# **State-Space Time Series Modeling of Structural Breaks**

J. Huston McCulloch

**Economics Department** 

Ohio State University

Working Paper # 00-11

June 18, 1997

Rev. Aug. 23, 2000

The author is grateful to participants in the Ohio State University Economics Department

Econometrics seminar for helpful comments and suggestions. Author's e-mail:

mcculloch.2@osu.edu

#### INTRODUCTION

The parameters of time series models often appear to change over time. A popular method of dealing with these changes is to introduce "structural breaks" or "regime shifts." These "breaks" are assumed to occur infrequently, but when they do happen, one or more parameter is customarily permitted to undergo a complete break with its past value. (See, e.g. Perron 1994; Bai, Lumsdaine and Stock 1998).

This note shows, using standard state-space time series modeling techniques, that if such a model is taken literally and estimated by maximum likelihood (ML), it does not give the desired results: The assumption that no new value of the parameter (or parameters) in question is any more likely than any other implies that the ML estimator of the parameter in question should ignore the possibility of structural breaks, no matter how evident they may be in the data. Furthermore, such a model has no predictive power, since the parameter(s) in question could make another complete break with its past at any moment, and take the time series anywhere.

The solution to this problem is to explicitly model the breaks as coming from a distribution that is more likely to give a new value in some vicinity of the old value than arbitrarily far from it. Once this is done, the rate of occurrence of the breaks and/or the parameter(s) of the breakgenerating distribution may be estimated by ML, and the time-changing parameter(s) estimated by an appropriate smoothing algorithm. The time-changing parameter(s) will make the future course of the time series more unpredictable than it would otherwise be, but at least this uncertainty will be probabilistically quantifiable.

#### A SIMPLE "STRUCTURAL BREAK" MODEL

The simplest model that gives rise to the possibility of "structural breaks" is one in which an observed time series  $y_t$  has an unobserved conditional mean  $x_t$  that occasionally undergoes permanent shifts. Let

$$y_t = x_t + \boldsymbol{e}_t, \qquad \boldsymbol{e}_t \sim f(\boldsymbol{e}_t),$$
  

$$x_t = x_{t-1} + \boldsymbol{h}_t, \qquad \boldsymbol{h}_t \sim g(\boldsymbol{h}_t),$$
(1)

where the observation errors  $\varepsilon_t$  and the random shifts  $\eta_t$  have mean zero and are serially and mutually independently distributed. The shift  $\eta_t$  may be thought of as representing a "break" if with some small positive probability  $\lambda$  it is drawn from some distribution with density  $h(\eta_t)$ , and otherwise, with probability  $(1 - \lambda)$ , is zero. This makes  $g(\eta_t)$  a compound distribution, with a mass point with weight  $(1 - \lambda)$  at 0, and density  $\lambda h(\eta_t)$  everywhere else.<sup>1</sup>

In general, system (1) can be estimated using the recursive filtering algorithm due to Sorenson and Alspach (1971) (see also Kitagawa 1987, Harvey 1989: 162-165). The implicit hyperparameter vector  $\theta$  includes  $\lambda$  and any other parameters of f(· ) and h(· ), such as their variances and any shape parameters. The filter is initialized with

$$p(x_1|y_1) = f(y_1 - x_1).$$
(2)

Then, given last period's filter density  $p(x_{t-1}|Y_{t-1})$ , the one-step-ahead predictive density is given by

$$p(x_{t} | Y_{t-1}) = \int_{-\infty}^{\infty} p(x_{t} | x_{t-1}) p(x_{t-1} | Y_{t-1}) dx_{t-1}$$
  

$$= \int_{-\infty}^{\infty} g(x_{t} - x_{t-1}) p(x_{t-1} | Y_{t-1}) dx_{t-1}$$
  

$$= (1 - I) p(x_{t-1} | Y_{t-1}) \Big|_{x_{t-1} = x_{t}} + I \int_{-\infty}^{\infty} h(x_{t} - x_{t-1}) p(x_{t-1} | Y_{t-1}) dx_{t-1}$$
(3)

