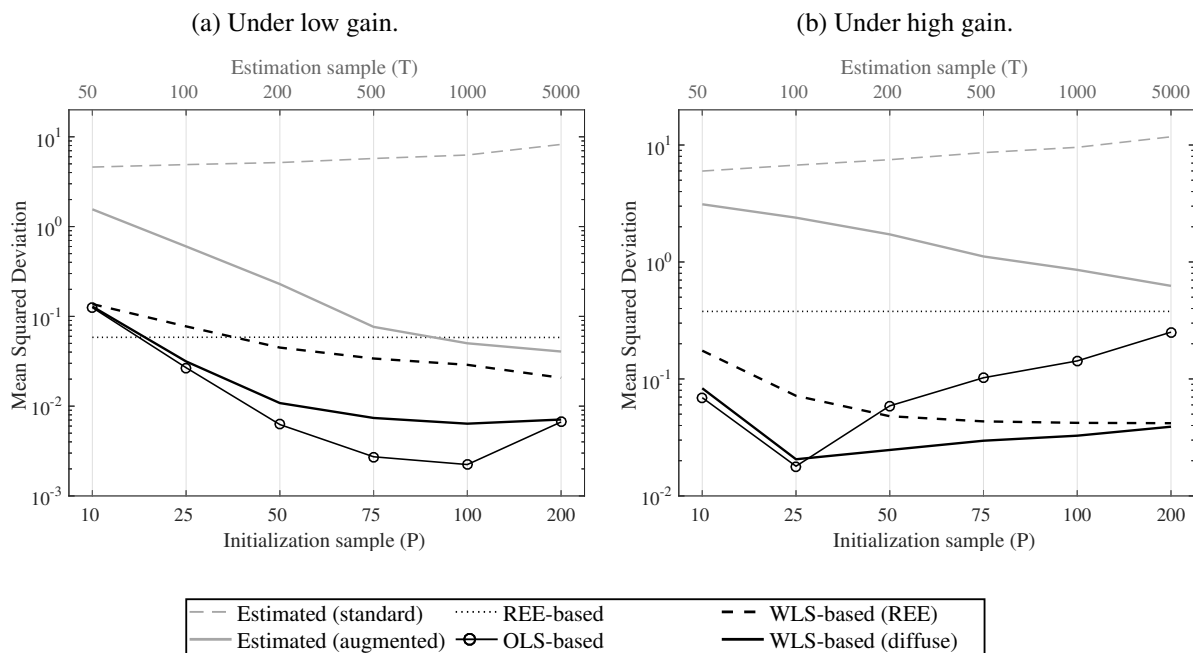
Notes: Statistics based on 10,000 simulations of the benchmark Phillips curve model under learning.

Table 3: Analytical and numerical MSDs for initials obtained under the example model.

