

# Model Comparisons in Unstable Environments

Raffaella Giacomini and Barbara Rossi

(*UCL and CeMMAP*)      (*ICREA – Univ. Pompeu Fabra,*  
*Barcelona GSE and CREI*)

January 8, 2015

## Abstract

The goal of this paper is to develop formal tests to evaluate the relative in-sample performance of two competing, misspecified, non-nested models in the presence of possible data instability. Compared to previous approaches to model selection, which are based on measures of global performance, we focus on the local relative performance of the models. We propose tests that are based on different measures of local performance and that correspond to different null and alternative hypotheses. The empirical application provides insights into the time variation in the performance of a representative DSGE model of the European economy relative to that of VARs.

**Keywords:** Model Selection Tests, Misspecification, Structural Change, Kullback-Leibler Information Criterion

**Acknowledgments:** We are grateful to D. Kristensen for helpful comments and M. del Negro, F. Smets, R. Wouters, W.B. Wu and Z. Zhao for sharing their codes. We also thank the editor, three anonymous referees, seminar participants at the Empirical Macro Study Group at Duke U., Atlanta Fed, UC Berkeley, UC Davis, U. of Michigan, NYU Stern, Boston U., U. of Montreal, UNC Chapel Hill, U. of Wisconsin, UCL, LSE, UCL, Stanford’s 2006 SITE workshop, the 2006 Cleveland Fed workshop, the 2006 Triangle Econometrics workshop, the Fifth ECB Workshop, the 2006 Cass Business School Conference on Structural Breaks and Persistence, the 2007 NBER Summer Institute, the 2009 NBER-NSF Time Series Conference and the 2012 AEA Meetings for useful comments and suggestions. Support by NSF grants 0647627 and 0647770 is gratefully acknowledged.

**J.E.L. Codes:** C22, C52, C53

# 1 Introduction

The problem of detecting time-variation in the parameters of econometric models has been widely investigated for several decades, and empirical applications have documented that structural instability is widespread.

In this paper, we depart from the literature by focusing on investigating instability in the performance of models, rather than focusing solely on instability in their parameters. The idea is simple: in the presence of structural change, it is plausible that the performance of a model may itself be changing over time, even if the model's parameters remain constant. In particular, when the problem is that of comparing the performance of competing models, it would be useful to understand which model performed better at which point in time.

The goal of this paper is therefore to develop formal techniques for conducting inference about the relative performance of two models over time, and to propose tests that can detect time variation in relative performance even when the parameters are constant. Existing model selection tests such as Rivers and Vuong (2002) are inadequate for answering this question, since they work under the assumption that there exists a globally best model. The central idea of our method is instead to propose a measure of the models' local relative performance: the "local relative Kullback-Leibler Information Criterion" (local relative KLIC), which represents the relative distance of the two (misspecified) likelihoods from the true likelihood at a particular point in time. We then investigate ways to conduct inference about the local relative KLIC and construct tests of the joint null hypothesis that the relative performance and the parameters of the models are constant over time.

We propose two tests, which correspond to different assumptions about the parameters and the relative performance under the null and alternative hypotheses: 1) a "one-time reversal" test against a one-time change in models' performance and parameters; and 2) a "fluctuation test" against smooth changes in both performance and parameters. The first test is based on estimating the parameters and the relative performance before and after potential change dates, whereas the second is based on nonparametric estimates of local performance and local parameters. The fluctuation test is based on a fixed-bandwidth approximation; unreported results show that it delivers a better finite-sample performance than a test based on a standard shrinking bandwidth approximation (e.g. Wu and Zhao, 2007).

For both tests, we show that the dependence of the local performance on unobserved parameters does not affect the asymptotic distribution of the test statistic, as long as the parameters are also estimated locally.

Our research is related to several papers in the literature, in particular Rossi (2005) and, more distantly, to Muller and Petalas (2009), Elliott and Muller (2006), Andrews and Ploberger (1994)

and Andrews (1993). Rossi (2005) proposes a test that is similar to our one-time reversal test but focuses on the case of nested and correctly specified models. Here we consider the more general case of non-nested and misspecified models and propose two additional tests. In a companion paper, Giacomini and Rossi (2010) investigate the problem of testing the time variation in the relative performance of models in an out-of-sample forecasting context. Even though some of the techniques are similar, the additional complication in the in-sample context considered in this paper is that the measure of relative performance depends on estimated parameters, which needs to be taken into account when performing inference. The dependence on parameter estimates can instead be ignored in an out-of-sample context, provided one adopts the asymptotic approximation with finite estimation window considered by Giacomini and Rossi (2010).

Our approach in this paper is also related to the literature on parameter instability testing (e.g., Brown, Durbin and Evans, 1975; Ploberger and Kramer, 1992; Andrews, 1993; Andrews and Ploberger, 1994; Elliott and Muller, 2006; Muller and Petalas, 2009) in that we adapt the tools developed in that literature to our different context where the null hypothesis of interest is a joint hypothesis that the relative performance of the models is equal at each point in time and that the parameters are constant.

The fact that parameters are constant under our null hypothesis means that we are not considering the potentially relevant case in which the performance of two models is equal in spite of their parameters changing over time. The reason for excluding this case is a pragmatic one. In principle, one could have developed versions of our tests that allow for some time variation in parameters under the null hypothesis. Doing so would however be costly in terms of the general applicability of our techniques, as it would require us to impose additional restrictions on the type of time variation under the null hypothesis, the properties of the data and the models that are compatible with the assumptions on which the tests are based. We illustrate this point more concretely when discussing the assumptions of each test in the body of the paper.

One important limitation of our approach is that our methods are not applicable when the competing models are nested, which is common in the literature on model selection testing based on Kullback-Leibler-type of measures. See Rivers and Vuong (2002) for an in-depth discussion of this issue.

The paper is structured as follows. The next section discusses a motivating example that illustrates the procedures proposed in this paper. Section 3 defines the null hypotheses and Section 4 describes the tests. Section 5 evaluates the small sample properties of our proposed procedures in a Monte Carlo experiment, and Section 6 presents the empirical results. Section 7 concludes. The proofs are collected in the appendix.

## 2 Motivating Example

Let  $y_t = \beta_t^0 x_t + \gamma_t^0 z_t + u_t$ , with  $u_t \sim i.i.d. N(0, 1)$ ,  $x_t, z_t$  independent  $N(0, \sigma_{x,t}^2)$  and  $N(0, \sigma_{z,t}^2)$ , respectively, independent of each other and of  $u_t$  for  $t = 1, \dots, T$ , so that the true conditional density of  $y_t$  is  $h_t : N(\beta_t^0 x_t + \gamma_t^0 z_t, 1)$ . Suppose the researcher's goal is to compare two misspecified models: model 1, which specifies a density  $f_t : N(\beta_t x_t, 1)$  and model 2, with density  $g_t : N(\gamma_t z_t, 1)$ . Here  $\beta_t$  and  $\gamma_t$  denote the pseudo-true parameters, defined as the parameters that maximize the expected log-density at time  $t$ ,  $\beta_t = \arg \max_{\beta} E[\ln f_t(\beta)]$  and  $\gamma_t = \arg \max_{\gamma} E[\ln g_t(\gamma)]$ . Even though under the assumptions considered in this example the pseudo-true parameters coincide with the true parameters, in general cases  $\beta_t$  and  $\beta_t^0$  will be different. For instance, introducing correlation between  $x_t$  and  $z_t$  will yield  $\beta_t = \beta_t^0 + \gamma_t^0 E[x_t z_t] / E[x_t^2]$ .

To measure the relative distance of  $f_t$  and  $g_t$  from  $h_t$  at time  $t$  we propose using the Kullback-Leibler Information Criterion at time  $t$ ,  $\Delta KLIC_t$ , (henceforth the “local relative KLIC”), defined as:

$$\text{Local relative KLIC : } \Delta KLIC_t(\theta_t) = E[\ln(h_t/g_t)] - E[\ln(h_t/f_t)] = E[\ln f_t(\beta_t) - \ln g_t(\gamma_t)], \quad (1)$$

where  $\theta_t = (\beta_t', \gamma_t')'$  and the expectation is taken with respect to the true density  $h_t$ . If  $\Delta KLIC_t(\theta_t) > 0$ , model 1 performs better than model 2 at time  $t$ . In our example, it can be shown that  $\beta_t = \beta_t^0$  and  $\gamma_t = \gamma_t^0$  and<sup>1</sup>

$$\Delta KLIC_t(\theta_t) = \frac{1}{2} [\beta_t^2 \sigma_{x,t}^2 - \gamma_t^2 \sigma_{z,t}^2]. \quad (2)$$

Intuitively,  $\Delta KLIC_t(\theta_t)$  measures the relative degree of mis-specification of the two models at time  $t$ . For model 2, the contribution of its mis-specification is reflected in the contribution of the omitted variable  $x_t$  to the variance of the error term, which equals  $\beta_t^2 \sigma_{x,t}^2$ . Similarly, the mis-specification of model 1 is measured by  $\gamma_t^2 \sigma_{z,t}^2$ . Thus, model 2 performs better than model 1 if the contribution of its mis-specification to the variance of the error is smaller than for model 1.

Importantly, equation (2) shows that the time variation in the relative KLIC reflects the time variation in the relative mis-specification of the two models. In particular, the time variation in relative performance might be due to the fact that the parameters change in ways that affect  $\Delta KLIC_t$  differently over time, but it might also be caused by the variances of the regressors  $\sigma_{x,t}^2$  and  $\sigma_{z,t}^2$  changing in different ways over time while the parameters remain constant. Moreover, time-variation in pseudo-true parameters does not correspond exactly to time-variation in true parameters, as can be seen from the expression  $\beta_t = \beta_t^0 + \gamma_t^0 E[x_t z_t] / E[x_t^2]$  obtained in the case of correlated regressors, as  $\beta_t$  could display different patterns of time variation depending on whether and how the different components change over time.

---

<sup>1</sup>We have  $\Delta KLIC_t = \frac{1}{2} E[(u_t + \beta_t x_t)^2 - (u_t + \gamma_t z_t)^2] = \frac{1}{2} E[\beta_t^2 x_t^2 - \gamma_t^2 z_t^2] = \frac{1}{2} (\beta_t^2 \sigma_{x,t}^2 - \gamma_t^2 \sigma_{z,t}^2)$

### 3 Null and Alternative Hypotheses

In this paper we construct tests of equal performance of two models over time that take into account the dependence of the relative performance on estimated parameters, and that allow for different types of time variation in both relative performance and in parameters over time. In the rest of the paper, we no longer make the distinction between true parameters and pseudo-true parameters, as all of our tests will be expressed in terms of pseudo-true parameters. For this reason, we adopt the convention of referring to pseudo-true parameters simply as parameters.

Recall that the measure of local performance is the local relative KLIC, defined as

$$\Delta KLIC_t(\theta_t) = E[\ln f_t(\beta_t) - \ln g_t(\gamma_t)], \quad (3)$$

where

$$\begin{aligned} \theta_t &= (\beta'_t, \gamma'_t)', \\ \beta_t &= \arg \max_b E[\ln f_t(b)], \gamma_t = \arg \max_c E[\ln g_t(c)]. \end{aligned} \quad (4)$$

We assume throughout that  $\theta_t \in \Theta$ , with  $\Theta$  compact.

We propose two different tests, which correspond to different null and alternative hypotheses.

The first test (“one-time reversal test”) assumes that under the null hypothesis the models perform equally well and the parameters are constant, whereas under the alternative hypothesis there is a one-time change in relative performance as well as (at most) a one-time change in parameters at the same time.

