# An Analysis of US Employment/Population

Rafael Amaral Porsani

## **1** Introductory Comments

The last US recession ended in the second-quarter of 2009, and the unemployment rate has since fallen considerably (see Figure 1). Indeed, unemployment is close to the lowest level it's been in the last 40 years, suggesting at first sight the predominance of very healthy labor-market conditions in the country. One could reasonably conjecture, however, that such precipitous drop may potentially be attributed to the simple fact that, more and more, individuals are dropping out of the labor-force. In truth, the unemployment rate - according to the Bureau of Labor Statistics (BLS)<sup>1</sup> - is measured as the estimated number of unemployed individuals (people who are simultaneously jobless, looking for a job, and available for work) divided by the estimated labor force (which is comprised of unemployed plus employed individuals<sup>2</sup>). A decrease in the number of individuals looking for a job leads to a decrease in both the number of unemployed individuals (the numerator of the unemployment rate), and to a similar decrease in the labor force (the denominator of such rate), ultimately producing a decline in the rate itself<sup>3</sup>.

Arguably, therefore, a more indicative measure of the overall health of the labor market is the employment-to-population ratio (illustrated in Figure 2). Here, we see that US workers have been, relatively speaking, having a hard time finding jobs - employment is at a much lower level than it was in the late 1990's. The 4% decline seen over the last recession was massive vis-a-vis declines registered in previous recessions, and recovery has been somewhat sluggish - i.e., if compared to the growth in employment seen in the 1980's. Figure 3 highlights in different colors what arguably could be considered different regimes in the employment-to-population series. This figure suggests that such

<sup>&</sup>lt;sup>1</sup>See https://www.bls.gov/cps\_htgm.htm.

<sup>&</sup>lt;sup>2</sup>In the section entitle "Brief Description of Data", below, we provide a more comprehensive summary of which individuals are considered to be employed by the BLS, and on how employment-related data is collected.

<sup>&</sup>lt;sup>3</sup>Note that this follows from a basic property of fractions. If a, b and x are non-zero positive integers, with a < b, then  $\frac{a-x}{b-x} < \frac{a}{b}$ .

series grew at a fast pace from about April 1975 to November 2000. In Figure 4, we show linear trends fit to the time periods corresponding to these potential regimes. It should be clear by assessing the slopes of these lines that employment did indeed grow at a stronger pace in the period that goes from 1975 to 2000.

In this study, we model and interpret the time series of employment-to-population from April 1975 to November 2000 (we leave the period that goes from December 2000 to November 2001 as a test set, which we later use to evaluate our forecasts). As we have argued above, this is a period when employment exhibited strong gains. Hopefully, by understanding it better, one could possibly develop recipes/policies for improving the current labor environment in the US. In our view, modeling it is a first step towards gaining such understanding<sup>4</sup>.

# 2 Brief Description of Data

Our data comes from the Federal Reserve Bank of Saint Louis (FRED)<sup>5</sup>. FRED provides, free-ofcharge, about 469,000 US and International time series, from numerous sources, in its website. We have actually downloaded our time series from FRED directly into R using the function "getSymbols" from the "quantmod" library. The data we use is monthly, and the observations originally come from the Current Population Survey (CPS) conducted early each month by the BLS. This survey, which measures employment and unemployment in the country, has been conducted in the US every month since 1940 (though the data that is available on FRED starts in 1948). The BLS reports that there are roughly 60,0000 eligible households in the sample for this survey. The sample is selected so as to be representative of the US population. Every month, government employees contact these 60,000 households and ask individuals living in them (aged 16 and above) questions about their labor activities. Individuals are considered employed if they: (i) "did any work for pay or profit during the week when the survey is conducted"; or (ii) have a job but couldn't work because of specific circumstances<sup>6</sup>. Note that part-time workers, according to this definition, are considered to be employed. And also note that

<sup>&</sup>lt;sup>4</sup>Side note: in Figure 5, we show the estimated total number of employees in the Manufacturing sector. This is included more as a curiosity here. Manufacturing jobs took a hard hit in the last recession, falling by about 30% then. Interestingly, the Trump campaign was apparently successful in identifying this movement - in the sense that they won the electoral vote in a number of states in which the discussion Manufacturing jobs seemed to have been critical.

<sup>&</sup>lt;sup>5</sup>https://fred.stlouisfed.org/.

<sup>&</sup>lt;sup>6</sup>In particular, if they have a job but were on vacation; ill; experiencing child care problems; on maternity or paternity leave; prevented from working due to bad weather; involved in a labor dispute; or taking care of personal of family obligation.

if a person has more than one job, he or she will simply show up in the survey as being employed - i.e., s/he won't be "counted twice". The employment-to-population ratio is then defined as the number of employed individuals, as a percentage of the total population (in the surveyed households).

# 3 Analysis

To model our data, the following steps will be taken: (i) we will detrend the data; (ii) we'll use spectral analysis to identify cycles/seasonal patterns in the detrended data, and then remove these cycles from it; (iii) the residuals (*original data - trend - seasonal components*) will then be analyzed and modeled as an ARMA process; (iv) finally, we'll evaluate how well we have modeled our time series by evaluating the residuals engendered by our approach, and by assessing the quality of our forecasts.

Before implementing these steps, we first assess the autocorrelation (ACF) and partial-autocorrelation functions (PACF) of the original data. On the top portion of Figure 6, we show the ACF and PACF of the original data when all possible lags are used; on the bottom of Figure 6, the ACF and PACF are depicted for a total of 50 lags. We see that the PACF hints towards the existence of cycles in the data; and so does the ACF plot which includes all possible lags. As we'll show below, the cyclical nature of our data will become more apparent once the series has been detrended.