<sup>&</sup>lt;sup>1</sup> If desired,  $g(\eta)$  may be written  $(1-\lambda)\delta(\eta) + \lambda h(\eta)$ , where  $\delta(\cdot)$  is the Dirac delta function.

where  $Y_t = \{y_1, y_2, ..., y_t\}$ . The new filter density is then given by

$$p(x_t | Y_t) = p(y_t | x_t) p(x_t | Y_{t-1}) / p(y_t | Y_{t-1}) = f(y_t - x_t) p(x_t | Y_{t-1}) / p(y_t | Y_{t-1}),$$
(4)

where the denominator, which is the probability of  $y_t$  conditional on  $Y_{t-1}$ , and therefore the likelihood contribution of  $y_t$ , is given by the integral of the numerator:

$$p(y_{t} | Y_{t-1}) = \int_{-\infty}^{\infty} p(y_{t} | x_{t}) p(x_{t} | Y_{t-1}) dx_{t}$$

$$= \int_{-\infty}^{\infty} f(y_{t} - x_{t}) p(x_{t} | Y_{t-1}) dx_{t}.$$
(5)

The first moment of the filter density gives the expectation of  $x_t$ , conditional on data up to time t, while the hyperparameter vector  $\theta$  may be estimated by numerically maximizing the likelihood, or equivalently, the probability of the entire sample through the last observed time T, conditional on the first observation:

$$L(\boldsymbol{q} \mid Y_T) = p(Y_T \mid y_1) = \prod_{t=2}^T p(y_t \mid Y_{t-1}).$$
(6)

The "smoother" density, which gives the distribution of  $x_t$  conditional on all the data in both directions, may be obtained by backward recursion using the following formula:

$$p(x_{t} | Y_{T}) = p(x_{t} | Y_{t}) \int_{-\infty}^{\infty} p(x_{t+1} | Y_{T}) p(x_{t+1} | x_{t}) / p(x_{t+1} | Y_{t}) dx_{t+1}$$

$$= p(x_{t} | Y_{t}) \int_{-\infty}^{\infty} p(x_{t+1} | Y_{T}) g(x_{t+1} - x_{t}) / p(x_{t+1} | Y_{t}) dx_{t+1}$$

$$= p(x_{t} | Y_{t}) \left( (1 - I) \frac{p(x_{t+1} | Y_{T}) |_{x_{t+1} = x_{t}}}{p(x_{t+1} | Y_{t}) |_{x_{t+1} = x_{t}}} + I \int_{-\infty}^{\infty} \frac{p(x_{t+1} | Y_{T}) h(x_{t+1} - x_{t}) dx_{t+1}}{p(x_{t+1} | Y_{t})} \right)$$
(7)

The expectation of this smoother density and its quantiles provide an ex post point estimate and confidence intervals for  $x_t$ . The estimated value of  $\lambda$  times T gives a point estimate of the number of "breaks" that are present in the sample.

### The Conventional "Structural Break" Model

The conventional "structural break" literature makes no assumption that any value of  $\eta_t$ , apart from 0, is more likely than any other value. The only way this can be true is if  $h(\eta)$  is everywhere 0. This tacit assumption implies that  $g(\eta)$  integrates to 1 -  $\lambda$ , which must be less than unity if breaks are a true possibility, and therefore that  $g(\eta)$  is an improper density.<sup>2</sup>

**Theorem 1:** If the observation errors  $\varepsilon_t$  in (1) are drawn from a normal distribution with mean 0 and some variance  $\sigma^2$ , i.e.

$$f(\boldsymbol{e}) = n(\boldsymbol{e}; 0, \boldsymbol{s}^2), \tag{8}$$

where

$$n(x; \boldsymbol{m}, \boldsymbol{s}^{2}) = \frac{1}{\sqrt{2\boldsymbol{p}\boldsymbol{s}}} \exp\left(\frac{-(x-\boldsymbol{m})^{2}}{2\boldsymbol{s}^{2}}\right),$$
(9)

and the parameter shifts  $\eta_t$  equal 0 with probability 1 -  $\lambda$ , and have zero or essentially zero density elsewhere, then

i. 
$$p(x_t|Y_t) = n(x_t; \overline{y}_t, \mathbf{s}^2 / t),$$
 (10)