This corresponds to the following null and alternative hypotheses:

$$H_0^{OT} : \{\Delta KLIC_t(\theta_t) = 0\} \cap \{\theta_t = \theta\} \text{ for } t = 1, \dots, T \text{ and some } \theta \in \Theta, \quad (5)$$

and

$$\begin{aligned} H_1^{OT} &: \cup_{\delta \in \Pi} \{\Delta KLIC_t(\theta_t) = \mu_1(\delta) 1(t \leq [T\delta]) + \mu_2(\delta) 1(t > [T\delta])\} \\ \cap \{\theta_t &= \theta_1(\delta) 1(t \leq [T\delta]) + \theta_2(\delta) 1(t > [T\delta])\}, t = 1, \dots, T, \end{aligned} \quad (6)$$

for some  $(\mu_1(\delta), \mu_2(\delta)) \neq (0, 0)$ , some  $\delta \in \Pi \subset (0, 1)$ , and/or some  $\theta_1(\delta) \neq \theta_2(\delta)$ ,  $\theta_i(\delta) = (\beta_i(\delta)', \gamma_i(\delta)')'$ ,  $i = 1, 2$ . Thus,  $\mu_i(\delta)$ ,  $i = 1, 2$  are the measures of local performance and  $\beta_i(\delta)$  and  $\gamma_i(\delta)$ ,  $i = 1, 2$  are the local parameters for the sub-samples before and after the reversal, which occurs at the unknown fraction of the total sample  $\delta$ . The one-time reversal test thus focuses on the models’ local relative performance by measuring it separately before and after the reversal. In case the null hypothesis is rejected, the time of the change,  $[T\delta]$ , can be estimated and the path of relative performance equals  $\mu_1(\delta)$  before the change and  $\mu_2(\delta)$  after the change.

The second test (“fluctuation test”) is based on a nonparametric estimator of relative performance. The fluctuation test uses a novel asymptotic approximation which assumes that the bandwidth is fixed. In this approximation, consistent estimation of the local relative performance  $\Delta KLIC_t(\theta_t)$  is not possible, but what can be consistently estimated is a different measure of relative performance, which is a smoothed version of the local relative KLIC, computed at the smoothed local parameter:

$$\text{Smoothed local relative KLIC : } \Delta KLIC_t^*(\theta_t^*) = E \left[ \frac{1}{Th} \sum_{j=1}^T K \left( \frac{t-j}{Th} \right) (\ln f_j(\beta_t^*) - \ln g_j(\gamma_t^*)) \right], \quad (7)$$

where  $\theta_t^* = (\beta_t^*, \gamma_t^*)'$  is defined as

$$\beta_t^* = \arg \max_b E \left[ \frac{1}{Th} \sum_{j=1}^T K \left( \frac{t-j}{Th} \right) \ln f_j(b) \right], \quad (8)$$

(and similarly for  $\gamma_t^*$ ), with  $K(\cdot)$  a kernel function and  $h$  the bandwidth.

The fluctuation test corresponds to different null and alternative hypotheses:

$$H_0^{FB} : \{\Delta KLIC_t^*(\theta_t^*) = 0\} \cap \{\theta_t^* = \theta^*\} \text{ for } t = 1, \dots, T \text{ and some } \theta^* \in \Theta, \quad (9)$$

and

$$H_1^{FB} : \Delta KLIC_t^*(\theta_t^*) \neq 0 \text{ at some } 1 \leq t \leq T. \quad (10)$$

In the example in Section 2, using a rectangular kernel with bandwidth  $h = m/T$  (and assuming for simplicity that  $m$  is an even number) we can see that  $\Delta KLIC_t^*(\theta_t^*)$  differs from  $\Delta KLIC_t(\theta_t)$  in (2), since, first of all,  $\beta_t^* = \left( \frac{1}{m} \sum_{j=t-m/2+1}^{t+m/2} \sigma_{x,j}^2 \beta_j^0 \right) / \left( \frac{1}{m} \sum_{j=t-m/2+1}^{t+m/2} \sigma_{x,j}^2 \right) \neq \beta_t^0$ , whereas  $\beta_t = \beta_t^0$ , and

$$\Delta KLIC_t^*(\theta_t^*) = \frac{1}{2} \left[ \frac{1}{m} \sum_{j=t-m/2+1}^{t+m/2} \left[ 2\beta_t^0 \beta_t^* - (\beta_t^*)^2 \right] \sigma_{x,j}^2 - \frac{1}{m} \sum_{j=t-m/2+1}^{t+m/2} \left[ 2\gamma_t^0 \gamma_t^* - (\gamma_t^*)^2 \right] \sigma_{z,j}^2 \right]. \quad (11)$$

An important point to emphasize is that for both tests the null hypothesis is a joint hypothesis of equal performance and constant parameters. This in practice rules out situations in which two models have equal performance but their parameters are changing over time. The assumption of constant parameters under the null hypothesis is in principle stronger than necessary, but it facilitates the statement and the verification of the assumptions on which the tests rely. For example, allowing for time-varying parameters under the null hypothesis would be difficult to reconcile with the assumption of constant variance of the loss differences that we make for all tests, since in general the variance of the loss differences will depend on the models' parameters. Both assumptions

of constant parameters and constant variance could be relaxed in the context of specific models and/or for a specific test. For example, one could allow for time variation in parameters that disappears asymptotically, or make local stationarity assumptions such as those considered in the nonparametric estimation literature (e.g., Kristensen, 2013). One could also relax the constant variance assumption and rely on bootstrap methods to derive the tests, along the lines of Cavaliere and Taylor (2005).

The difference between the various alternative hypotheses as well as the difference between  $\Delta KLIC_t(\theta_t)$  and  $\Delta KLIC_t^*(\theta_t^*)$  is clarified by Figure 1, which shows an example of two different types of time variation in relative performance that could arise in the context of the simple example considered in this section. In the first scenario (left panels of Figure 1) the time variation in relative performance is due to  $\beta_t$  varying smoothly as a random walk whereas  $\gamma_t, \sigma_{x,t}^2, \sigma_{z,t}^2$  are constant,  $t = 1, \dots, 100$ . In the second scenario (right panels of Figure 1),  $\beta_t, \gamma_t, \sigma_{z,t}^2$  are constant but the relative performance is time-varying because  $\sigma_{x,t}^2$  has a break at  $T/2$ .

INSERT FIGURE 1 HERE

Figures 1(a) and 1(b) report the local relative  $KLIC_t$  in equation (1) in the two scenarios. Figures 1(c) and 1(d) show  $\Delta KLIC_t$  as well as  $\Delta KLIC_t^*$  in equation (7) computed using a bandwidth  $m/T = 1/5$ . Note that Figures 1(a-d) report population quantities (that is, they assume that the parameters and variances are known). Finally, Figures 1(e) and 1(f) show the measure of relative performance that arises as a result of testing (9) and (5). One can see that the measures of relative performance that we propose capture the time variation in the relative performance of the models over time.

In contrast, the large dot reported in panels (a-d) of Figure 1 shows the global relative KLIC ( $T^{-1} \sum_{t=1}^T \Delta KLIC_t$ ), which compares the average performance of the models over the whole sample and which is the object of interest of existing tests in the literature (e.g., Rivers and Vuong (2002)). One can see that the global relative KLIC is very close to zero, which means that Rivers and Vuong's (2002) test would not reject the null hypothesis that the models perform equally well. This occurs because in our example there are reversals in the relative performance of the models during the time period considered. Since model 1 is better than model 2 in the first part of the sample, but model 2 is better than model 1 in the second part of the sample by a similar magnitude, on average over the full sample the two models have similar performance. However, the figure shows that the relative performance did change over time, and that the existing approaches would miss this important feature of the data, whereas our approach would be able to reveal which model performed best at different points in time.

In the following section, we develop the theory for the two statistical tests. The one-time reversal test of hypothesis (5) can be intuitively viewed as performing a Rivers and Vuong's (2002)

test of equal performance allowing for one structural break under the alternative. The fluctuation test of hypothesis (9) relies on constructing confidence bands for the different object  $\Delta KLIC_t^*$  in (7) under the null hypothesis by using a fixed-bandwidth approximation. We refer to this test as the fluctuation test in analogy with the literature on parameter stability testing (Brown et al. 1975 and Ploberger and Kramer 1992). Even though one can see that our tests draw on the existing literature on parameter instability testing, we face additional challenges in particular due to the fact that we are testing joint hypotheses of equal performance and stability and that the measure of performance depends on unknown parameters.

The two tests involve trade-offs, some of which are highlighted by Figure 1. The first consideration is what type of alternative hypothesis seems more appropriate in a given situation. If the type of variation under the alternative hypothesis is a one-time change, the one-time reversal test (Figure 1(f)) will in principle capture it; conversely, the fluctuation test (Figure 1(d)), which relies on the smoothed relative KLIC, will smooth out the time variation: the time variation will thus be more difficult to detect, lowering the power of the test (again, depending on the choice of bandwidth). This is also the case when one postulates a smooth change under the alternative hypothesis, in which case the fluctuation test (Figure 1(c)) should have lower power than the other test because of its smoothing out of the time variation. The one-time reversal test would also be suboptimal in this context because it is based on an approximate measure of time variation, as can be seen in Figure 1(c).

How would the tests that we propose be implemented in practice? We provide an example in Figures 1(e-h). For the fluctuation test we provide boundary lines that would contain the time path of the models' smoothed local relative  $KLIC$  with a pre-specified probability level under the null hypothesis that the relative performance is equal. Figures 1(e,f) depict such boundary lines. Clearly, the test rejects the hypothesis that the relative performance is the same. When this happens, researchers can rely on visual inspection of the local average  $\Delta KLIC$  to ascertain which model performed best at any point in time.

Figures 1(g,h) illustrate the one-time reversal test<sup>2</sup> for the two cases. The procedure estimates the time of the largest change in the relative performance, and then fits measures of average performance separately before and after the reversal. Figure 1(h) shows that when the true underlying relative performance has a sharp reversal, such as in the second scenario, then the procedure will accurately estimate its time path. However, when the true underlying relative performance evolves smoothly over time, then the procedure will approximate it with a sharp reversal, as depicted in Figure 1(g). In both cases, the one-time reversal test strongly rejects the null hypothesis of equal performance.

---

<sup>2</sup>The One-time Reversal test is implemented as a Sup-type test. See Section 4.1 for more details.



## 4 Tests of Stability in the Relative Performance of Models

In this section, we derive the two classes of tests assessing the stability in the relative performance of two models over time.

Each test assumes that the user has available two possibly misspecified parametric models for the variable of interest  $y_t$ . The models can be multivariate, dynamic and nonlinear. In line with the literature (e.g., Vuong (1989) and Rivers and Vuong (2002)), an important restriction is that the models must be non-nested, which, loosely speaking, means that the models' likelihoods cannot be obtained from each other by imposing parameter restrictions.

### 4.1 The One-time Reversal Test

The first test is a test of the null hypothesis (5) against (6), which draws from the literature on testing the stability of the mean of a time series (e.g. Andrews, 1993). These tests are designed for a specific form of time variation in the relative performance of the models under the alternative hypothesis, namely a one-time reversal in the relative performance and in the parameters, which occur at the same time.