### 3.1 Detrending the Data

Firstly, we tried fitting a linear regression model to the data. The fitted model takes on the following form:  $\hat{Y}_t = -520 + 0.29t$ , where  $\hat{Y}_t$  denotes the fitted employment-to-population ratio at time t. The t - statistic for our slope coefficient is 35.03, indicating that the trend in this period is statistically significant (even at the 1% level of significance). This fitted linear trend suggests that employment-topopulation tended to increase by 0.29% per year from April 1975 to November 2000 - a rate which we have argued is higher than those evidenced in neighboring time-periods.

Figure 7 shows the detrended data, when such linear trend is used. Interestingly, this figure suggests that the resulting detrended series is not stationary: the variance on the left-hand site of our plot seems markedly higher than that on the right-hand side. To confirm our visual intuition, we conduct an augmented Dickey-Fuller (ADF) test on the detrended series<sup>7</sup>. The p-value from our test is 0.36 - therefore, we fail to reject the null hypotheses that there exists a unit root in this data. This in turn

<sup>&</sup>lt;sup>7</sup>The adf.test function in R selected a lag order of 6 for this test.

suggests that the detrended series is indeed non-stationary.

One could attempt to get around this by taking first-differences of the original data; but our main interest here is in the employment-to-population ratio itself, not in first-differences of this quantity. Thus, we will refrain from taking first-differences of our time-series. Instead, we opt to fit a non-linear trend to the data, using a smoothing-spline regression (i.e., a spline regression where the knots are tied to the observations themselves).

We have written a function in R that fits a non-linear trend to our data through smoothing-splines. Note that we could also have utilized the lowess or loess functions in R to fit a non-linear trend here, as well as other methods. One of the advantages of using our own code to fit a non-linear trend is that, by doing so, we know exactly how the trend is being fit - and this allows us to more easily forecast the trend component of our series.

The fitted trend is of the form  $f(x) = \alpha_0 + \sum_{j=1}^p \alpha_j max(0, x - k_j)$ , where x here is equivalent tot (time), and the  $k'_j s$  are the times associated with the observations in the employment-to-population series<sup>8</sup>. The fitted  $\alpha' s$  are found by minimizing  $\sum_{t=1}^n ||y_t - \alpha_0 - \sum_{j=1}^p \alpha_j max(0, x_t - k_j)||^2 + \lambda \sum_{j=1}^p \alpha_j^2$ , with  $y_t$  denoting the employment-to-population ratio at time  $t, x_t$  is again simply the time t associated with observation  $y_t$ ;  $\lambda$  is a tuning parameter that penalizes the complexity of our model<sup>9</sup>. In the appendix, we provide a brief derivation of the spline solution.

Figure 8 shows our time series and the non-linear trend. The tuning parameter  $\lambda$  was selected so as to yield a detrended time series that looked stationary, while producing a trend component that was reasonably smooth<sup>10</sup>. The detrended data is shown in Figure 9. We can see that after removing the non-linear trend, we obtain a series that seems stationary. An ADF test was conducted on this time series, producing a very low p - value (p < 0.01)<sup>11</sup>- the null hypothesis that there exists a unit root in this series can now be rejected. This is the detrended series that we use in our analysis.

Some comments are in order. Firstly, one could speculate that the overall positive trend seen in this period is associated with an increased integration between the US economy and other economies (globalization), which may have led to increasing exports and the creation of local job openings. It is also possible that an increasing participation of women in the labor-force contributed to the steep increase in employment we saw from 1975 to 2000. Lastly, automation and product innovations may

<sup>&</sup>lt;sup>8</sup>Note that p = 308, the total number of observations.

<sup>&</sup>lt;sup>9</sup>Also, n = 308 and p = 308; and  $k_j$  denotes, as we mentioned, the time associated with observation j.

<sup>&</sup>lt;sup>10</sup>We tried different values for  $\lambda$ , ultimately choosing  $\lambda = 100$ .

<sup>&</sup>lt;sup>11</sup>Again a lag order of 6 was used in the test.

destroy jobs in some industries, but can also create jobs in others (and note that these can also make companies more efficient, potentially leading to greater shareholder wealth, which can also spur consumption and jobs) - thus, automation by itself (which one could proxy through labor productivity) may also have helped create jobs in this period. We believe these should be subject to scrutiny in further studies.

### 3.2 Identifying and Removing Cycles (Spectral Analysis)

Figure 10 depicts the ACF and PACF of the detrended data - we can clearly see, from this figure, that the detrended series is influenced by cycles. Figure 11 illustrates the raw periodogram of the detrended data. A large portion of the total variation in the data is explained by the 1.013 frequency (which corresponds to a period of roughly one year). There are also peaks at  $\omega = 1.988 \approx 2$  and  $\omega = 3$ , harmonics of  $\omega = 1$ . This indicates that the one-year cycle present in the series is not a perfect sinusoid. Also note that there is considerable power at lower frequencies. In particular, we see a small peak at the frequency of 0.118 (which corresponds to a period of 5.33 years, or 64 months). One could associate this peak with long-term business-cycle fluctuations not captured by our non-linear trend.