Method	$\sigma_u^2 = 1$				$\sigma_u^2 = 4$			
	$\beta = 0.9$		$\beta = -5$		$\beta = 0.9$		$\beta = -5$	
	$\gamma \Rightarrow$	0.02	0.25	0.02	0.25	0.02	0.25	0.02
REE-based:								
	0.100	1.269	0.002	0.083	0.400	5.082	0.007	0.331
	(0.100)	(1.266)	(0.002)	(0.083)	(0.400)	(5.063)	(0.007)	(0.333)
WLS-based:								
$P = 10$	0.067	0.004	0.001	0.000	0.267	0.016	0.005	0.001
	(0.067)	(0.004)	(0.001)	(0.000)	(0.267)	(0.016)	(0.005)	(0.001)
$P = 25$	0.037	0.000	0.001	0.000	0.145	0.000	0.003	0.000
	(0.036)	(0.000)	(0.001)	(0.000)	(0.146)	(0.000)	(0.003)	(0.000)
$P = 50$	0.013	0.000	0.000	0.000	0.053	0.000	0.001	0.000
	(0.013)	(0.000)	(0.000)	(0.000)	(0.053)	(0.000)	(0.001)	(0.000)
$P = 100$	0.002	0.000	0.000	0.000	0.007	0.000	0.000	0.000
	(0.002)	(0.000)	(0.000)	(0.000)	(0.007)	(0.000)	(0.000)	(0.000)
OLS-based:								
$P = 10$	0.079	0.066	0.051	0.045	0.317	0.261	0.205	0.180
	(0.079)	(0.065)	(0.051)	(0.045)	(0.316)	(0.261)	(0.204)	(0.182)
$P = 25$	0.022	0.284	0.008	0.063	0.089	1.128	0.033	0.252
	(0.022)	(0.283)	(0.008)	(0.063)	(0.089)	(1.132)	(0.032)	(0.253)
$P = 50$	0.006	0.580	0.002	0.072	0.025	2.304	0.007	0.288
	(0.006)	(0.577)	(0.002)	(0.072)	(0.025)	(2.308)	(0.007)	(0.290)
$P = 100$	0.003	0.905	0.001	0.077	0.013	3.609	0.003	0.308
	(0.003)	(0.903)	(0.001)	(0.078)	(0.013)	(3.612)	(0.003)	(0.311)
SSR estimation-based:								
$T = 10$	0.147	0.543	0.005	0.017	0.591	2.152	0.019	0.070
	(0.147)	(0.542)	(0.005)	(0.018)	(0.588)	(2.167)	(0.019)	(0.070)
$T = 100$	0.050	0.541	0.002	0.017	0.199	2.146	0.006	0.070
	(0.050)	(0.540)	(0.002)	(0.018)	(0.199)	(2.160)	(0.006)	(0.070)
$T = 250$	0.049	0.541	0.002	0.017	0.196	2.146	0.006	0.070
	(0.049)	(0.540)	(0.002)	(0.018)	(0.196)	(2.160)	(0.006)	(0.070)
$T = 1000$	0.049	0.541	0.002	0.017	0.196	2.146	0.006	0.070
	(0.049)	(0.540)	(0.002)	(0.018)	(0.196)	(2.160)	(0.006)	(0.070)

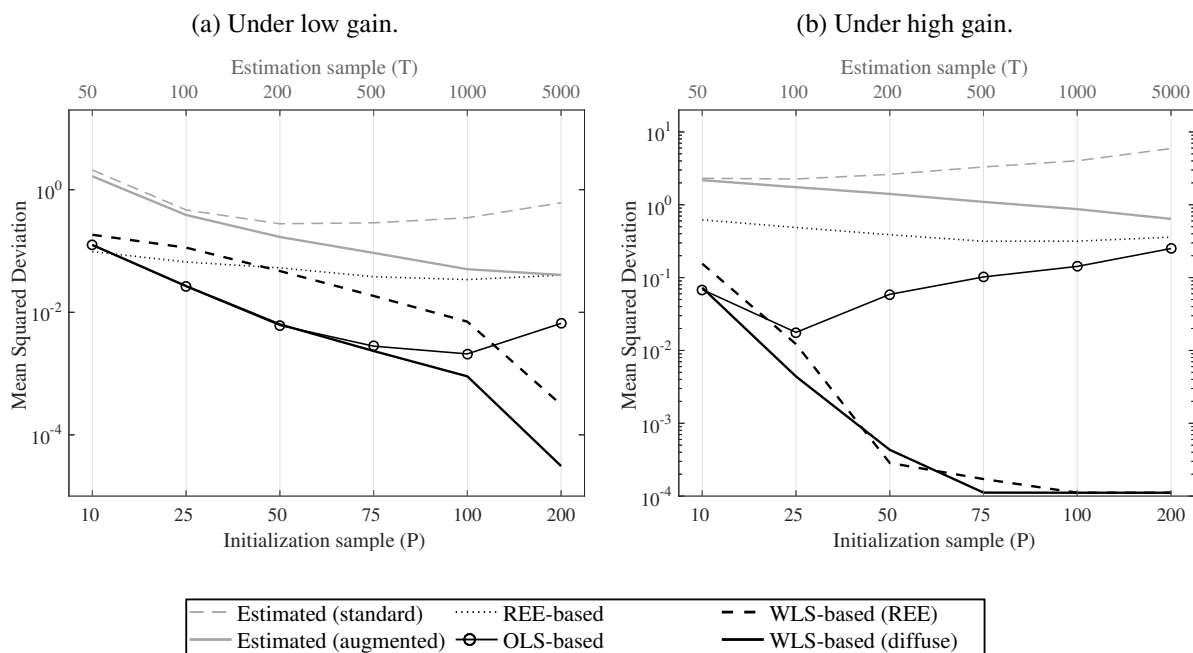
Notes: All statistics refer to MSD estimates. Values in parentheses come from the analytical expressions derived in the text, whereas those without parentheses come from the numerical simulation.