The test is implemented as follows. For a given  $\delta \in \Pi \subset (0, 1)$ , let  $\hat{\mu}(\delta) \equiv [\hat{\mu}_1(\delta), \hat{\mu}_2(\delta)]$ , where:

$$\hat{\mu}_1(\delta) = \frac{1}{[T\delta]} \sum_{t=1}^{[T\delta]} \Delta L_t(\hat{\theta}_1(\delta)), \quad \hat{\mu}_2(\delta) = \frac{1}{[T(1-\delta)]} \sum_{t=[T\delta]+1}^T \Delta L_t(\hat{\theta}_2(\delta)) \quad (12)$$

and  $\hat{\theta}_1(\delta) = (\hat{\beta}_1(\delta)', \hat{\gamma}_1(\delta)')'$ ,  $\hat{\theta}_2(\delta) = (\hat{\beta}_2(\delta)', \hat{\gamma}_2(\delta)')'$ ,  $\Delta L_t(\cdot) = \ln f_t(\cdot) - \ln g_t(\cdot)$ , with

$$\begin{aligned} \hat{\beta}_1(\delta) &= \arg \max_b \left( \frac{1}{[T\delta]} \sum_{t=1}^{[T\delta]} \ln f_t(b) \right) \\ \hat{\beta}_2(\delta) &= \arg \max_b \left( \frac{1}{[T(1-\delta)]} \sum_{t=[T\delta]+1}^T \ln f_t(b) \right), \end{aligned}$$

(and similarly for  $\hat{\gamma}_1(\delta), \hat{\gamma}_2(\delta)$ ). Also, let  $\hat{\beta}_T = \arg \max_b \left( \frac{1}{T} \sum_{t=1}^T \ln f_t(b) \right)$  (and similarly for  $\hat{\gamma}_T$ ), and  $\hat{\theta}_T \equiv [\hat{\beta}_T', \hat{\gamma}_T']'$ .

The test relies on the following assumptions:

*Assumptions OT:* Let  $\theta$  be the constant value of  $\theta_t$  under the null hypothesis. The following holds: (1)  $\left\{ T^{-1/2} \sum_{j=1}^{[\lambda T]} \Delta L_j(\theta) \right\}$  obeys a Functional Central Limit Theorem (FCLT) under  $H_{0,T}^{OT}$  for  $0 \leq \lambda \leq 1$ , such that:  $\sigma^{-1} T^{-1/2} \sum_{j=1}^{[\lambda T]} \Delta L_j(\theta) \Rightarrow \mathcal{B}(\lambda)$  (a standard Brownian Motion process); (2)  $\sup_{\delta \in \Pi} \|\hat{\beta}_1(\delta) - \beta\| = o_p(1)$  and  $\sup_{\delta \in \Pi} \|T^{1/2} (\hat{\beta}_1(\delta) - \beta)\| = O_p(1)$  under  $H_{0,T}^{OT}$  as  $T \rightarrow \infty$  (and similarly for  $\hat{\gamma}_1(\delta), \hat{\beta}_2(\delta), \hat{\gamma}_2(\delta), \hat{\beta}_T, \hat{\gamma}_T$ ); (3) (a) the log-likelihoods of both models,  $\ln f_T(\beta, \delta) = \sum_{t=1}^T \ln f_t(\beta, \delta)$  and  $\ln g_T(\gamma, \delta) = \sum_{t=1}^T \ln g_t(\gamma, \delta)$ , do not depend on  $\delta$  for all  $\beta, \gamma$  in the

null hypothesis; (b)  $\theta$  is an interior point of the parameter space  $\Theta$ ; (c)  $f_T(\tilde{\beta}, \delta), g_T(\tilde{\gamma}, \delta)$  are twice continuously partially differentiable in  $\tilde{\beta}, \tilde{\gamma}$  for all  $\delta \in \Pi$  and  $\tilde{\beta}, \tilde{\gamma}$  in some neighborhood of the null,  $\Theta_0$ ; (d)  $-B_T^{-1} \nabla^2 \ln f_T(\tilde{\beta}, \delta) B_T^{-1} \xrightarrow{p} \Phi_\beta(\tilde{\beta}, \delta), -B_T^{-1} \nabla^2 \ln g_T(\tilde{\gamma}, \delta) B_T^{-1} \xrightarrow{p} \Phi_\gamma(\tilde{\gamma}, \delta)$  uniformly over  $\delta \in \Pi$  and  $\tilde{\beta}, \tilde{\gamma} \in \Theta_0$  under  $\beta, \gamma$  for some nonrandom matrix functions  $\Phi_\beta(\beta, \delta), \Phi_\gamma(\gamma, \delta)$  and some sequence of nonrandom diagonal matrices  $\{B_T : T \geq 1\}$  whose elements diverge to infinity as  $T \rightarrow \infty$ ;<sup>3</sup> (e)  $\Phi_\beta(\tilde{\beta}, \delta), \Phi_\gamma(\tilde{\gamma}, \delta)$  are uniformly continuous in  $(\tilde{\theta}, \delta)$  over  $\Theta_0 \times \Pi$ ; (f)  $\Phi_\beta(\beta, \delta), \Phi_\gamma(\gamma, \delta)$  are uniformly positive definite over  $\delta \in \Pi$ ; (4)  $\lim_{T \rightarrow \infty} \text{var} \left( T^{-1/2} \sum_{t=1}^T \Delta L_t(\theta) \right) = \sigma^2 > 0$  is constant and finite, and  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$ .

Assumption OT(1) assumes a FCLT for partial sum processes. Assumptions OT(2,3) are standard assumptions that guarantee that the estimated parameters as well as the score functions obey regularity conditions ensuring their convergence. In particular, Assumption OT(3) follow from assumptions similar to Andrews and Ploberger (1994). Assumption OT(4) imposes that the variance is constant under the null hypothesis, and a consistent estimator under the null hypothesis is for example a standard HAC estimator

$$\hat{\sigma}^2 = \sum_{i=-q(T)+1}^{q(T)-1} (1 - |i/q(T)|) T^{-1} \sum_{j=q(T)}^{T+1-q(T)} \Delta L_j^d(\hat{\theta}_T) \Delta L_{j-i}^d(\hat{\theta}_T), \quad (13)$$

where  $q(T)$  is a bandwidth that grows with  $T$  (e.g., Newey and West, 1987, Andrews, 1991), and  $\Delta L_j^d(\cdot)$  indicates demeaned  $\Delta L_j(\cdot)$ .

Under Assumption OT, we provide Sup-type tests for the one-time reversal in the following proposition:<sup>4</sup>

**Proposition 1 (Sup-type Test)** *Suppose Assumption OT holds. Let  $QLR_T^* = \sup_{\delta \in \Pi} \Phi_T(\delta)$ ,  $\Phi_T(\delta) = LM_1 + LM_2(\delta)$ , where  $\Pi \subset (0, 1)$  and*

$$\begin{aligned} LM_1 &= \hat{\sigma}^{-2} \left[ T^{-1/2} \sum_{t=1}^T \Delta L_t(\hat{\theta}_T) \right]^2, \\ LM_2(\delta) &= \hat{\sigma}^{-2} \frac{1}{\delta(1-\delta)} \left[ (1-\delta) T^{-1/2} \sum_{t=1}^{[T\delta]} \Delta L_t(\hat{\theta}_1(\delta)) - \delta T^{-1/2} \sum_{t=[T\delta]+1}^T \Delta L_t(\hat{\theta}_2(\delta)) \right]^2. \end{aligned}$$

Under the null hypothesis  $H_0^{OT}$ , we have:  $QLR_T^* \Rightarrow \sup_{\delta \in \Pi} \left[ \frac{\mathcal{B}\mathcal{B}(\delta)^2}{\delta(1-\delta)} + \mathcal{B}(1)^2 \right]$ , and  $\mathcal{B}(\cdot)$  and  $\mathcal{B}\mathcal{B}(\delta) \equiv \mathcal{B}(\delta) - \delta \mathcal{B}(1)$  are, respectively, a standard univariate Brownian motion and a Brownian bridge. The null hypothesis is thus rejected when  $QLR_T^* > k_\alpha$ . The critical values  $(\alpha, k_\alpha)$  are: (0.05, 9.8257), (0.10, 8.1379).

<sup>3</sup>  $\nabla^2$  denotes the second derivative with respect to the parameter.

<sup>4</sup> Sup-type tests have been used in the parameter instability literature since Andrews (1993). Note that the sup-type test could alternatively be implemented as  $\sup_{\delta \in \Pi} W_T(\delta)$ , where  $W_T(\delta)$  is defined in eq. (14).

Among the advantages of the Sup-type approach, we have that: (i) when the null hypothesis is rejected, it is possible to evaluate whether the rejection is due to instabilities in the relative performance or to a model being constantly better than its competitor; (ii) if such instability is found, it is possible to estimate the time of the switch in the relative performance; (iii) the test is designed against one time breaks in the relative performance. Here below is a step by step procedure to implement the approach suggested in Proposition 1 with an overall significance level  $\alpha$ :

(i) test the hypothesis of equal performance at each time by using the statistic  $QLR_T^*$  from Proposition 1 at  $\alpha$  significance level;

(ii) if the null is rejected, compare  $LM_1$  and  $\sup_{\delta \in \Pi} LM_2(\delta)$ , with the following critical values: (3.84, 8.85) for  $\alpha = 0.05$ , (2.71, 7.17) for  $\alpha = 0.10$ , and (6.63, 12.35) for  $\alpha = 0.01$ . If only  $LM_1$  rejects then there is evidence in favor of the hypothesis that one model is constantly better than its competitor. If only  $\sup_{\delta \in \Pi} LM_2(\delta)$  rejects, then there is evidence that there are instabilities in the relative performance of the two models but neither is constantly better over the full sample. Note that the latter corresponds to Andrews' (1993) Sup-test for structural break. If both reject then it is not possible to attribute the rejection to a unique source.<sup>5</sup>

(iii) estimate the time of the reversal by  $\delta^* = \arg \sup_{\delta \in \{0.15, \dots, 0.85\}} LM_2(\delta)$  and let  $t^* \equiv [\delta^*/T]$ .

(iv) to extract information on which model to choose, we suggest to plot the time path of the underlying relative performance as:

$$\begin{cases} \frac{1}{t^*} \sum_{t=1}^{t^*} \left( \ln f_t(\hat{\beta}_1(\delta^*)) - \ln g_t(\hat{\gamma}_1(\delta^*)) \right) & \text{for } t \leq t^*; \\ \frac{1}{(T-t^*)} \sum_{t=t^*+1}^T \left( \ln f_t(\hat{\beta}_2(\delta^*)) - \ln g_t(\hat{\gamma}_2(\delta^*)) \right) & \text{for } t > t^*. \end{cases}$$

We also provide tests similar in spirit to those proposed by Andrews and Ploberger (1994). The tests rely on the Wald-type test statistic, rather than LM-type.<sup>6</sup>

**Corollary 2 (AP test)** *Suppose Assumption OT holds. Consider the test statistics*

$$W_T(\delta) = T\hat{\mu}(\delta)' H' \left( H\mathcal{I}_{T,\delta}^{-1} H' \right)^{-1} H\hat{\mu}(\delta) \quad (14)$$

$$ExpW_{\infty,T}^* = \ln \frac{1}{1-2\delta_0} \int_{\delta_0}^{1-\delta_0} \exp \left( \frac{1}{2} W_T(\delta) \right) d\delta; \quad (15)$$

$$MeanW_T^* = \frac{1}{1-2\delta_0} \int_{\delta_0}^{1-\delta_0} W_T(\delta) d\delta, \quad (16)$$

---

<sup>5</sup>This procedure is justified by the fact that the two components  $LM_1$  and  $LM_2$  are asymptotically independent – see Rossi (2005). Performing two separate tests does not result in an optimal test, but it is nevertheless useful to heuristically disentangle the causes of rejection of equal performance. The critical values for  $LM_1$  are from a  $\chi_1^2$  distribution whereas those for  $LM_2$  are from Andrews (1993).