ii. 
$$p(x_t|Y_T) = n(x_t; \bar{y}, s^2 / T),$$
 (11)

iii. 
$$L(\boldsymbol{q} \mid \boldsymbol{Y}_T) = p(\boldsymbol{Y}_T \mid \boldsymbol{y}_1) = \frac{1}{\sqrt{T}} \left( \frac{1-\boldsymbol{l}}{\sqrt{2\boldsymbol{p}\boldsymbol{s}}} \right)^{T-1} \exp\left( -\frac{1}{2\boldsymbol{s}^2} \sum_{t=1}^T (\boldsymbol{y}_t - \overline{\boldsymbol{y}})^2 \right),$$
 (12)

where

$$\overline{y}_{t} = \frac{1}{t} \sum_{i=1}^{t} y_{i} , \qquad (13)$$

<sup>&</sup>lt;sup>2</sup> Alternatively,  $h(\eta)$  may be thought of as having say a normal distribution, with a very large yet finite variance. In this case,  $g(\eta)$  is technically proper yet has a density that is essentially 0 except for its mass point at the origin.

and  $\overline{y} = \overline{y}_T$ .

### **Proof:**

See Appendix.

### **Corollary:**

It follows immediately from parts i) and ii) of the Theorem that

$$E(x_t|Y_t) = \overline{y}_t, \tag{14}$$

and

$$E(x_t|Y_T) = \overline{y}, \qquad (15)$$

regardless of the value of  $\lambda$ . In other words, the appropriate ex post estimate of the mean of the distribution that  $y_t$  is being drawn from is simply the average of all the observations, regardless of any appearance of structural breaks. Furthermore, the likelihood iii) above is maximized by

$$\hat{I}_{ML} = 0, \tag{16}$$

again regardless of any appearance of structural breaks.<sup>3</sup>

### **DISCUSSION AND CONCLUSION**

It has been demonstrated that proper ML estimation of a traditional "structural break" model, taken literally, leads to an estimator of the relevant parameter that ignores the possibility of breaks, and to an estimator of the probability of breaks that makes them impossible, no matter how numerous or obvious "regime shifts" may be in the data. This problem arises because of the tacit assumption that any particular sized non-zero regime shift has zero density. This makes a

 $<sup>^{3}</sup>$  It is interesting that the ML estimator of the variance using (6) is the sum of squared residuals about the mean divided by (T-1), rather than by T.

finite regime shift infinitely less likely than a draw of long runs of  $\varepsilon$ 's at different levels, no matter how astronomically unlikely the latter may be.

The solution to this impasse is to explicitly model the breaks as coming from a distribution that gives finite density to  $h(\eta_t)$ . Once this is done, the appropriate estimator of the "breaking" parameter in question is the mean of its smoother density, as computed by structural time series methods.

If  $\lambda$  is simply assumed to be 1 and  $f(\eta_t)$  and  $g(\eta_t)$  are both taken to be Gaussian with mean 0 and some variances  $\sigma^2$  and  $\tau^2$ , a "regime shift" of sorts occurs in every period, with Gaussian density  $h(\eta_t) = g(\eta_t)$ . In this case, the Sorenson/Alspach filter (4) reduces to the familiar Kalman filter, which may be estimated quickly in closed form. The two variances may then be estimated by ML. Unfortunately, this most easily computed case does not allow dramatic regime shifts of the type that are desired.

Once either  $f(\varepsilon_t)$  or  $g(\eta_t)$  is non-Gaussian, as is the case if  $\lambda < 1$ , there is no computational advantage to a Gaussian assumption for the remainder of the system, since the Sorenson/Alspach filter and/or Kitagawa smoother must still be computed numerically. A Laplace (back-to-back exponential) distribution for the breaks, conditional on their occurrance, is easily computed and would allow a wider variety of break sizes to take place than if  $h(\eta)$  were taken to be Gaussian. Naturally, if the observed breaks are sufficiently few in number, it will be nearly impossible to determine anything about the precise shape of their distribution empirically. Nevertheless, any proper distribution is better than  $h(\eta) \equiv 0.^4$ 