Figure 10: Mean Squared Deviations of initials with ϕ_0 and γ jointly estimated.



Notes: See notes of Figure 1. The estimates of γ used for the WLS-based initials are those obtained under the augmented estimation approach with $T = 200$.

Figure 11: MSDs - ϕ_0 and θ jointly estimated.



Notes: See notes of Figure 1. The estimates of θ used for the REE-based initials are those obtained under the augmented estimation approach.

1 reached that value, the simulated inflation series was already exploding. The upper bounds are
2 then obtained by averaging the maximum gain values across the 1,000 simulations for each
3 combination of parameters. For the empirical estimation of the model, we obtain an interpolant
4 function, $G(\theta, \omega)$, using cubic splines based on the gridded gain upper bounds presented in
5 Figure 7. This function is then supplied to the numerical estimation routine as a nonlinear
6 restriction on the estimates.

7 **C.4 Empirical estimates for alternative sample**

8 Estimates of macroeconomic models may be sensitive to the historical period underlying the
9 sample used for estimation. U.S. post-WWII inflation history over our sample of data has
10 been marked by two main episodes of interest: (i) the Great Inflation period, starting around
11 early-1970s and ending by mid-1980s; and, (ii) the more recent period that followed the 2007-
12 08 financial crisis. To check for the sensitiveness of our results, we re-estimated the model
13 discussed in the text by focusing on samples that exclude these periods. Results are presented
14 in Tables 4 and 5, the former excluding the Great Inflation period and the latter also excluding
15 the period since the financial crisis. Finally, we also re-estimated the model with an earlier
16 sample beginning in order to check for the sensitiveness of the implied inflation expectations at
17 the start of the Great Inflation period. These results are presented in Table 6.

Table 4: Empirical estimates of US NKPC - 1985q1 - 2016q3.

Exercise	Parameter estimates						Implied reduced form parameters		
	Structural		Learning				ψ_f	ψ_b	δ
	θ	ω	γ	a_0	b_0	c_0			
Rational expectations									
	0.86	0.37					0.70	0.30	0.01
	(0.11)	(0.06)							
Learning with jointly estimated initials									
- Standard	0.87	0.92	0.05	3.29	-0.02	-3.73	0.48	0.52	0.00
	(3.21)	(3.39)	(0.02)	(1.16)	(0.52)	(1.13)			
- Augmented	0.99	0.09	0.00	1.12	-1.10	-0.35	0.91	0.08	0.00
	(7.07)	(1.18)	(0.01)	(0.79)	(0.68)	(0.70)			
Learning with REE-based initials									
- Standard	0.77	0.17	0.08	0.00	0.22	0.15	0.81	0.18	0.05
	(0.10)	(0.08)	(0.02)						
- Augmented	0.65	0.30	0.06	0.00	0.46	0.33	0.68	0.32	0.09
	(0.08)	(0.04)	(0.02)						
Under learning with WLS-based initials									
- Standard	0.47	0.37	0.07	0.30	0.77	-0.04	0.56	0.44	0.21
	(0.03)	(0.05)	(0.04)						
- Augmented	0.66	0.53	0.00	0.31	0.81	0.02	0.55	0.45	0.05
	(0.11)	(0.06)	(0.00)						
Under learning with OLS-based initials									
- Standard	0.99	0.01	0.09	0.31	0.81	0.02	0.98	0.01	0.00
	(2.48)	(0.17)	(0.03)						
- Augmented	0.64	0.53	0.01	0.31	0.81	0.02	0.54	0.45	0.05
	(0.08)	(0.07)	(0.01)						

Notes: See notes to Table 1.