<sup>6</sup>The Wald-type test allows for a more general variance estimator. One could also implement the Sup-type test in Proposition 1 as  $QLR_T^* = \sup_{\delta \in \Pi} W_T(\delta)$ .

where  $\delta_0 = 0.15$ ,  $H \equiv \begin{pmatrix} 1 & -1 \\ \delta & 1 - \delta \end{pmatrix}$ ,  $\mathcal{I}_{T,\delta}^{-1} = \begin{pmatrix} \delta^{-1}\hat{\sigma}_1^2 & 0 \\ 0 & (1 - \delta)^{-1}\hat{\sigma}_2^2 \end{pmatrix}$ ,  $\hat{\sigma}_1^2$  is a HAC estimator of the asymptotic variance of  $\Delta L_t(\hat{\theta}_1(\delta))$ ,  $t = 1, \dots, [T\delta]$  and  $\hat{\sigma}_2^2$  is a HAC estimator of the variance of  $\Delta L_t(\hat{\theta}_2(\delta))$ ,  $t = [T\delta] + 1, \dots, T$ .<sup>7</sup> Under the null hypothesis  $H_0^{OT}$ ,

$$ExpW_{\infty,T}^* \Rightarrow \ln \left[ \frac{1}{1 - 2\delta_0} \int_{\delta_0}^{1-\delta_0} \exp \left( \frac{1}{2} \frac{\mathcal{BB}(\delta)^2}{\delta(1-\delta)} + \frac{1}{2} \mathcal{B}(1)^2 \right) d\delta \right], \quad (17)$$

$$MeanW_T^* \Rightarrow \frac{1}{1 - 2\delta_0} \int_{\delta_0}^{1-\delta_0} \left( \frac{\mathcal{BB}(\delta)^2}{\delta(1-\delta)} + \mathcal{B}(1)^2 \right) d\delta, \quad (18)$$

where  $t = [\delta T]$  and  $\mathcal{B}(\cdot)$  and  $\mathcal{BB}(\cdot)$  are, respectively, a standard univariate Brownian motion and a Brownian bridge, where  $\mathcal{BB}(\delta) \equiv \mathcal{B}(\delta) - \delta \mathcal{B}(1)$ . The null hypothesis is rejected when  $ExpW_{\infty,T}^* > \kappa_\alpha$  and  $MeanW_T^* > \nu_\alpha$ . Simulated values of  $(\alpha; \kappa_\alpha, \nu_\alpha)$  are:  $(0.05; 3.13, 5.36)$  and  $(0.10; 2.44, 4.26)$ .

The power properties of these tests will be evaluated in Section 5.

## 4.2 The Fluctuation Test

In this section, we derive a test of the hypothesis (9) against the alternative (10). The test relies on constructing a nonparametric estimate of the local relative performance; however we consider an asymptotic approximation in which the bandwidth is fixed instead of the common approximation used in the nonparametric literature which requires that the bandwidth shrinks as the sample grows. For simplicity, we restrict attention to a rectangular kernel, but the analysis could be easily extended to the case of a general kernel. For a particular choice of fixed bandwidth  $m = [hT]$ , we thus define the smoothed local relative KLIC using a rectangular kernel as:

$$\Delta KLIC_t^*(\theta_t^*) = m^{-1} \sum_{j=t-m/2+1}^{t+m/2} E[\Delta L_j(\theta_t^*)], \quad t = m/2, \dots, T - m/2. \quad (19)$$

Therefore, the null and alternative hypotheses become:

$$H_0^{FB} : \{ \Delta KLIC_t^*(\theta_t^*) = 0 \} \cap \{ \theta_t^* = \theta^* \} \text{ for } t = m/2, \dots, T - m/2 \text{ and some } \theta^* \in \Theta \quad (20)$$

$$\text{against} \quad (21)$$

$$H_1^{FB} : \Delta KLIC_t^*(\theta_t^*) \neq 0 \text{ at some } m/2 \leq t \leq T - m/2,$$

We estimate the smoothed local relative KLIC in eq. (19) as:

$$m^{-1} \sum_{j=t-m/2+1}^{t+m/2} \Delta L_t(\hat{\theta}_t^*), \quad (22)$$

---

<sup>7</sup>The formula is similar to eq. (13) except that it applies to the relevant sub-sample.

where  $\hat{\theta}_t^* = [\hat{\beta}_t^{*'}, \hat{\gamma}_t^{*'}]'$  are nonparametric estimates of the local parameters obtained as the solution to (e.g., for the first model)

$$m^{-1} \sum_{j=t-m/2+1}^{t+m/2} \nabla \ln f_t(\hat{\beta}_t^*) = 0, \quad (23)$$

with  $\nabla \ln f_t(\cdot)$  denoting the first derivative of the log-likelihood at time  $t$ . The fact that under the null hypothesis the parameters are constant means that one could in principle obtain a valid test by letting the measure of local performance (22) depend on any other estimator of the parameters, and not necessarily a local estimator. The reason for considering a local estimator of the parameters is to obtain a test that has power against smooth time variation in parameters under the alternative hypothesis (although we do not formally derive the properties of the test under the alternative hypothesis, but only those under the null hypothesis).

Our proposed test, which we call the fluctuation test, can be derived under the following assumptions:<sup>8</sup>

*Assumptions FB:* Let  $\theta^*$  be the constant value of  $\theta_t^*$  under the null hypothesis and  $t = [\lambda T]$ . The following holds: (1)  $\left\{ T^{-1/2} \sum_{j=1}^{[\lambda T]} \Delta L_j(\theta^*) \right\}$  obeys a Functional Central Limit Theorem (FCLT) under  $H_{0,T}^{FB}$  for  $\lambda \in \Pi \subset (0, 1)$ , such that:  $\sigma^{-1} T^{-1/2} \sum_{j=1}^{[\lambda T]} \Delta L_j(\theta^*) \Rightarrow \mathcal{B}(\lambda)$  (a standard Brownian Motion process); (2)  $\sup_{\lambda} \|\hat{\beta}_{[T\lambda]}^* - \beta^*\| = o_p(1)$  and  $\sup_{\lambda} \|T^{1/2} (\hat{\beta}_{[T\lambda]}^* - \beta^*)\| = O_p(1)$  under  $H_{0,T}^{FB}$  as  $T \rightarrow \infty$  (and similarly  $\hat{\gamma}_{[T\lambda]}^*(\delta)$ ); (3)  $\sup_{\lambda} \|m^{-1} \sum_{j=[T\lambda]-[T\mu]/2+1}^{[T\lambda]+[T\mu]/2} \nabla^2 \ln f_j(\hat{\beta}_{[T\lambda]}^*) - E(\nabla^2 \ln f_j(\beta^*))\| \xrightarrow{p} 0$  whenever  $\ddot{\beta}_{[T\lambda]}$  satisfies  $\sup_{\lambda} \|\ddot{\beta}_{[T\lambda]} - \beta^*\| \xrightarrow{p} 0$  as  $m, T \rightarrow \infty$ , and  $E(\nabla^2 \ln f_j(\beta^*))$  is positive and finite; (4)  $\lim_{T \rightarrow \infty} \text{var} \left( T^{-1/2} \sum_{t=1}^T \Delta L_t(\theta^*) \right) = \sigma^2 > 0$  is constant and finite, and  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$ ; (5)  $m/T = h$ , with  $h \in (0, \infty)$  and  $m, T \rightarrow \infty$ .

The assumptions underlying the fluctuation test are similar to those considered for the one-time reversal test. In particular we require the loss differences to satisfy a FCLT when evaluated at the pseudo-true parameters, which are assumed to be constant under the null hypothesis. We also assume the asymptotic variance to be constant under the null hypothesis, which again is a stronger requirement than necessary, but it facilitates the statement of the FCLT. Assumptions FB(2,3) are high-level but standard; more primitive conditions can be specified in the context of specific models. The main difference between assumption FB and assumptions used in the non-parametric literature is that the fluctuation test considers a bandwidth that is a fixed proportion of the total sample size (assumption (5)).

---

<sup>8</sup>See Brown et al. (1975) and Ploberger and Kramer (1992) for fluctuation tests in the context of parameter instability.

One can verify that Assumption FB(1) is satisfied in the example of Section 2, where  $\Delta L_t(\theta^*) = \frac{1}{2} \left\{ \left[ 2\beta_t^0 \beta^* - (\beta^*)^2 \right] x_t^2 - \left[ 2\gamma_t^0 \gamma^* - (\gamma^*)^2 \right] z_t^2 \right\}$ . Since  $x_t$  and  $z_t$  are i.i.d., under the further assumption that the true parameters  $\beta_t^0$  and  $\gamma_t^0$  are also constant under the null hypothesis,  $\Delta L_t(\theta^*)$  is i.i.d. and  $\sigma_t^2$  is constant, which satisfies the assumptions of Donsker's FCLT theorem. It is easy to verify that Assumptions FB (2,3) hold in the linear model case that we consider.

The following proposition provides a justification for the fluctuation test.

**Theorem 3 (Fluctuation Test)** *Suppose Assumption FB holds. Consider the test statistic*

$$\max_{t=m/2, \dots, T-m/2} |F_t| = \max_{t=m/2, \dots, T-m/2} \left| \hat{\sigma}^{-1} m^{-1/2} \sum_{j=t-m/2+1}^{t+m/2} \Delta L_j(\hat{\theta}_t^*) \right|, \quad (24)$$

where  $\hat{\theta}_t^*$  is as in discussed in equation (23) and  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$ , e.g. as in equation (13). Under the null hypothesis (20)

$$F_t \implies [\mathcal{B}(\lambda + h/2) - \mathcal{B}(\lambda - h/2)] / \sqrt{h}, \quad (25)$$

where  $t = \lfloor \lambda T \rfloor$ ,  $h = m/T$  and  $\mathcal{B}(\cdot)$  is a standard univariate Brownian motion. The critical values for a significance level  $\alpha$  are  $\pm k_\alpha$ , where  $k_\alpha$  solves

$$\Pr \left\{ \max_{\lambda} \left| [\mathcal{B}(\lambda + h/2) - \mathcal{B}(\lambda - h/2)] / \sqrt{h} \right| > k_\alpha \right\} = \alpha. \quad (26)$$

The null hypothesis is rejected when  $\max_t |F_t| > k_\alpha$ . Simulated values of  $(\alpha, k_\alpha)$  are reported in Table 1 for various choices of  $h = m/T$ .

INSERT TABLE 1 HERE

The fluctuation test relies on a specific choice of bandwidth  $h$ , and the results of the test will be different for different choices of  $h$ . As a practical recommendation, we suggest assessing the sensitivity of the test to a few different choices of  $h$ , which is easy to do as we only provide critical values for several possible choices of  $h$ .

## 5 A Small Monte Carlo Analysis

This section investigates the finite-sample size and power properties of the tests for equal performance introduced in the previous section. We consider two designs for the Data Generating Processes (DGPs), which are representative of the features discussed in the main example in Section 2. In particular, as mentioned before, the time variation in the relative KLIC might be due to the fact that the parameters change in ways that affect  $\Delta KLIC_t$  differently over time; design 1

focuses on this situation. However, time variation in the relative KLIC might also occur when the parameters are constant but some other aspects of the distribution of the data change in different ways over time, which will be described by design 2.

More in details, the true DGP is:

$$y_t = \beta_t x_t + \gamma_t z_t + \varepsilon_t, \quad \varepsilon_t \sim i.i.d.N(0, 1),$$

where  $x_t \sim N(0, \sigma_{x,t}^2)$ ,  $z_t \sim N(0, \sigma_{z,t}^2)$ ,  $t = 1, 2, \dots, T$ ,  $T = 200$ . The two competing models are: Model 1:  $y_t = \beta_t x_t + \varepsilon_{1,t}$  and Model 2:  $y_t = \gamma_t z_t + \varepsilon_{2,t}$ . We consider the following designs:

*Design 1.*  $\sigma_{x,t}^2 = \sigma_{z,t}^2 = 1$ ,  $\gamma_t = 1$ ,  $\beta_t = 1 + \beta_A \cdot 1(t \leq 0.5T) - \beta_A \cdot 1(t > 0.5T)$ . In this design, we let the parameter  $\beta$  change over time, and this affects the relative performance of the models over time.