In Figure 12, we show in detail the yearly seasonal component, which, as we have argued, is responsible for a good portion of the total variation in the detrended series. The graph suggests that there is a strong employment season from March to May, and a weak one from October to January. Generally speaking, April seems to be the best month to find a job, while October is apparently the worst. The yearly seasonal component declines, in general terms, after April, and starts to increase after October. July seems to actually be a stronger month for employment-population than June, August and September. Also note that the annual seasonal component is negative from September to February, staying positive in other months. Summer vacations of hiring managers may be a possible reason for the slowdown in Summer months; a desire by some CFO's to boost calendar-year-end bottom lines<sup>12</sup>, coupled with the holiday season, might help to explain the low year-end figures.

We have also computed approximate approximate confidence intervals for the spectral densities associated with the one-year cycle, the cycles corresponding to  $\omega = 2$  and  $\omega = 3$  (harmonics), and the 5.33 year cycle. The lower values of these intervals are all higher than most of the other periodogram ordinates, indicating that these cycles are all significant.

<sup>&</sup>lt;sup>12</sup>Note that the calendar year is used as the fiscal year by most publicly-traded US firms.

Figure 13 contains the smoothed periodogram of our detrended series - obtained using a modified Daniell (2,2) kernel. This periodogram is smoother than the raw one, as expected. The peaks occur at exactly the same frequencies we identified previously ( $\omega = 0.118, 1.013, 1.988$  and 3). The smoothed periodogram is specially helpful to confirm the low-frequency peak at  $\omega = 0.118$  - given in the raw periodogram we saw reasonable power around this frequency.

Moreover, we have utilized an AR spectral estimator to further substantiate our findings. Figure 14 demonstrates that the optimal number of lags to use with such estimator, when using AIC as a criteria, is 39<sup>13</sup>. Figure 15 shows that this procedure produces peaks at nearly the same frequencies we had identified previously. There are clear peaks at  $\omega \approx 1,2$  and 3, as before. The low frequency peak is now located at  $\omega = 0.192$ , indicating a cycle with period of roughly 5.2 years - which is very close to the 5.33 years period that was uncovered with the raw and smoothed periodograms.

The ACF of the detrended data, after removal of the one-year cycle, is provided in Figure 16. We can see from this figure that the series still exhibits a cyclical pattern, as expected, when only the one-year cycle is removed from it. Figure 16 also illustrates the ACF of the detrended data after the removal of both the one-year and the 5.33 years (64 months) cycles. Even when influences from these two cycles are removed, we still see that the resulting series possesses a cyclical behavior. Thus, in the analysis that follows, we have opted to remove the influences of all of the four cycles that are evident in our periodograms, and which have found to be significant: the one-year cycle, the 5.33 years one, and the cycles associated with the second and third harmonics of  $\omega = 1$  (i.e.,  $\omega = 2$  and  $\omega = 3$ ; note that removing these helps us to account for the non-sinusoidal behavior of the yearly component).

#### 3.3 Modeling Residuals with ARMA models, Diagnostics and Forecasting

Next, we turn our attention to our residual data (the series obtained after the trend and all seasonal components have been removed). Figure 17 shows the ACF and PACF of the residual data. On the top portion of this figure, we see that the ACF seems to tail off (albeit it alternates between positive and negative territory), while the PACF apparently cuts off after a certain number of lags. Closer inspection (see bottom portion of Figure 17) of the PACF suggests that the estimated partial-autocorrelation function cuts off at about lag 14 or 15. When combined, these findings suggest fitting an AR(14) or AR(15) model to the data. We confirm our visual intuition by assessing what would be

<sup>&</sup>lt;sup>13</sup>Note that the optimal number of lags, when AIC is used, is not that clear here. There are many lags with AIC values very close to the one obtained with p = 39. We have used other values for p in our analysis, obtaining similar results.

the best AR model, from an AIC standpoint, to fit to the residual data. Figure 18 demonstrates that the lowest AIC (-2.29) is obtained when 14 lags are used, confirming the intuition we developed from evaluating the ACF and PACF of our series.

In Figure 19, model diagnostics are presented for the AR(14) model. For the most part, an AR(14) model does a good job in describing our residual data. An inspection of the normal Q-Q plot and the histogram of residuals produced by this model indicates that model residuals are normally distributed. A Shapiro-Wilk normality test conducted on these yields a p-value of 0.58, further supporting this assertion. The Ljung-Box test, however, suggests there may be serial correlation in the residuals from lags 12 to 16.

As a result, we have tried to include MA terms in our model, fitting different ARMA models to the data. Good results were obtained when an ARMA(14,9) is used (see Figure 20). Residuals from the ARMA(14,9) model look normal<sup>14</sup>; they resemble a white-noise process (see ACF in Figure 20); and Ljung-Box tests provide no evidence of serial correlation in them. Furthermore, the AIC associated with this model was -2.64 - a value that is lower than the one we obtained with the AR(14) model.

Lastly, Figures 21 to 23 show forecasts obtained when an AR(14) model and an ARMA(14,9) are used to model the residual series. Forecasts were computed as the sum of the trend and seasonal components, plus predictions from either the AR(14) or ARMA(14,9) model. For the AR(14) model, we also show one-step-ahead forecasts (Figure 23). Both models seem to perform well in the testing period (December 2000 to November 2001), with most realizations falling withing two-standard errors from our predictions.