Table 5: Empirical estimates of US NKPC - 1985q1 - 2007q4.

Exercise	Parameter estimates						Implied reduced form parameters		
	Structural		Learning				ψ_f	ψ_b	δ
	θ	ω	γ	a_0	b_0	c_0			
Rational expectations									
	0.91	0.20					0.81	0.18	0.01
	(0.13)	(0.06)							
Learning with jointly estimated initials									
- Standard	0.59	0.10	0.01	1.19	-0.74	6.70	0.85	0.15	0.22
	(0.38)	(0.14)	(0.00)	(0.14)	(0.09)	(1.01)			
- Augmented	0.33	0.57	0.04	4.43	-0.39	-3.77	0.36	0.63	0.22
	(0.05)	(0.07)	(0.01)	(0.99)	(0.25)	(1.43)			
Learning with REE-based initials									
- Standard	0.95	0.17	0.08	0.00	0.18	0.01	0.84	0.15	0.00
	(0.56)	(0.11)	(0.02)						
- Augmented	0.69	0.25	0.07	0.00	0.35	0.26	0.73	0.27	0.08
	(0.09)	(0.06)	(0.02)						
Under learning with WLS-based initials									
- Standard	0.66	0.51	0.00	0.31	0.81	0.02	0.56	0.44	0.05
	(0.22)	(0.14)	(0.01)						
- Augmented	0.66	0.31	0.06	0.31	0.78	-0.04	0.68	0.32	0.08
	(0.06)	(0.07)	(0.04)						
Under learning with OLS-based initials									
- Standard	0.97	0.94	0.00	0.31	0.78	-0.04	0.51	0.49	0.00
	(15.87)	(15.78)	(0.01)						
- Augmented	0.66	0.34	0.06	0.31	0.78	-0.04	0.65	0.34	0.08
	(0.06)	(0.06)	(0.03)						

Notes: See notes to Table 1.

Table 6: Empirical estimates of US NKPC - 1966q4 - 2016q3.

Exercise	Parameter estimates						Implied reduced form parameters		
	Structural		Learning				ψ_f	ψ_b	δ
	θ	ω	γ	a_0	b_0	c_0			
Rational expectations									
	0.84	0.40					0.67	0.32	0.01
	(0.07)	(0.06)							
Learning with jointly estimated initials									
- Standard	0.96	0.29	0.06	1.45	0.11	0.56	0.76	0.23	0.00
	(1.66)	(0.70)	(0.03)	(2.01)	(1.71)	(0.81)			
- Augmented	0.72	0.22	0.07	4.89	-1.16	0.03	0.76	0.23	0.07
	(0.06)	(0.17)	(0.02)	(3.03)	(0.60)	(2.21)			
Learning with REE-based initials									
- Standard	0.78	0.04	0.10	0.00	0.05	0.15	0.94	0.05	0.06
	(0.07)	(0.09)	(0.02)						
- Augmented	0.72	0.38	0.05	0.00	0.53	0.17	0.65	0.35	0.05
	(0.11)	(0.07)	(0.01)						
Under learning with WLS-based initials									
- Standard	0.77	0.03	0.10	0.22	0.56	-0.06	0.95	0.04	0.07
	(0.06)	(0.09)	(0.02)						
- Augmented	0.85	0.35	0.05	0.23	0.50	-0.04	0.70	0.29	0.01
	(0.17)	(0.06)	(0.02)						
Under learning with OLS-based initials									
- Standard	0.77	0.01	0.11	0.23	0.50	-0.04	0.98	0.01	0.07
	(0.06)	(0.09)	(0.02)						
- Augmented	0.88	0.32	0.06	0.23	0.50	-0.04	0.73	0.27	0.01
	(0.21)	(0.08)	(0.01)						

Notes: See notes to Table 1.