*Design 2.*  $\sigma_{x,t}^2 = 1 + \sigma_A^2 \cdot 1(t > 0.75T)$ ,  $\sigma_{z,t}^2 = 1$ ,  $\beta_t = 1$ ,  $\gamma_t = 1$ . In this design, the parameters in the conditional mean are constant but one of the variances ( $\sigma_{x,t}^2$ ) changes over time, thus resulting in a change in the relative performance over time.

Tables 2 and 3 show the empirical rejection frequencies of the various tests for a nominal size of 5%. The Fluctuation test is implemented with  $h = 0.3$ . Size properties are obtained by setting  $\beta_A = 0$  and  $\sigma_A = 0$ . Table 2 demonstrates that all tests have good size properties. It also shows that the tests with highest power against a one-time reversal are the  $ExpW_{\infty,T}^*$  and  $QLR_T^*$  tests; the  $MeanW_T^*$  test has slightly lower power than the former. The fluctuation test has worse power properties relative to them. Note that a standard full-sample likelihood ratio test would have power equal to size in design 1. Regarding design 2, Table 3 shows that the  $ExpW_{\infty,T}^*$  and  $QLR_T^*$  tests have quite similar performance in terms of power, although the Sup-type test has slightly better power properties than the other tests, and the fluctuation test has slightly worse power properties.

INSERT TABLES 2 AND 3 HERE

## 6 Empirical Application: Time-variation in the Performance of DSGE vs. BVAR Models

In a highly influential paper, Smets and Wouters (2003) (henceforth SW) show that a DSGE model of the European economy - estimated using Bayesian techniques over the period 1970:2-1999:4 - fits the data as well as a-theoretical Bayesian VARs (BVARs). Furthermore, they find that the parameter estimates from the DSGE model have the expected sign. Perhaps for these reasons, this new generation of DSGE models has attracted a lot of interest from forecasters and central banks. SW's model features include sticky prices and wages, habit formation, adjustment costs in capital accumulation and variable capacity utilization, and the model is estimated using

seven variables: GDP, consumption, investment, prices, real wages, employment, and the nominal interest rate. Their conclusion that the DSGE fits the data as well as BVARs is based on the fact that the marginal data densities for the two models are of comparable magnitudes over the full sample. However, given the changes that have characterized the European economy over the sample analyzed by SW - for example, the creation of the European Union in 1993, changes in productivity and in the labor market, to name a few - it is plausible that the relative performance of theoretical and a-theoretical models may itself have varied over time. In this section, we apply the techniques proposed in this paper to assess whether the relative performance of the DSGE model and of BVARs was stable over time. We extend the sample considered by SW to include data up to 2004:4, for a total sample of size  $T = 145$ .

In order to compute the local measure of relative performance, (the local  $\Delta KLIC$ ), we estimate both models recursively over a moving window of size  $m = 70$  using Bayesian methods. As in SW, the first 40 data points in each sample are used to initialize the estimates of the DSGE model and as training samples for the BVAR priors. We consider a BVAR(1) and a BVAR(2), both of which use a variant of the Minnesota prior, as suggested by Sims (2003).<sup>9</sup> We present results for two different transformations of the data. The first applies the same de-trending of the data used by SW, which is based on a linear trend fitted on the whole sample (we refer to this as “full-sample de-trending”). As cautioned by Sims (2003), this type of pre-processing of the data may unduly favour the DSGE, and thus we further consider a second transformation of the data, where de-trending is performed on each rolling estimation window (“rolling-sample de-trending”).

Figure 2 displays the evolution of the posterior mode of some representative parameters. Figure 2(a) shows parameters that describe the evolution of the persistence of some representative shocks (productivity, investment, government spending, and labor supply); Figure 2(b) shows the estimates of the standard deviation of the same shocks; and Figure 2(c) plots monetary policy parameters. Overall, Figure 2 reveals evidence of parameter variation. In particular, the figures show some decrease in the persistence of the productivity shock, whereas both the persistence and the standard deviation of the investment shock seem to increase over time. The monetary policy parameters appear to be overall stable over time.

FIGURE 2 HERE

We then apply our fluctuation test to test the hypothesis that the DSGE model and the BVAR have equal performance at every point in time over the historical sample.

Figure 3 shows the implementation of the fluctuation test for the DSGE vs. a BVAR(1) and BVAR(2), using full-sample de-trending of the data. The estimate of the local relative  $KLIC$  is

---

<sup>9</sup>The BVAR’s were estimated using software provided by Chris Sims at [www.princeton.edu/~sims](http://www.princeton.edu/~sims). As in Sims (2003), for the Minnesota prior we set the decay parameter to 1 and the overall tightness to .3. We also included sum-of-coefficients (with weight  $\mu = 1$ ) and co-persistence (with weight  $\lambda = 5$ ) prior components.



evaluated at the posterior modes of the models' parameters computed over the rolling windows, using the fact that they are consistent estimates of the pseudo-true parameters  $\beta_t$  and  $\gamma_t$  (see, e.g., Fernandez-Villaverde and Rubio-Ramirez, 2004).

FIGURE 3 HERE

Figure 3 suggests that the DSGE has comparable performance to both a BVAR(1) and BVAR(2) up until the early 1990s, at which point the performance of the DSGE dramatically improves relative to that of the reduced-form models.

To assess whether this result is sensitive to the data filtering, we implement the fluctuation test for the DSGE vs. a BVAR(1) and BVAR(2), this time using rolling-window de-trended data.

FIGURE 4 HERE

The results confirm the suspicion expressed by Sims (2003) that the pre-processing of the data utilized by SW penalizes the reduced-form models in favour of the DSGE. As we see from Figure 4, once the de-trending is performed on each rolling window, the advantage of the DSGE at the end of the sample disappears, and the DSGE performs as well as a BVAR(1) on most of the sample, whereas it is outperformed by a BVAR(2) for all but the last few dates in the sample (when the two models perform equally well).

## 7 Conclusions

This paper developed statistical testing procedures for evaluating the relative performance of two competing models in unstable environments. We proposed two tests: 1) a one-time reversal test; and 2) a fluctuation test. We investigated the advantages and limitations of the two approaches and compared the quality of the approximation that they deliver in finite samples. Based on the results of the latter, the choice between the one-time reversal and the fluctuation test should be driven by the type of alternative hypothesis of interest in a given application. Finally, an empirical application to the European economy points to the presence of instabilities in the models' parameters, and suggests that a VAR fitted the last two decades of data better than a standard DSGE model, a conclusion that is however sensitive to the de-trending method utilized.

## References

- [1] Andrews, D.W.K. (1991), "Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimation", *Econometrica* 59, 817-858.

- [2] Andrews, D.W.K. (1993), "Tests for Parameter Instability and Structural Change with Unknown Change Point", *Econometrica* 61, 821–856.
- [3] Andrews, D.W.K. and W. Ploberger (1994), "Optimal Tests When a Nuisance Parameter is Present only under the Alternative", *Econometrica* 62(6), 1383-1414.
- [4] Brown, R.L., J. Durbin and J.M. Evans (1975), "Techniques for Testing the Constancy of Regression Relationships over Time with Comments", *Journal of the Royal Statistical Society, Series B*, 37, 149-192.
- [5] Cavaliere, G. and R. Taylor (2005), "Stationarity Tests Under Time-Varying Second Moments", *Econometric Theory* 21, 1112-1129.
- [6] Elliott, G. and U. Muller (2006), "Efficient Tests for General Persistent Time Variation in Regression Coefficients", *The Review of Economic Studies* 73, 907-940.
- [7] Fernández-Villaverde, J. and J. Rubio-Ramírez (2004), "Comparing Dynamic Equilibrium Models to Data: A Bayesian Approach.", *Journal of Econometrics* 123.
- [8] Giacomini, R. and B. Rossi (2010), "Forecast Comparisons in Unstable Environments", *Journal of Applied Econometrics* 25(4), 595-620.
- [9] Kristensen, D. (2013), "Nonparametric Estimation of Time-Varying Parameters", mimeo
- [10] Muller, U. and P. Petalas (2009), "Efficient Estimation of the Parameter Path in Unstable Time Series Models", *The Review of Economic Studies*, forthcoming.
- [11] Newey, W. and K. West (1987), "A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix", *Econometrica* 55, 703-708.
- [12] Ploberger, W. and W. Kramer (1992), "The Cusum Test with Ols Residuals", *Econometrica* 60(2), 271-285.
- [13] Rivers, D. and Q. Vuong (2002), "Model Selection Tests for Nonlinear Dynamic Models", *Econometrics Journal*, 5, 1-39.
- [14] Rossi, B. (2005), "Optimal Tests for Nested Model Selection with Underlying Parameter Instabilities", *Econometric Theory* 21(5), 962-990.
- [15] Sims, C. (2003), "Comment on Smets and Wouters", *mimeo*, available at: <https://sims.princeton.edu/yftp/Ottawa/SWcommentSlides.pdf>
- [16] Smets, F. and R. Wouters (2003), "An Estimated Stochastic Dynamic General Equilibrium Model of the Euro Area", *Journal of the European Economic Association*, 1, 1123-1175.

- [17] Vuong, Q. H. (1989), “Likelihood Ratio Tests for Model Selection and Non-nested Hypotheses”, *Econometrica*, 57, 307-333.
- [18] Wu, W.B. and Z. Zhao (2007), “Inference of Trends in Time Series”, *Journal of the Royal Statistical Society* B69, 391-410.

## 8 Appendix A - Proofs

**Lemma 4** Let  $X_T(\delta) = o_{p,\delta}(1)$  denote  $\sup_{\delta \in \Pi} \|X_T(\delta)\| = o_p(1)$  and  $X_T(\delta) = O_{p,\delta}(1)$  denote  $\sup_{\delta \in \Pi} \|X_T(\delta)\| = O_p(1)$ . Under Assumption OT and  $H_0^{OT}$ ,

$$T^{-1/2} \sum_{t=1}^{[T\delta]} \left( \ln f_t(\hat{\beta}_1(\delta)) - \ln g_t(\hat{\gamma}_1(\delta)) \right) = T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta(\delta)) - \ln g_t(\gamma(\delta))) + o_{p,\delta}(1).$$