# 4 Concluding Remarks

We speculate that different variables may have contributed to the robust gains in employment evidenced in the period we studied. Firstly, the US economy became more integrated with other economies during this period (globalization). Growing exports may have helped spur job growth over these years. Secondly, it is possible that an increase in the participation of women (and other groups) in the labor force may have also helped boost jobs then. We have also speculated that automation itself (which could be measured through labor productivity, i.e., output per hours worked) may also have contributed to an

 $<sup>^{14}</sup>$  The Shapiro-Wilk normality test here yielded a p-value of 0.105 - so we fail to reject the null that the data comes from a normal distribution, at a 5% level of significance.

increase in the level of employment-to-population in this period: while automation may destroy jobs. it also has the potential to make companies more efficient, contributing to an increase in shareholder wealth, possibly thus having an impact in household consumption (via a wealth effect - households may "feel richer" and consume more), leading then to more employment. Lastly, this was also a period in which a good number of skilled foreign workers came to the US. These may have helped the economy grow further, imparting their knowledge on local workers, and going on to start-up new companies by themselves. In future studies, it would be interesting to investigate how these forces may have driven employment in this period. This could be done, at first, by for example looking at the cross-spectrum and coherence of employment and other variables, such as exports, labor productivity, and women labor-force participation. Potentially, having an immigration policy that incentivizes skilled workers to join the work force in the US might be a good idea? Maybe creating incentives for women to become entrepreneurs and/or business leaders should be a priority of our government? These questions warrant further scrutiny. Indeed, one of the limitations of our study is that we did not use other time series to explain employment-to-population rates. As a final remark, it is also interesting to see that the detrended series can be largely explained by a one-year seasonal component/cycle. Having more time, we would like to investigate the reasons underscoring the predominance of a weak employment season during the last months of the year (we have speculated as to why this may be the case)<sup>15</sup>, and study why hiring in March, April and May is apparently stronger than in other months of the year.

<sup>&</sup>lt;sup>15</sup>Note that the BLS argues that, in their view, extreme weather conditions do not affect employment by much - though they likely affect the statistic "number of hours worked". In our view, a more likely reason for these low numbers in the last months of the year is financial: CFO's may be looking to improve their end-of-year numbers, thus curtailing hiring in this period.

# 5 APPENDIX

## 5.1 Figures

Figure 1: Unemployment Rate in the US - Recessions in Blue. The unemployment rate suggests the economy has almost never been better.



Figure 2: Employment-to-Population Ratio - Recessions in Blue. Employment is at a much lower level than it was in 2001.





Figure 3: Different Regimes in Employment/Population. We essentially model the data in Green - April 1975 to November 2000 - leaving December 2000 to November 2001 as a test set.

Figure 4: Comparing Trends. There is a steep increase in employment from 1975 to 2000. Blue Area Green Area (Until Nov. 2000)





Figure 5: Total Number of Employees in Manufacturing (in thousands)- Recessions in Blue.





Figure 7: Detrended Data - Linear Trend. This figure suggests that the detrended data (when a linear trend is used) is non-stationary - its variance is not constant. Trend: Y = -520 + 0.29t; Augmented Dickey-Fuller Test - p-value = 0.36.



Figure 8: Data + Non-Linear Trend. We estimate a non-linear trend using smoothing splines:  $f(x) = \alpha_0 + \sum_{j=1}^p \alpha_j max(0, x-k_j)$ . The solution is found by minimizing:  $\sum_{i=1}^n ||y_i - \alpha_0 - \sum_{j=1}^p \alpha_j max(0, x_i - k_j)||^2 + \lambda \sum_{j=1}^p \alpha_j^2$ . Observations themselves are the  $k'_j s$ .



Figure 9: Detrended Data - Non-Linear Trend. The data now looks stationary. Also: Augmented Dickey-Fuller Test - p-value < 0.01.

![](_page_12_Figure_1.jpeg)

Detrended data - non-linear trend

![](_page_12_Figure_3.jpeg)

![](_page_12_Figure_4.jpeg)

![](_page_12_Figure_5.jpeg)

Detrended Data - 50 Lags

300

![](_page_12_Figure_7.jpeg)

Figure 11: Raw Periodogram. We see a one Year Cycle + 5.3 Year Cycle, and Harmonics of  $\omega = 1$ . Red dotted line indicates lower bound of the confidence interval for the spectrum corresponding to  $\omega = 1$ . Peaks are located at  $\omega = 0.188, 1.013, 1.988, 3$ . These correspond to the following periods (in years), respectively: 5.33, 0.99, 0.50, 0.33.

![](_page_13_Figure_1.jpeg)

![](_page_13_Figure_2.jpeg)

![](_page_13_Figure_3.jpeg)

Figure 13: Smoothed Periodogram - Modified Daniell(2,2). Low frequency peak at  $\omega = 0.188$  - same as in the raw periodogram.

![](_page_14_Figure_1.jpeg)

Figure 14: Autoregressive Spectral Estimator (1). Lowest AIC: p = 39.

![](_page_14_Figure_3.jpeg)

Figure 15: Autoregressive Spectral Estimator (2). Low frequency peak at  $\omega = 0.192$ ; Corresponds to period of 5.2 years.

![](_page_15_Figure_1.jpeg)

Figure 16: ACF and PACF of Detrended Data, After Removal of Select Seasonal Components. ACF still indicates there are seasonal patters left in the data, if we don't remove all cycles from it.

![](_page_15_Figure_3.jpeg)

![](_page_15_Figure_4.jpeg)

![](_page_15_Figure_5.jpeg)

Cycle removed: 1 year and 5.33 years (64 months)

Figure 17: ACF and PACF of Residual Data (Trend and All Seasonal Components Have been Removed). ACF seems to tail off and PACF seems to cut off at lag 14 or 15, suggesting fitting either an AR(14) or AR(15) model to the residual data.