<sup>&</sup>lt;sup>4</sup> Another approach that allows occasional dramatic regime shifts, alongside continual minor regime shifts, is to set  $\lambda = 1$  but then to model  $g(\eta_t) = h(\eta_t)$  as a leptokurtic distribution such as the symmetric stable or Student. The degree of leptokurtosis, i.e. the stable characteristic exponent or Student degrees of freedom, is then an additional hyperparameter that needs to be estimated, but this is offset by the fact that  $\lambda$  does not. Such a model, with stable

Once system (1) is properly estimated, a proper out-of-sample predictive densities for  $x_t$ may be obtained by repeated application of (3) for t > T, so long as  $h(\eta_t)$  is non-zero. The predictive density for  $y_t$  is then given by

$$p(y_t | Y_T) = \int_{-\infty}^{\infty} p(y_t | x_t) p(x_t | Y_T) dx_t$$

$$= \int_{-\infty}^{\infty} f(y_t - x_t) p(x_t | Y_T) dx_t.$$
(17)

This of course reflects uncertainty that increases with horizon, and that is greater than would be the case if  $x_t$  were a constant. But if  $h(\eta_t)$  is tacitly assumed to be 0, and the ML estimator of  $\lambda$ overriden with a value  $\lambda^*$  greater than 0,<sup>5</sup> the predictive densities for out-of-sample values of  $x_t$ and therefore  $y_t$  become improper and increasingly uninformative, with a probability equal to only  $(1-\lambda^*)^{(t-T)}$  of any finite value being observed. The mean of such an improper distribution is necessarily undefined.

Having adopted a proper model and estimated it, the presence of breaks can be tested for by constraining  $\lambda$  to be 0 and/or h( $\eta$ ) to have zero scale, and constructing a likelihood ratio statistic for this null hypothesis. Unfortunately, this null hypothesis is on the boundary of the parameter space, and there may also be unidentified nuisance parameters under the null, so that the LR statistic will not necessarily have its customary  $c_{(1)}^2$  distribution. Nevertheless, its distribution may be ascertained by Monte Carlo simulation.<sup>6</sup>

If breaks are present, they may be approximately dated by examining the quantiles of the computed smoother densities. However, the smoother density will not pinpoint the breaks, and,

distributions for both  $f(\varepsilon_t)$  and  $g(\eta_t)$ , has been implemented numerically for bond returns by Oh (1994) and for U.S. inflation by Bidarkota and McCulloch (1998).

<sup>&</sup>lt;sup>5</sup> Equal, say, to the number of apparent breaks divided by T.

like the filter density, it will in fact ordinarily widen in the vicinity of apparent breaks.<sup>7</sup> If desired, a binary break indicator may be added to the system as an additional state variable, and the probability of a break in each period in addition to the filter and smoother densities for the breaking parameter.

The present note deals exclusively with the simple local level model (1), with a single additive unobserved state variable. Nevertheless, its considerations extend directly to more complicated models, in which one or more parameter may undergo permanent or semi-permanent shifts. If there is more than one unobserved continuous state variable, however, direct numerical integration quickly becomes intractable, since it requires evaluating a multiple integral for each point of a multi-dimensional grid at each t = 1, ... T, and then repeating this for each iteration in the likelihood maximization search. In such a case, Monte Carlo integration along the lines proposed by Kitagawa (1996) is the most promising avenue at present.<sup>8</sup>

<sup>&</sup>lt;sup>6</sup> Cp. McCulloch (1997), where LR critical values are simulated for the null hypothesis that the stable distribution characteristic exponent  $\alpha$  takes on its Gaussian boundary value of 2.

<sup>&</sup>lt;sup>7</sup> Cp. Bidarkota and McCulloch (1998).

<sup>&</sup>lt;sup>8</sup> A binary additional state variable, as suggested in the preceding paragraph, would not require an unwieldy double integral, but merely a manageable pair of single integrals.

### APPENDIX

# Lemma 1 (elementary):

The product of two normal densities is a third normal density, times a scaling factor.