**Proof of Lemma (4).** By applying a second order Taylor expansion around  $\theta = [\beta', \gamma']'$ :

$$\begin{aligned} & T^{-1/2} \sum_{t=1}^{[T\delta]} \left( \ln f_t(\hat{\beta}_1(\delta)) - \ln g_t(\hat{\gamma}_1(\delta)) \right) \\ = & T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta) - \ln g_t(\gamma)) + \\ & + \frac{1}{2} T^{-1} \sum_{t=1}^{[T\delta]} \nabla \ln f_t(\hat{\beta}_1(\delta)) (\hat{\beta}_1(\delta) - \beta) T^{1/2} \end{aligned} \quad (27)$$

$$- \frac{1}{2} T^{-1} \sum_{t=1}^{[T\delta]} \nabla \ln g_t(\hat{\gamma}_1(\delta)) (\hat{\gamma}_1(\delta) - \gamma) T^{1/2} \quad (28)$$

$$+ \frac{1}{2} (\hat{\beta}_1(\delta) - \beta)' \left[ T^{-1} \sum_{t=1}^{[T\delta]} \nabla^2 \ln f_t(\ddot{\beta}_{1,T}(\delta)) \right] (\hat{\beta}_1(\delta) - \beta) T^{1/2} \quad (29)$$

$$- \frac{1}{2} (\hat{\gamma}_1(\delta) - \gamma)' \left[ T^{-1} \sum_{t=1}^{[T\delta]} \nabla^2 \ln g_t(\ddot{\gamma}_{1,T}(\delta)) \right] (\hat{\gamma}_1(\delta) - \gamma) T^{1/2} \quad (30)$$

$$= T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta) - \ln g_t(\gamma)) + o_{p,\delta}(1) + o_{p,\delta}(1), \quad (31)$$

where  $\ddot{\beta}_{1,T}(\delta)$  is an intermediate point between  $\hat{\beta}_1(\delta)$  and  $\beta$  (similarly for  $\ddot{\gamma}_{1,T}(\delta)$ ). By definition of the ML estimator,  $\sum_{t=1}^{[T\delta]} \nabla \ln f_t(\hat{\beta}_1(\delta)) = 0$  (similarly for  $\nabla \ln g_t$ ). By Assumption OT(2),  $T^{1/2}(\hat{\beta}_1(\delta) - \beta) = O_{p,\delta}(1)$  which proves that (27) and (28) are  $o_{p,\delta}(1)$ . Furthermore, Assumptions OT(2,3) ensure that, under the null hypothesis,  $T^{-1} \sum_{t=1}^{[T\delta]} \nabla^2 \ln f_t(\ddot{\beta}_{1,T}(\delta)) = O_{p,\delta}(1)$  (and similarly for the component in  $\gamma$ ); by Assumption OT(2),  $T^{1/2}(\hat{\beta}_1(\delta) - \beta) = O_{p,\delta}(1)$  and  $(\hat{\beta}_1(\delta) - \beta) = o_{p,\delta}(1)$  (and similarly for the components in  $\gamma$ ), proving that (29) and (30) are  $o_{p,\delta}(1)$ . ■

**Proof of Proposition 1.** By Lemma 4 and similar arguments, as well as the consistency of  $\hat{\sigma}$ , we have:

$$(i) \quad LM_1 = \sigma^{-2} \left[ T^{-1/2} \sum_{t=1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right]^2 + o_{p,\delta}(1)$$

and

$$\begin{aligned}
(ii) \quad LM_2(\delta) &= \sigma^{-2} \delta^{-1} (1 - \delta)^{-1} \\
&\quad \left[ (1 - \delta) T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) + \right. \\
&\quad \left. - \delta T^{-1/2} \sum_{t=[T\delta]+1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right]^2 + o_{p,\delta}(1).
\end{aligned}$$

By Assumption OT, under the null hypothesis:

$$\sigma^{-1} T^{-1/2} \sum_{t=1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \implies \mathcal{B}(1) \quad (32)$$

$$\begin{aligned}
&\sigma^{-1} \delta^{-1/2} (1 - \delta)^{-1/2} \left[ T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right. \\
&\quad \left. - \delta T^{-1/2} \sum_{t=1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right] \\
&\implies \delta^{-1/2} (1 - \delta)^{-1/2} [\mathcal{B}(\delta) - \delta \mathcal{B}(1)] = \delta^{-1/2} (1 - \delta)^{-1/2} \mathcal{B}\mathcal{B}(\delta), \quad (33)
\end{aligned}$$

where the limiting distributions in (32) and (33) are asymptotically uncorrelated.<sup>10</sup> Then:

$$\begin{aligned}
LM_1 + LM_2(\delta) &= \sigma^{-2} \left[ T^{-1/2} \sum_{t=1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right]^2 \\
&\quad + \sigma^{-2} \delta^{-1} (1 - \delta)^{-1} \left[ T^{-1/2} \sum_{t=1}^{[T\delta]} (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right. \\
&\quad \left. - \delta T^{-1/2} \sum_{t=1}^T (\ln f_t(\beta_t) - \ln g_t(\gamma_t)) \right]^2 + o_{p,\delta}(1) \\
&\implies \mathcal{B}(1)^2 + \delta^{-1} (1 - \delta)^{-1} \mathcal{B}\mathcal{B}(\delta)^2 \quad (34)
\end{aligned}$$

and the result follows by the Continuous Mapping Theorem. ■

**Proof of Corollary 2.** The proof follows from  $\sigma^{-1} T^{1/2} \hat{\mu}(\delta) \implies \left( \delta^{-1} \mathcal{B}(\delta), (1 - \delta)^{-1} [\mathcal{B}(1) - \mathcal{B}(\delta)] \right)'$  using Lemma 4 and arguments similar to those in Proposition 1. Thus,  $\sigma^{-1} H T^{1/2} \hat{\mu}(\delta) \implies \left[ \frac{\mathcal{B}\mathcal{B}(\delta)}{\delta(1-\delta)}, \mathcal{B}(1) \right]'$ , which implies  $W_T(\delta) \implies \frac{\mathcal{B}\mathcal{B}(\delta)^2}{\delta(1-\delta)} + \mathcal{B}(1)^2$ , and the result obtains by applying the Continuous Mapping Theorem. ■

---

<sup>10</sup>See Rossi (2005).

**Proof of Theorem 3.** Let  $X_T(\lambda) = o_{p,\lambda}(1)$  denote  $\sup_\lambda \|X_T(\lambda)\| = o_p(1)$  and  $X_T(\lambda) = O_{p,\lambda}(1)$  denote  $\sup_\lambda \|X_T(\lambda)\| = O_p(1)$ . Let  $\sum_j \equiv \sum_{j=t-m/2+1}^{t+m/2}$  for  $t = m/2, \dots, T - m/2$ . We first show that, under the null hypothesis,  $\sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\hat{\theta}_t^*) = \sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\theta^*) + o_{p,\lambda}(1)$ . Applying a second order Taylor expansion, we have:

$$\sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\hat{\theta}_t^*) \quad (35)$$

$$= \sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\theta^*) + \sigma^{-1} \frac{1}{2} \left\{ \left[ m^{-1} \sum_j \nabla \ln f_j(\hat{\beta}_t^*) \right] \sqrt{m} (\hat{\beta}_t^* - \beta^*) \right. \quad (36)$$

$$\left. - \left[ m^{-1} \sum_j \nabla \ln g_j(\hat{\gamma}_t^*) \right] \sqrt{m} (\hat{\gamma}_t^* - \gamma^*) \right\} \quad (37)$$

$$+ \sigma^{-1} (\hat{\beta}_t^* - \beta^*)' \left[ m^{-1} \sum_j \nabla^2 \ln f_j(\ddot{\beta}_t^*) \right] \sqrt{m} (\hat{\beta}_t^* - \beta^*) \quad (38)$$

$$- (\hat{\gamma}_t^* - \gamma^*)' \left[ m^{-1} \sum_j \nabla^2 \ln g_j(\ddot{\gamma}_t^*) \right] \sqrt{m} (\hat{\gamma}_t^* - \gamma^*), \quad (39)$$

where  $\ddot{\beta}_t^*$  is an intermediate point between  $\hat{\beta}_t^*$  and  $\beta^*$  (and similarly for  $\ddot{\gamma}_t^*$ ). By construction,  $m^{-1} \sum_j \nabla \ln f_j(\hat{\beta}_t^*) = 0$  (and similarly for  $\hat{\gamma}_t^*$ ); by Assumption FB(2),  $\sqrt{T} (\hat{\beta}_t - \beta)$  is  $O_{p,\lambda}(1)$ ; therefore, by Assumption FB(5), (36) is  $o_{p,\lambda}(1)$ , and similarly for (37). Note that

$$\begin{aligned} & \sigma^{-1} (\hat{\beta}_t^* - \beta^*)' \left[ m^{-1} \sum_j \nabla^2 \ln f_j(\ddot{\beta}_t^*) \right] \sqrt{m} (\hat{\beta}_t^* - \beta^*) \\ = & \sigma^{-1} (\hat{\beta}_t^* - \beta^*)' \left[ m^{-1} \sum_j \nabla^2 \ln f_j(\ddot{\beta}_t^*) - E(\nabla^2 \ln f_j(\beta^*)) \right] \sqrt{m} (\hat{\beta}_t^* - \beta^*) \end{aligned} \quad (40)$$

$$+ \sigma^{-1} (\hat{\beta}_t^* - \beta^*)' E(\nabla^2 \ln f_j(\beta^*)) \sqrt{m} (\hat{\beta}_t^* - \beta^*); \quad (41)$$

by Assumptions FB(2,3), (40) and (41) are both  $o_{p,\lambda}(1)$ , and similarly for (39). Thus,

$$\sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\hat{\theta}_t^*) = \sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\theta^*) + o_{p,\lambda}(1). \quad (42)$$

Now write

$$\sigma^{-1}m^{-1/2} \sum_j \Delta L_j(\theta^*) = (m/T)^{-1/2} \left( \sigma^{-1}T^{-1/2} \sum_{j=1}^{t+m/2} \Delta L_j(\theta^*) - \sigma^{-1}T^{-1/2} \sum_{j=1}^{t-m/2} \Delta L_j(\theta^*) \right).$$

By (42) and Assumptions FB(1), FB (4) and FB(5), we have

$$\sigma^{-1}m^{-1/2}\sum_j\Delta L_j(\widehat{\theta}_t^*)\Longrightarrow[\mathcal{B}(\lambda+h/2)-\mathcal{B}(\lambda-h/2)]/\sqrt{h}.$$

The statement in the proposition then follows from consistency of  $\widehat{\sigma}$  for  $\sigma$ . ■

## 9 Tables and Figures

**Table 1. Critical Values for the  
Fluctuation Test ( $k_\alpha$ )**

$h$	$\alpha$	
	0.05	0.10
0.1	3.393	3.170
0.2	3.179	2.948
0.3	3.012	2.766
0.4	2.890	2.626
0.5	2.779	2.500
0.6	2.634	2.356
0.7	2.560	2.252
0.8	2.433	2.130
0.9	2.248	1.950



**Table 2. Monte Carlo: Design 1**

$\beta_A$	Fluctuation	QLR $_T^*$	Break	ExpW $_{\infty,T}^*$	MeanW $_T^*$
0	0.04	0.04	0.05	0.04	0.05
0.1	0.06	0.09	0.10	0.09	0.08
0.2	0.15	0.20	0.24	0.19	0.16
0.3	0.32	0.44	0.51	0.42	0.34
0.4	0.53	0.69	0.75	0.66	0.56
0.5	0.72	0.86	0.90	0.84	0.76
0.6	0.87	0.96	0.97	0.95	0.90
0.7	0.94	0.99	0.99	0.99	0.96
0.8	0.98	1	1	1	0.98
0.9	0.99	1	1	1	1
1.0	1	1	1	1	1
1.1	1	1	1	1	1
1.2	1	1	1	1	1
1.3	1	1	1	1	1
1.4	1	1	1	1	1
1.5	1	1	1	1	1
1.6	1	1	1	1	1
1.7	1	1	1	1	1
1.8	1	1	1	1	1
1.9	1	1	1	1	1
2	1	1	1	1	1

**Table 3. Monte Carlo: Design 2**

$\sigma_A^2$	Fluctuation	QLR $_T^*$	Break	ExpW $_{\infty,T}^*$	MeanW $_T^*$
0	0.04	0.05	0.05	0.05	0.05
0.1	0.05	0.06	0.06	0.06	0.06
0.2	0.06	0.10	0.10	0.10	0.09
0.3	0.08	0.18	0.16	0.16	0.14
0.4	0.10	0.27	0.25	0.25	0.22
0.5	0.18	0.40	0.37	0.38	0.34
0.6	0.25	0.53	0.49	0.50	0.45
0.7	0.34	0.69	0.64	0.66	0.60
0.8	0.46	0.78	0.73	0.76	0.71
0.9	0.55	0.85	0.81	0.83	0.80
1.0	0.64	0.90	0.87	0.89	0.86
1.1	0.74	0.95	0.93	0.94	0.92
1.2	0.81	0.97	0.96	0.97	0.95
1.3	0.88	0.98	0.98	0.98	0.97
1.4	0.92	0.99	0.99	0.99	0.98
1.5	0.95	1	1	1	0.99
1.6	0.97	1	1	1	1
1.7	0.98	1	1	1	1
1.8	0.99	1	1	1	1
1.9	0.99	1	1	1	1
2.0	1	1	1	1	1

## Notes to Tables and Figures

Notes to Table 1. The table reports critical values for the fluctuation test in Proposition 3. Values of  $k_\alpha$  in Table 1 are obtained by Monte Carlo simulations (based on 8,000 Monte Carlo replications and by approximating the Brownian motion with 400 observations) using  $\Pi = [0.15, \dots, 0.85]$ .