1 **References**

2 Barucci, E. and L. Landi (1997). Least mean squares learning in self-referential linear stochastic
3 models. *Economics Letters* 57(3), 313–317.

4 Berardi, M. and J. K. Galimberti (2013). A note on exact correspondences between adaptive
5 learning algorithms and the kalman filter. *Economics Letters* 118(1), 139–142.

6 Berardi, M. and J. K. Galimberti (2014). A note on the representative adaptive learning algo-
7 rithm. *Economics Letters* 124(1), 104 – 107.

8 Berardi, M. and J. K. Galimberti (2016, December). On the Initialization of Adaptive Learning
9 in Macroeconomic Models. Technical report, KOF Working Papers, Zürich.

10 Berardi, M. and J. K. Galimberti (2017). Smoothing-based initialization for learning-to-
11 forecast algorithms. *Macroeconomic Dynamics* in press.

12 Bray, M. M. and N. E. Savin (1986). Rational expectations equilibria, learning, and model
13 specification. *Econometrica* 54(5), 1129–1160.

14 Bullard, J. and S. Eusepi (2005). Did the great inflation occur despite policymaker commitment
15 to a taylor rule? *Review of Economic Dynamics* 8(2), 324 – 359.

16 Carboni, G. and M. Ellison (2009). The great inflation and the greenbook. *Journal of Monetary*
17 *Economics* 56(6), 831 – 841.

18 Carceles-Poveda, E. and C. Giannitsarou (2007). Adaptive learning in practice. *Journal of*
19 *Economic Dynamics and Control* 31(8), 2659–2697.

20 Carceles-Poveda, E. and C. Giannitsarou (2008). Asset pricing with adaptive learning. *Review*
21 *of Economic Dynamics* 11(3), 629 – 651.

22 Chevillon, G., M. Massmann, and S. Mavroeidis (2010). Inference in models with adaptive
23 learning. *Journal of Monetary Economics* 57(3), 341–351.

24 Eusepi, S. and B. Preston (2011, October). Expectations, learning, and business cycle fluctua-
25 tions. *American Economic Review* 101(6), 2844–72.

26 Evans, G. W. and S. Honkapohja (1998). Stochastic gradient learning in the cobweb model.
27 *Economics Letters* 61(3), 333–337.

28 Evans, G. W. and S. Honkapohja (2001). *Learning and expectations in macroeconomics*. Fron-
29 tiers of Economic Research. Princeton, NJ: Princeton University Press.