![](_page_16_Figure_1.jpeg)

Figure 18: AIC and BIC of possible AR models. Minimum AIC at p = 14 (AIC = -2.29).

![](_page_16_Figure_3.jpeg)

![](_page_17_Figure_0.jpeg)

Figure 19: Diagnostics - AR(14) Model. Looks like we are mostly good, except for a bit of serial correlation in the residuals.

Figure 20: Diagnostics - ARMA(14,9) Model. Residuals resemble normally-distributed white noise. AIC = -2.64.

![](_page_17_Figure_3.jpeg)

![](_page_18_Figure_0.jpeg)

Figure 21: Forecasts using AR(14):red line indicates forecasts; blue lines indicate 2 standard-error bounds.

Figure 22: Forecasts using ARMA(14,9):red line indicates forecasts; blue lines indicate 2 standard-error bounds.

![](_page_18_Figure_3.jpeg)

![](_page_19_Figure_0.jpeg)

![](_page_19_Figure_1.jpeg)

### 5.2 Math Appendix - Derivation of Solution to Penalized Spline Regression

Suppose we have training examples of the type  $(y_i, x_i)$ , i = 1, ..., n. Let x be one-dimensional (for example, x can be represent t = time). Also assume we have knots  $k_j$ , j = 1, ..., p. Suppose we try to fit a linear spline of the form  $f(x) = \alpha_0 + \sum_{j=1}^p \alpha_j max(0, x - k_j)$  to the data, by minimizing  $\sum_{i=1}^n ||y_i - \alpha_0 - \sum_{j=1}^p \alpha_j max(0, x_i - k_j)||^2 + \lambda \sum_{j=1}^p \alpha_j^2$ .

We look for:

$$\hat{\alpha} = argmin_{\alpha_0, \{\alpha_j\}_{j=1...p}} \sum_{i=1}^n ||y_i - \alpha_0 - \sum_{j=1}^p \alpha_j max(0, x_i - k_j)||^2 + \lambda \sum_{j=1}^p \alpha_j^2$$

Let  $\hat{\alpha}$  now denote the  $(p+1) \times 1$  vector containing our parameter estimates.

Notice that we can write this problem as:

$$\hat{\alpha} = argmin_{\alpha}||Y - Z\alpha||^2 + \lambda \alpha^T D\alpha,$$

Where D is a  $(p+1) \times (p+1)$  matrix, with 1 in all its diagonal elements aside from the first one, which is 0, and with remaining elements equal to 0. In other words,  $D_{rc} = 1$  if both r = c and  $r, c \neq 1$ ; and 0 otherwise (for r,c $\in$ {1...p+1}):

$$D = \begin{bmatrix} 0 & 0 \\ 0 & I_p \end{bmatrix}$$

And Z is a  $n \times (p+1)$  matrix, whose entries are:

$$Z = \begin{bmatrix} 1 & max(0, x_1 - k_1) & max(0, x_1 - k_2) & \dots & max(0, x_1 - k_j) \\ 1 & max(0, x_2 - k_1) & max(0, x_2 - k_2) & \dots & max(0, x_2 - k_j) \\ \dots & \dots & \dots & \dots & \dots \\ 1 & max(0, x_{n-1} - k_1) & max(0, x_{n-1} - k_2) & \dots & max(0, x_{n-1} - k_j) \\ 1 & max(0, x_n - k_1) & max(0, x_n - k_2) & \dots & max(0, x_n - k_j) \end{bmatrix}$$

Taking the derivative with respect to  $\alpha$  and setting it to zero, we get the FOC:

$$-2Z^{T}(Y - Z\alpha) + 2\lambda D\alpha = 0$$
$$-Z^{T}Y + Z^{T}Z\alpha + \lambda D\alpha = 0$$
$$(Z^{T}Z + \lambda D)\alpha = Z^{T}Y$$

Which implies:

1

$$\hat{\alpha} = (Z^T Z + \lambda D)^{-1} Z^T Y$$

### 5.3 Appendix: R Code

```
\mathbf{2}
 3 library(ggplot2) # nice-looking plots
 4 library(quantmod) # reads in data from Fred
 5 library(scales) # customize x axis in ggplots
 6 library (reshape2) # plot multiple ts with ggplot
 7 library(gridExtra) # plot multiple graphs with ggplot (like par(mfrow) command)
 8 library (tseries)
 9 library (forecast) # Acf plots that ommit lag 0
10 library (data.table)
1\,1
12
13 employment = getSymbols('LNU02300000', src='FRED', auto.assign=F)
14 employment.df = data.frame(date=time(employment), coredata(employment))
15 employment.df[,1] = as.Date(employment.df[,1])
16 colnames(employment.df)[2] = "Employment"
17
18 unemployment = getSymbols('UNRATENSA', src='FRED', auto.assign=F)
   unemployment.df = data.frame(date=time(unemployment), coredata(unemployment) )
19
20
   colnames (unemployment.df) [2] = "Unemployment"
21
22 manufacturing = getSymbols('CEU300000001',src='FRED', auto.assign=F)
23 manufacturing.df = data.frame(date=time(manufacturing), coredata(manufacturing))
24 colnames(manufacturing.df)[2] = "Manufacturing"
25 manufacturing.df = subset(manufacturing.df, date >= "1948-01-01")
26
27
28 # reading-in recession dates - We will just copy-paste the dates
29 # here from Fred's website, for simplicity pourposes
30 recessions.df = read.table(textConnection(
31 "Peak, Trough
32 1857-06-01, 1858-12-01
33 1860-10-01, 1861-06-01
34 1865-04-01, 1867-12-01
35 1869-06-01, 1870-12-01
36 1873-10-01, 1879-03-01
37 1882-03-01, 1885-05-01
38 1887-03-01, 1888-04-01
39 1890-07-01, 1891-05-01
```

```
41 1895-12-01, 1897-06-01
42 1899-06-01, 1900-12-01
43 1902-09-01, 1904-08-01
44 1907-05-01, 1908-06-01
45 1910-01-01, 1912-01-01
46 1913-01-01, 1914-12-01
47 1918-08-01, 1919-03-01
48
   1920-01-01, 1921-07-01
49 1923-05-01, 1924-07-01
50 1926-10-01, 1927-11-01
51 1929-08-01, 1933-03-01
52 1937-05-01, 1938-06-01
53 1945-02-01, 1945-10-01
54 1948-11-01, 1949-10-01
55 1953-07-01, 1954-05-01
56 1957-08-01, 1958-04-01
57 1960-04-01, 1961-02-01
58 1969-12-01, 1970-11-01
59 1973-11-01, 1975-03-01
60 1980-01-01, 1980-07-01
61 1981-07-01, 1982-11-01
62 1990-07-01, 1991-03-01
63
   2001-03-01, 2001-11-01
64 2007-12-01, 2009-06-01"), sep=',',
65 colClasses=c('Date', 'Date'), header=TRUE)
66
67 recessions.trim = subset(recessions.df, Peak >= "1948-01-01")
68
69
70 # Creating initial plots
71 # Employment - Population
72 emp = ggplot(employment.df) + geom_line(aes(x=date, y=Employment)) +
73 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
74 geom_rect (data=recessions.trim, aes(xmin=Peak, xmax=Trough, ymin=-Inf, ymax=+Inf),
75 fill='blue', alpha=0.1) +ylab("Employment-Population Ratio")
76
   emp +scale_x_date(date_breaks = '5 years', date_labels = "%Y %b")
77
78
79 # Unemployment
80 unemp = ggplot(unemployment.df) + geom_line(aes(x=date, y=Unemployment)) +
81 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
82 geom_rect (data=recessions.trim, aes(xmin=Peak, xmax=Trough, ymin=-Inf, ymax=+Inf),
83 fill='blue', alpha=0.1) +ylab("Unemployment Rate")
84
85
  unemp + scale_x_date(date_breaks = '5 years', date_labels = "%Y %b")
86
87 # Manufacturing
88 manuf = ggplot(manufacturing.df) + geom_line(aes(x=date, y=Manufacturing)) +
89 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
90 geom_rect(data=recessions.trim, aes(xmin=Peak, xmax=Trough, ymin=-Inf, ymax=+Inf),
  fill='blue', alpha=0.1) +ylab("Number of Employees in Manufacturing")
91
92
93
   manuf +scale_x_date(date_breaks = '5 years', date_labels = "%Y %b")
94
95
96
97 # Dividing-up data into regimes
```