Specifically,

$$n(x; \mathbf{m}_{1}, \mathbf{s}_{1}^{2}) n(x; \mathbf{m}_{2}, \mathbf{s}_{2}^{2}) = \frac{\mathbf{s}_{3}}{\sqrt{2\mathbf{p}}\mathbf{s}_{1}\mathbf{s}_{2}} \exp\left(-\frac{1}{2}\left(\frac{\mathbf{m}_{1}^{2}}{\mathbf{s}_{1}^{2}} + \frac{\mathbf{m}_{2}^{2}}{\mathbf{s}_{2}^{2}} - \frac{\mathbf{m}_{3}^{2}}{\mathbf{s}_{3}^{2}}\right)\right) n(x; \mathbf{m}_{3}, \mathbf{s}_{3}^{2}), \quad (A.1)$$

where

$$\frac{1}{\boldsymbol{s}_{3}^{2}} = \frac{1}{\boldsymbol{s}_{1}^{2}} + \frac{1}{\boldsymbol{s}_{2}^{2}}, \qquad \boldsymbol{m}_{3} = \boldsymbol{s}_{3}^{2} \left( \boldsymbol{m}_{1} \frac{1}{\boldsymbol{s}_{1}^{2}} + \boldsymbol{m}_{2} \frac{1}{\boldsymbol{s}_{2}^{2}} \right)$$
(A.2)

**Proof of Lemma 1:** 

$$n(x; \mathbf{m}_{1}, \mathbf{s}_{1}^{2}) n(x; \mathbf{m}_{2}, \mathbf{s}_{2}^{2}) = \frac{1}{2\mathbf{p}\mathbf{s}_{1}\mathbf{s}_{2}} \exp\left(-\frac{1}{2}\left(\frac{(x-\mathbf{m}_{1})^{2}}{\mathbf{s}_{1}^{2}} + \frac{(x-\mathbf{m}_{2})^{2}}{\mathbf{s}_{2}^{2}}\right)\right)$$
(A.3)

while

$$\frac{(x-\mathbf{m}_{1})^{2}}{\mathbf{s}_{1}^{2}} + \frac{(x-\mathbf{m}_{2})^{2}}{\mathbf{s}_{2}^{2}} = x^{2} \left(\frac{1}{\mathbf{s}_{1}^{2}} + \frac{1}{\mathbf{s}_{2}^{2}}\right) - 2x \left(\frac{\mathbf{m}_{1}}{\mathbf{s}_{1}^{2}} + \frac{\mathbf{m}_{2}}{\mathbf{s}_{2}^{2}}\right) + \frac{\mathbf{m}_{1}^{2}}{\mathbf{s}_{1}^{2}} + \frac{\mathbf{m}_{2}^{2}}{\mathbf{s}_{2}^{2}} = \frac{(x-\mathbf{m}_{3})^{2}}{\mathbf{s}_{3}^{2}} + \frac{\mathbf{m}_{1}^{2}}{\mathbf{s}_{1}^{2}} + \frac{\mathbf{m}_{2}^{2}}{\mathbf{s}_{2}^{2}} - \frac{\mathbf{m}_{3}^{2}}{\mathbf{s}_{3}^{2}}.$$
(A.4)

Multiplying and dividing the RHS of (A.3) by  $\sqrt{2ps}_3$  and substituting in (A.4) yields (A.1).

///

# **Proof of Theorem 1:**

We have

$$p(x_1|y_1) = n(x_1; y_1, \boldsymbol{s}^2)$$
(A.5)

This satisfies i) for t = 1. Assume tentatively that i) is true for t-1. The predictive density (3) is then

$$p(x_{t} | Y_{t-1}) = (1 - \mathbf{I})p(x_{t-1} | Y_{t-1})$$
  
=  $(1 - \mathbf{I})n(x_{t}; \overline{y}_{t-1}, \frac{\mathbf{s}^{2}}{t-1})$  (A.6)

By (4), the new filter density for time t is

$$p(x_t | Y_t) = n(y_t; x_t, \boldsymbol{s}^2)(1 - \boldsymbol{l})n(x_t; \overline{y}_{t-1}, \frac{\boldsymbol{s}^2}{t-1}) / p(y_t | Y_{t-1}).$$
(A.7)