Note to Tables 2-3. The tables report empirical rejection probabilities for the ("Fluctuation"), one-time reversal Sup-type (" $QLR_T^*$ "), the  $ExpW_{\infty,T}^*$  and  $MeanW_T^*$  tests. The table also reports empirical rejection probabilities for a standard QLR test for breaks ("Break"). Table 2 reports results for design 1 and Table 3 for design 2 – see Section 5 for details.

Notes to Figure 1. The figure refers to the example in Section 2. The solid line in the figure plots the path of time variation of  $\Delta KLIC$  arising from  $\beta_t$  varying smoothly while keeping the other parameters constant (time varying parameter case, panel a) and from  $\sigma_{x,t}^2$  experiencing a one-time break while keeping the other parameters constant (break in the variance of the regressor case, panel b). Panels c,d show the smoothed  $\Delta KLIC$  (dotted line). In all panels, the large dot shows the global (or average)  $\Delta KLIC$ . Panels e,f show the Fluctuation test statistic (solid line) and the boundary lines (dotted lines) in the time varying parameter case and in the break in the variance of the regressor cases, respectively. Panels g,h plot  $\Delta KLIC$  (solid line) and the one time reversal estimate (dotted line) for the time varying parameter case and in the break in the variance of the regressor cases, respectively.

Notes to Figure 2(a). The figure plots rolling estimates of some parameters in Smets and Wouter's (2002) model. See Smets and Wouters' Table 1, p. 1142 for a description.

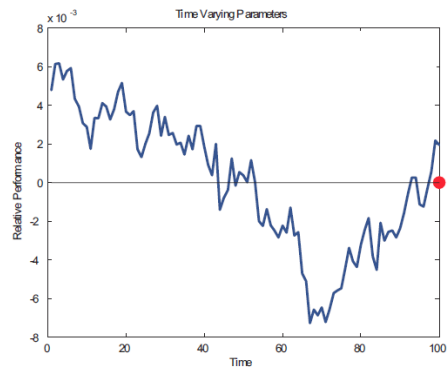
Notes to Figure 2(b). The figure plots rolling estimates of some parameters in Smets and Wouter's (2002) model using full-sample de-trended data. See Smets and Wouters' Table 1, p. 1142 for a description.

Notes to Figure 2(c). The figure plots rolling estimates of the parameters in the monetary policy reaction function described in Smets and Wouters' (2002) eq. (36), given by:  $\hat{R}_t = \rho \hat{R}_{t-1} + (1 - \rho) \left\{ \bar{\pi}_t + r_\pi (\hat{\pi}_{t-1} - \bar{\pi}_t) + r_Y (\hat{Y}_{t-1} - \hat{Y}_t^p) \right\} + r_{\Delta\pi} (\hat{\pi}_t - \hat{\pi}_{t-1}) + r_{\Delta Y} ((\hat{Y}_t - \hat{Y}_t^p) - (\hat{Y}_{t-1} - \hat{Y}_{t-1}^p)) + \eta_t^R$ ,  $\bar{\pi}_t = \rho_\pi \bar{\pi}_{t-1} + \eta_t^\pi$ . The figure plots: inflation coefficient ( $r_\pi$ ), d(inflation) coefficient ( $r_{\Delta\pi}$ ), lagged interest rate coefficient ( $\rho$ ), output gap coefficient ( $r_Y$ ), d(output gap) coefficient ( $r_{\Delta Y}$ ), and standard deviation of the interest rate shock ( $\sqrt{var(\eta_t^\pi)}$ ).

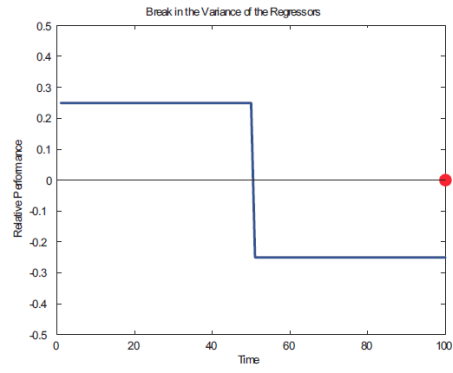
Notes to Figure 3. The figure plots the Fluctuation test statistic for testing equal performance of the DSGE and BVARs, using a rolling window of size  $m = 70$  (the horizontal axis reports the central point of each rolling window). The 10% boundary lines are derived under the hypothesis that the local  $\Delta KLIC$  equals zero at each point in time. The data is de-trended by a linear trend computed over the full sample. The top panel compares the DSGE to a BVAR(1) and the lower panel compares the DSGE to a BVAR(2).

Notes to Figure 4. The figure plots the Fluctuation test statistic for testing equal performance of the DSGE and BVARs, using a rolling window of size  $m = 70$  (the horizontal axis reports the central point of each rolling window). The 10% boundary lines are derived under the hypothesis that the local  $\Delta KLIC$  equals zero at each point in time. The data is de-trended by a linear trend computed over each rolling window. The top panel compares the DSGE to a BVAR(1) and the lower panel compares the DSGE to a BVAR(2).

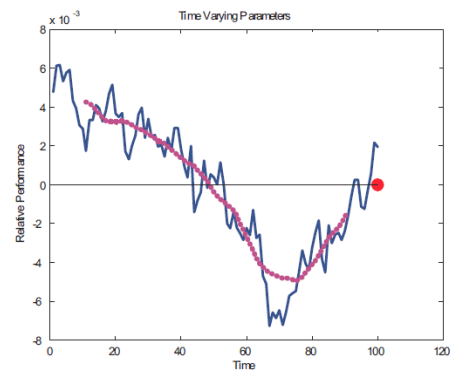
**Figure 1(a)**



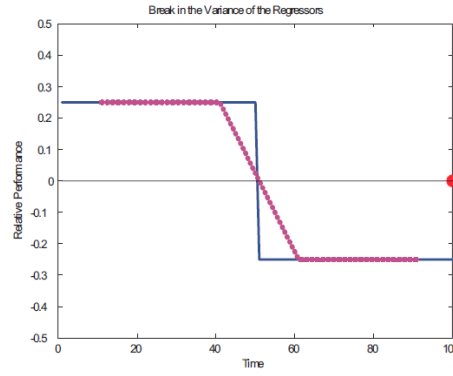
**Figure 1(b)**



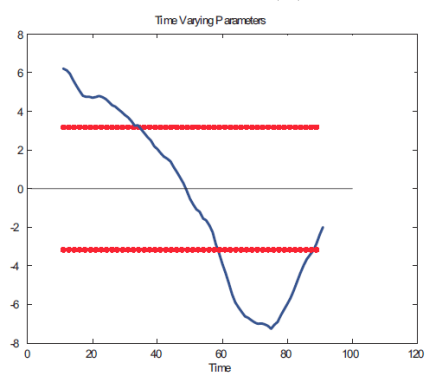
**Figure 1(c)**



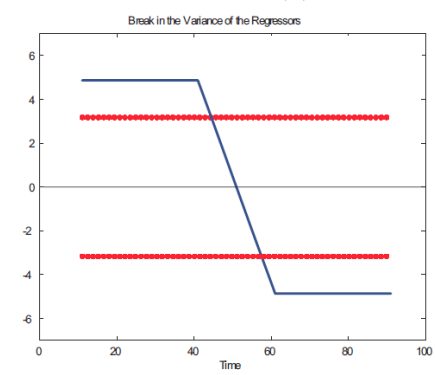
**Figure 1(d)**



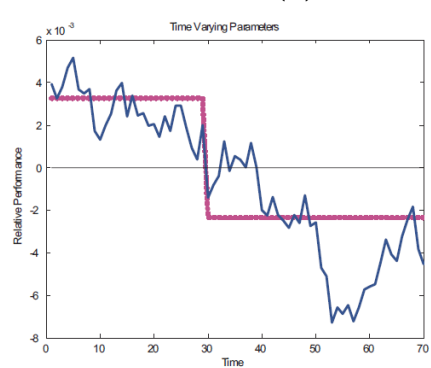
**Figure 1(e)**



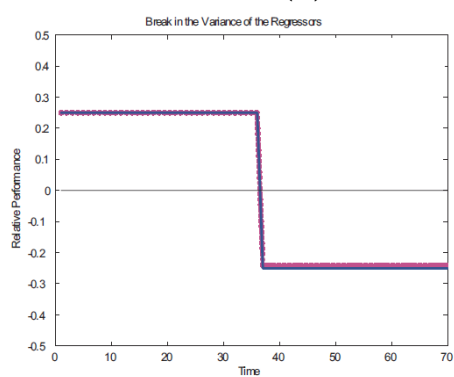
**Figure 1(f)**



**Figure 1(g)**



**Figure 1(h)**



**Figure 2(a). Rolling estimates of DSGE parameters (persistence of the shocks).**

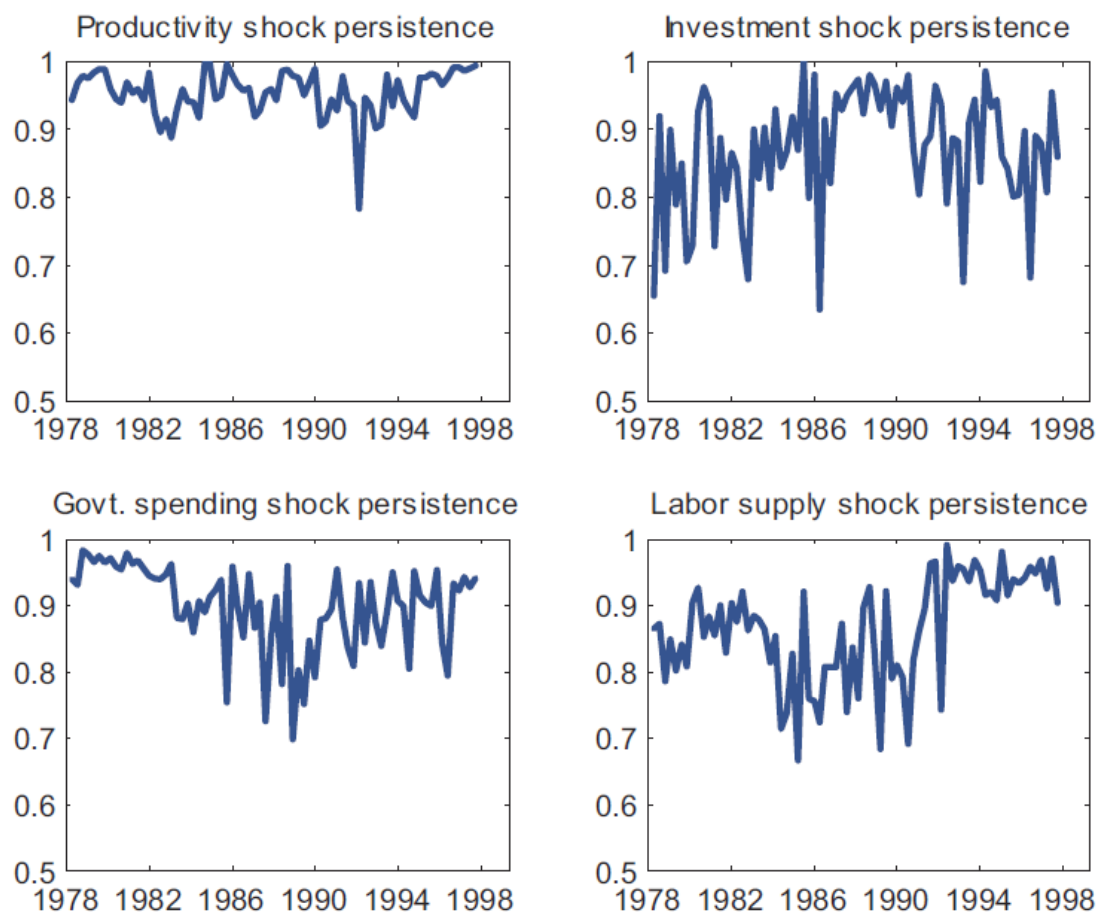


Figure 2(b). Rolling estimates of DSGE parameters (standard deviation of the shocks).