40 1893-01-01, 1894-06-01

```
98 window1 = data.frame(Beggining = as.Date("1948-01-01"), End = as.Date("1975-03-01"))
99 window2 = data.frame(Beggining = as.Date("1975-04-01"), End = as.Date("2001-11-01"))
100 window3 = data.frame(Beggining = as.Date("2001-12-01"), End = as.Date("2008-09-01"))
101 window4 = data.frame(Beggining = as.Date("2009-7-01"), End = as.Date("2016-06-01"))
102
103
104 emp_piecewise = ggplot(employment.df) + geom_line(aes(x=date, y=Employment)) +
105
   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
106
   geom_rect(data=window1, aes(xmin=Beggining, xmax=End, ymin=-Inf, ymax=+Inf),
107 fill='blue', alpha=0.1) +
108 geom_rect(data=window2, aes(xmin=Beggining, xmax=End, ymin=-Inf, ymax=+Inf),
109 fill='green', alpha=0.1) +
110 geom_rect (data=window3, aes (xmin=Beggining, xmax=End, ymin=-Inf, ymax=+Inf),
111 fill='yellow', alpha=0.1) +
112 geom_rect (data=window4, aes (xmin=Beggining, xmax=End, ymin=-Inf, ymax=+Inf),
113 fill='red', alpha=0.1) +
114 ylab("Employment - Population Ratio")
115
116 emp_piecewise +scale_x_date(date_breaks = '5 years', date_labels = "%Y %b")
117
118
119 # fitting trend lines
120 window1.index = employment.df$date >= window1$Beggining &
121
    employment.df $date <= window1$End</pre>
122
123 window2.index = employment.df$date >= window2$Beggining &
124 employment.df $date <= window2$End
125
126 window3.index = employment.df$date >= window3$Beggining &
127
    employment.df $date <= window3$End</pre>
128
129 window4.index = employment.df$date >= window4$Beggining &
130 employment.df $date <= window4 $End
131
132
133 emp_ts1 = ts (employment.df [window1.index,2], start = c(1948,1),end = c(1975,3),frequency = 12)
134
    emp_ts2 = ts(employment.df[window2.index,2], start = c(1975,4),end = c(2000,11),frequency = 12)
135
136
    emp_ts3 = ts(employment.df[window3.index,2], start = c(2001,12),end = c(2008,09),frequency = 12)
137
138
139 emp_ts4 = ts (employment.df [window4.index,2], start = c(2009,7),end = c(2016,6),frequency = 12)
140
141
142 summary(fit1 <- lm(emp_ts1~time(emp_ts1))) #
143 summary (fit2 <- lm (emp_ts2~time(emp_ts2))) #
144 summary(fit3 <- lm(emp_ts3~time(emp_ts3))) #
145 summary(fit4 <- lm(emp_ts4~time(emp_ts4))) #
146
147
148 # Plotting linear trends
    data.plot <- data.frame(date = time(emp_ts1), Employment_Population = as.numeric(emp_ts1),Trend = fit1$fitted.values)</pre>
149
150
151 data_long <- melt(data.plot, id="date")</pre>
152 plot1 = ggplot(data=data_long,
153 aes(x=date, y=value, colour=variable)) +
154 theme(legend.position="none", axis.text.x = element_text(angle = 45, hjust = 1)) +
155 geom_line() +
```

```
156 labs(title="Blue Area") +ylab("Ratio")
157
158 data.plot <- data.frame(date = time(emp_ts2), Employment_Population = as.numeric(emp_ts2), Trend = fit2$fitted.values)
159
160 data_long <- melt(data.plot, id="date")</pre>
161 plot2 = ggplot(data=data_long,
162 aes(x=date, y=value, colour=variable)) +
163
   theme(legend.position="none",axis.text.x = element_text(angle = 45, hjust = 1))+
164
    geom_line() +
165 labs(title="Green Area (Until Nov. 2000)") +ylab("Ratio")
166
167 data.plot <- data.frame(date = time(emp_ts3), Employment_Population = as.numeric(emp_ts3), Trend = fit3$fitted.values)
168
169 data_long <- melt(data.plot, id="date")
170 plot3 = ggplot(data=data_long,
171
    aes(x=date, y=value, colour=variable)) +
172 theme(legend.position="none", axis.text.x = element_text(angle = 45, hjust = 1))+
173 geom_line() +
174 labs(title="Yellow Area")+ylab("Ratio")
175
176 data.plot <- data.frame(date = time(emp_ts4), Employment_Population = as.numeric(emp_ts4), Trend = fit4$fitted.values)
177
178 data_long <- melt(data.plot, id="date")
179
   plot4 = ggplot(data=data_long,
180 aes(x=date, y=value, colour=variable)) +
181 theme(legend.position="none", axis.text.x = element_text(angle = 45, hjust = 1))+
182 geom_line() +
183 labs(title="Red Area")+ylab("Ratio")
184
185
186 grid.arrange(plot1, plot2,plot3,plot4,nrow = 2, ncol=2)
187
188
189 # From now on we focus on the 4-1975 to 11-2000 period
190 # redefining our variables
191 window = data.frame(Beggining = as.Date("1975-04-01"), End = as.Date("2000-11-01"))
192 emp_ts = emp_ts2
193
194 # ACF and PACF
195 par(mfrow = c(2,2))
196
197 # Maximum number of lags
198 Acf(emp_ts,length(emp_ts),main = "Original Data - All Lags")
199 Pacf(emp_ts,length(emp_ts),main = "Original Data - All Lags")
200
201 # ACF and PACF
202 # Maximum number of lags
203 Acf(emp_ts,50,main = "Original Data - 50 lags")
204 Pacf(emp_ts,50,main = "Original Data - 50 Lags")
205
206 ## The detrended data
207
    detrended = resid(fit2)
208
209 # Plotting detrended data
210 data.plot <- data.frame(date = time(emp_ts), detrended = detrended)
211 ggplot(data.plot) + geom_line(aes(x=date, y=detrended)) +
212 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
213 labs(title="Detrended data - linear trend")
```

```
214
215
216 adf.test(detrended)
217
218 # using smoothing splines to extract trend
219 # we will use a function that we have created
220
\left| 221 \right| # Creating the knots, note that we will use the xi's
222 # themselves as the knots
223 x = as.numeric(time(emp_ts))
224 knots = x
225 n = length(x)
226 Y = as.numeric(emp_ts)