By Lemma 1 above,

$$n(y_{t};x_{t},\boldsymbol{s}^{2})n(x_{t};\bar{y}_{t-1},\frac{\boldsymbol{s}^{2}}{t-1}) = n(x_{t};y_{t},\boldsymbol{s}^{2})n(x_{t};\bar{y}_{t-1},\frac{\boldsymbol{s}^{2}}{t-1})$$
$$= \left(\frac{t-1}{2\boldsymbol{p}t\boldsymbol{s}^{2}}\right)^{1/2} \exp\left(-\frac{1}{2\boldsymbol{s}^{2}}\left(y_{t}^{2}+(t-1)\bar{y}_{t-1}^{2}-t\bar{y}_{t}^{2}\right)\right)n(x_{t};\bar{y}_{t},\frac{\boldsymbol{s}^{2}}{t}).$$
(A.8)

Since the third normal density integrates to unity, it follows that the likelihood contribution is just  $(1-\lambda)$  times the scaling factor:

$$p(y_t | Y_{t-1}) = (1 - I) \left( \frac{t - 1}{2pts^2} \right)^{1/2} \exp\left( -\frac{1}{2s^2} \left( y_t^2 + (t - 1)\overline{y}_{t-1}^2 - t\overline{y}_t^2 \right) \right),$$
(A.9)

and that the filter density is this third normal density itself:

$$p(x_t|Y_t) = n(x_t; \overline{y}_t, \frac{\boldsymbol{s}^2}{t}).$$
(A.10)

This completes the proof of i) by induction. Part iii) follows immediately from (5).

The proof of ii) proceeds by a backward induction that begins by confirming that the proposition in question is valid is true for t = T, and then demonstrates that if it is valid for any t+1, it is also valid for t. By i), ii) is valid for t = T. Suppose that ii) is valid for t+1. Then by (7) with  $h(\cdot) \equiv 0$ , and using (A.6),

$$p(x_t | Y_T) = n(x_t; \overline{y}_t, \frac{\boldsymbol{s}^2}{t})(1 - \boldsymbol{I}) \frac{n(x_t; \overline{y}, \frac{\boldsymbol{s}^2}{T})}{(1 - \boldsymbol{I})n(x_{t+1}; \overline{y}_t, \frac{\boldsymbol{s}^2}{t})\Big|_{x_{t+1} = x_t}}$$
(A.11)  
$$= n(x_t; \overline{y}, \frac{\boldsymbol{s}^2}{T}).$$

This completes the proof of ii), and therefore of Theorem 1.

///

#### REFERENCES

- Bai, Jushan, Robin L. Lumsdaine, and James H. Stock (1998). Testing for and dating common breaks in multivariate time series. *Review of Economic Studies* 63 : 395-432.
- Bidarkota, Prasad V., and J. Huston McCulloch (1998). Optimal univariate inflation forecasting with symmetric stable shocks. *Journal of Applied Econometrics* **13** : 659:70.
- Harvey, Andrew C. (1989). Forecasting, Structural Time Series Models and the Kalman Filter.Cambridge University Press.
- Kitagawa, Genshiro (1987). Non-Gaussian state-space modeling of nonstationary time series. *Journal of the American Statistical Association* 82: 1032-41. See also comments and rejoinder, pp. 1041-63.
- \_\_\_\_\_\_ (1996). Monte Carlo filter and smoother for non-Gaussian nonlinear state space models. *Journal of Computational and Graphical Statistics* **5**: 1-25.
- McCulloch, J. Huston (1997). Measuring tail thickness to estimate the stable index α: A critique. Journal of Business & Economic Statistics 15: 74-81.
- Oh, Chang-Seok (1994). Estimation of Time Varying Term Premia of U.S. Treasury Securities: Using a STARCH Model with Stable Distributions. Ph.D. dissertation, Ohio State Univ.
- Perron, Pierre (1989). The great crash, the oil price shock, and the unit root hypothesis. *Econometrica* **57**: 1361-1401.
- Sorenson, H.W., and D.L. Alspach (1971). Recursive Bayesian estimation using Gaussian sums. *Automatica* **7**: 465-79.