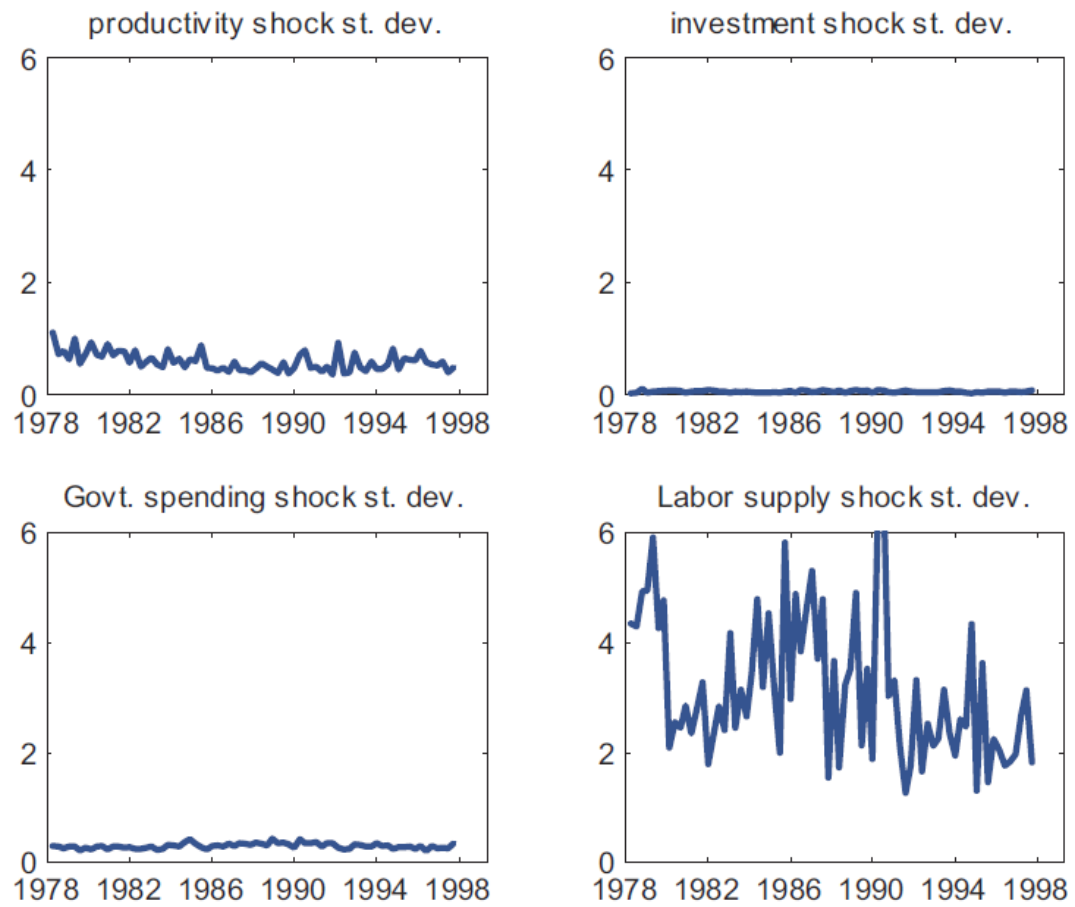
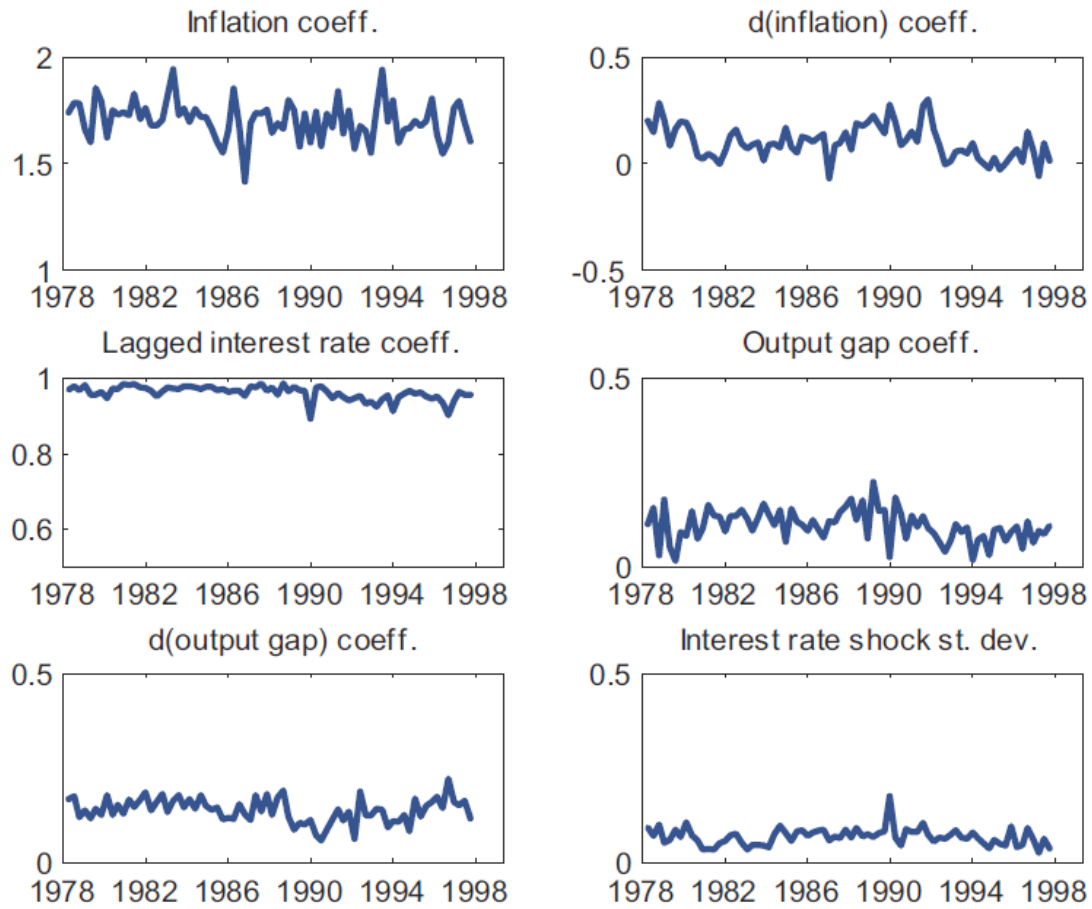
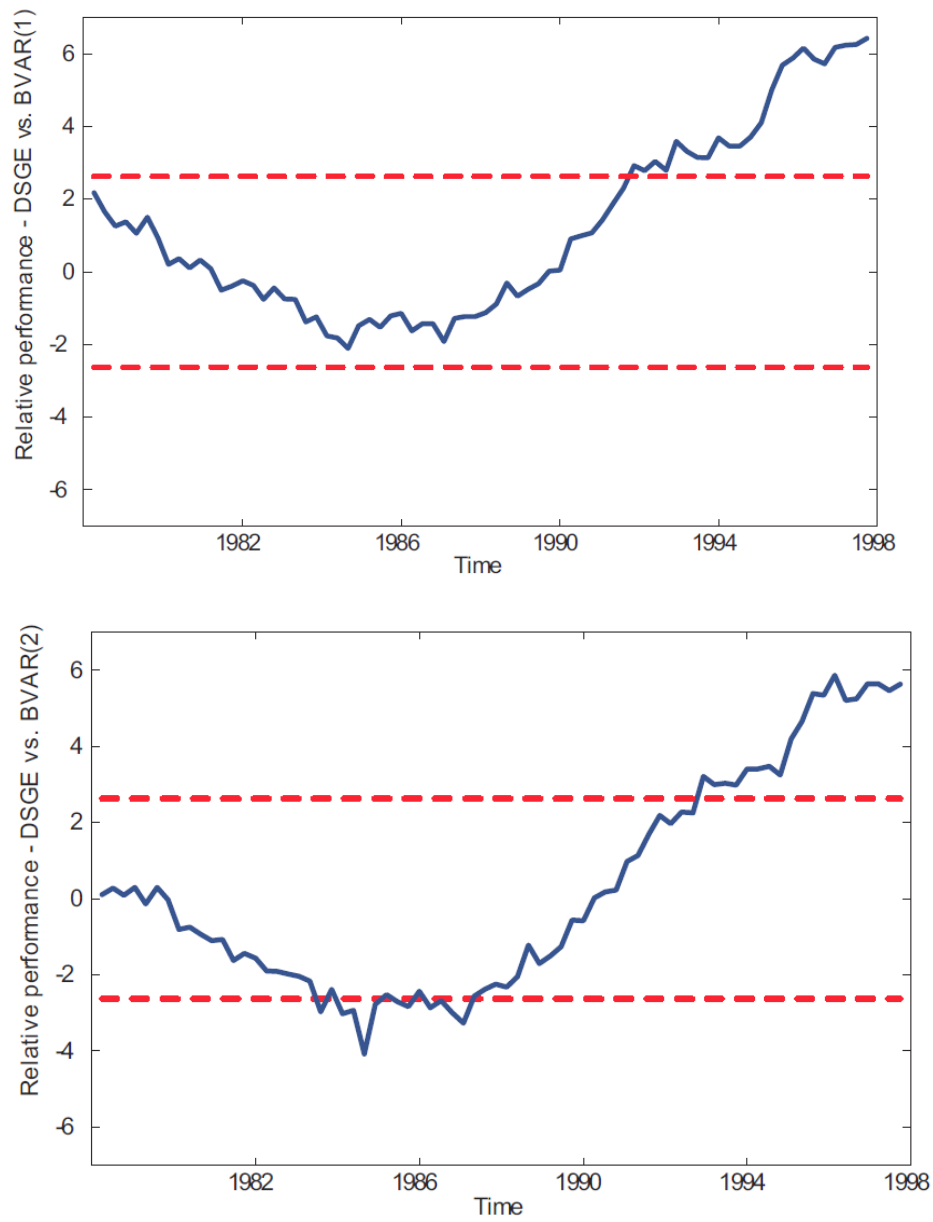




Figure 2(c). Rolling estimates of DSGE parameters (monetary policy parameters).



**Figure 3. Fluctuation test DSGE vs. BVARs. Full-sample detrending**



**Figure 4. Fluctuation test DSGE vs. BVARs. Rolling sample detrending**