- 1 Evans, G. W. and S. Honkapohja (2009, 06). Robust Learning Stability with Operational Mone-
2 tary Policy Rules. In K. Schmidt-Hebbel, C. E. Walsh, N. L. S. Editor), and K. S.-H. (Series
3 (Eds.), *Monetary Policy under Uncertainty and Learning*, Volume 13 of *Central Banking,*
4 *Analysis, and Economic Policies Book Series*, Chapter 5, pp. 145–170. Central Bank of Chile.
- 5 Fuhrer, J. C., G. R. Moore, and S. D. Schuh (1995). Estimating the linear-quadratic inventory
6 model maximum likelihood versus generalized method of moments. *Journal of Monetary*
7 *Economics* 35(1), 115 – 157.
- 8 Gali, J. and M. Gertler (1999, October). Inflation dynamics: A structural econometric analysis.
9 *Journal of Monetary Economics, Elsevier* 44(2), 195–222.
- 10 Gaus, E. and S. Ramamurthy (2014, August). Estimation of constant gain learning models.
11 Working Papers 12-01, Ursinus College, Department of Economics.
- 12 Hall, A. R. (2005). *Generalized Method of Moments*. Advanced Texts in Econometrics. Oxford
13 University Press.
- 14 Huang, K. X., Z. Liu, and T. Zha (2009). Learning, adaptive expectations and technology
15 shocks. *The Economic Journal* 119(536), 377–405.
- 16 Ljung, L. and T. Soderstrom (1983). *Theory and Practice of Recursive Identification*. The MIT
17 Press.
- 18 Lubik, T. A. and C. Matthes (2014, January). Indeterminacy and learning: An analysis of mon-
19 etary policy in the great inflation. Working Paper 14-2, Federal Reserve Bank of Richmond.
- 20 Marcet, A. and J. P. Nicolini (2003). Recurrent hyperinflations and learning. *American Eco-*
21 *nomics Review* 93(5), 1476–1498.
- 22 Mavroeidis, S., M. Plagborg-Møller, and J. H. Stock (2014). Empirical evidence on inflation
23 expectations in the new keynesian phillips curve. *Journal of Economic Literature* 52(1),
24 124–88.
- 25 Milani, F. (2007, October). Expectations, learning and macroeconomic persistence. *Journal of*
26 *Monetary Economics* 54(7), 2065–2082.
- 27 Milani, F. (2008). Learning, monetary policy rules, and macroeconomic stability. *Journal of*
28 *Economic Dynamics and Control* 32(10), 3148 – 3165.
- 29 Milani, F. (2011). Expectation shocks and learning as drivers of the business cycle. *The Eco-*
30 *nomics Journal* 121(552), 379–401.
- 31 Milani, F. (2014). Learning and time-varying macroeconomic volatility. *Journal of Economic*
32 *Dynamics and Control* 47(0), 94 – 114.

- 1 Newey, W. K. and K. D. West (1987, May). A simple, positive semi-definite, heteroskedasticity
2 and autocorrelation consistent covariance matrix. *Econometrica* 55(3), 703–08.
- 3 Ormeño, A. and K. Molnár (2015, June). Using survey data of inflation expectations in the
4 estimation of learning and rational expectations models. *Journal of Money, Credit and Bank-*
5 *ing* 47(4), 673–699.
- 6 Orphanides, A. and J. C. Williams (2005a, November). The decline of activist stabilization pol-
7 icy: Natural rate misperceptions, learning, and expectations. *Journal of Economic Dynamics*
8 *and Control* 29(11), 1927–1950.
- 9 Orphanides, A. and J. C. Williams (2005b). Inflation scares and forecast-based monetary policy.
10 *Review of Economic Dynamics* 8(2), 498 – 527.
- 11 Pfajfar, D. and E. Santoro (2010). Heterogeneity, learning and information stickiness in infla-
12 tion expectations. *Journal of Economic Behavior & Organization* 75(3), 426–444.
- 13 Primiceri, G. E. (2006). Why inflation rose and fell: Policy-makers’ beliefs and u. s. postwar
14 stabilization policy. *The Quarterly Journal of Economics* 121(3), 867–901.
- 15 Sargent, T., N. Williams, and T. Zha (2006). Shocks and government beliefs: The rise and fall
16 of american inflation. *American Economic Review* 96(4), 1193–1224.
- 17 Sargent, T. J. (1999). *The Conquest of American Inflation*. Princeton, NJ: Princeton University
18 Press.
- 19 Slobodyan, S. and R. Wouters (2012a). Learning in a medium-scale dsge model with ex-
20 pectations based on small forecasting models. *American Economic Journal: Macroeco-*
21 *nomics* 4(2), 65–101.
- 22 Slobodyan, S. and R. Wouters (2012b). Learning in an estimated medium-scale dsge model.
23 *Journal of Economic Dynamics and Control* 36(1), 26 – 46.
- 24 Williams, N. (2003, January). Adaptive learning and business cycles. Mimeo.