227
228 filloutZ = function(x,knots){
229
230 p = length (knots)
231
232 Z = matrix (nrow = n, ncol = p+1)
233 ZMax = matrix(nrow = n, ncol = p)
234
235 for (i in 1:n){
236
237 for (j in 1:p){
238
239 ZMax[i,j] = max(0,x[i]-knots[j])
240 }
241
242 }
243
244 Z = as.matrix(cbind(1,ZMax))
245
246 return (Z)
247 }
248
249
\left| 250 \right| # Next, we declare a function that solves for alpha.hat
251
    solveforalpha.hat = function(Z,Y,lambda){
252
253 \mod (Z) = \dim (Z) [2]
254
255 D = diag(ncolZ)
256 D[1,1] = 0
257
258 return (solve(t(Z) % * % Z + lambda*D,t(Z) % * % Y))
259
260 }
261
262
263 # Let's fit the model for a fixed level of lambda
264 # We'll manually tweak lambda, until we get a
265 # trend that is smooth enough, and a detrended
266 # series that looks stationary
267 lambda = 100
268
269 Z = filloutZ(x,knots)
270 alpha.hat = solveforalpha.hat(Z,Y,lambda)
271 Y.hat = Z%*%alpha.hat
```

```
272
273 detrended = as.numeric(emp_ts) - as.numeric(Y.hat)
274
275 # Plotting non-linear trend
276 data.plot <- data.frame(date = x, Employment_Population = Y,Trend = Y.hat)
277
278 data_long <- melt(data.plot, id="date")
279
    ggplot (data=data_long,
280
    aes(x=date, y=value, colour=variable)) +
281 theme(legend.position="none",axis.text.x = element_text(angle = 90, hjust = 1)) +
282 geom_line() +
283 labs(title="Data + Non-Linear Trend") +ylab("Ratio")
284
285
286 # Plotting detrended data
287 data.plot <- data.frame(date = x, detrended = detrended)
288 ggplot(data.plot) + geom_line(aes(x=date, y=detrended)) +
289 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
290 labs(title="Detrended data - non-linear trend")
291
292
293 #adf test
294 adf.test(detrended)
295
detrended = ts (detrended, start = c(1975,4), end = c(2000,11), frequency = 12)
297
298 par(mfrow = c(2,2))
299 # ACF and PACF
300 # Maximum number of lags
301 Acf(detrended,length(detrended),main = "Detrended Data - All Lags")
302 Pacf (detrended, length (detrended), main = "Detrended Data - All Lags")
303
304 # ACF and PACF
305 # Maximum number of lags
306 Acf(emp_ts,50,main = "Detrended Data - 50 lags")
307 Pacf(emp_ts,50,main = "Detrended Data - 50 Lags")
308
309
    # Creating Periodograms
310
311 \, \text{par}(\text{mfrow} = c(1,1))
312
314 detrended.per = spec.pgram(detrended, taper=0, log="no",detrend = FALSE,main = "")
315
316 frequencies = detrended.per$freq
317 spectrum = detrended.per$spec
318 rank (-spectrum)
319
320 # The frequencies corresponding to the
321 # largest spectral densities values
322
323 f = NULL
324
325 f[1] = frequencies[rank(-spectrum) == 1]
326 f[2] = frequencies[rank(-spectrum) == 2]
327 f[3] = frequencies[rank(-spectrum) == 3]
328 f[4] = frequencies[rank(-spectrum) == 4]
329 f[5] = frequencies[rank(-spectrum) == 5]
```

```
330 f[6] = frequencies[rank(-spectrum) == 6]
331
332
    # these correspond to the following periods:
333 1/f
334
335 # conf intervals - returned value:
336 | U = qchisq(.025,2) # 0.05063
337 L = qchisq(.975,2) # 7.37775
338 11 = 2*detrended.per$spec[frequencies == f[1]]/L
339 u1 = 2*detrended.per$spec[frequencies == f[1]]/U
340 12 = 2*detrended.per$spec[frequencies == f[3]]/L
341 u2 = 2*detrended.per$spec[frequencies == f[3]]/U
342 13 = 2*detrended.per$spec[frequencies == f[4]]/L
343 u4 = 2*detrended.per$spec[frequencies == f[4]]/U
344 14 = 2*detrended.per$spec[frequencies == f[5]]/L
345 u4 = 2*detrended.per$spec[frequencies == f[5]]/U
346
347
348 abline(v=f[1], lty="dotted")
349 abline (v=f[3], lty="dotted")
350 abline(v=f[4], lty="dotted")
351 abline(v=f[5], lty="dotted")
352
353
    abline(h = l1, lty="dotted", col = 'red')
354
    #abline(h = 12,lty="dotted", col = 'red')
355 #abline(h = 13, lty="dotted", col = 'red')
356 #abline(h = 14, lty="dotted", col = 'red')
357
358 # The spectrums of relevance:
359 spectrum_rel = NULL
360 spectrum_rel[1] = detrended.per$spec[frequencies == f[1]]
361 spectrum_rel[2] = detrended.per$spec[frequencies == f[3]]
362 spectrum_rel[3] = detrended.per$spec[frequencies == f[4]]
363 spectrum_rel[4] = detrended.per$spec[frequencies == f[5]]
364
365 freq = c(f[1],f[3],f[4],f[5])
366
   period = 1/freq
    summary = as.matrix(cbind(freq,period))
367
368
369
371
372 k = kernel("modified.daniell", c(2.2))
373
374 smooth = spec.pgram(detrended, k, taper=0, log="no", main = "Smoothed Periodogram - Modified Daniell(2,2)")
375 frequencies = smooth$freq
376 spectrum = smooth$spec
377 rank (-spectrum)
378
379 # The frequencies corresponding to the
380 # largest spectral densities values
381
382 f = NULL
383
384 f[1] = frequencies[rank(-spectrum) == 1]
385 f[2] = frequencies[rank(-spectrum) == 2]
386 f[3] = frequencies[rank(-spectrum) == 3]
387 f[4] = frequencies[rank(-spectrum) == 4]
```

```
388 f[5] = frequencies[rank(-spectrum) == 5]
389 f[6] = frequencies[rank(-spectrum) == 6]
390 f[7] = frequencies[rank(-spectrum) == 7]
391 f[8] = frequencies[rank(-spectrum) == 8]
392 f[9] = frequencies[rank(-spectrum) == 9]
393 f[10] = frequencies [rank(-spectrum) == 10]
394 f[11] = frequencies [rank(-spectrum) == 11]
395 f[12] = frequencies [rank(-spectrum) == 12]
396
   # these correspond to the following periods:
397
398 1/f
399
    400
401
402
    403
    404
405 \, \text{par}(\text{mfrow} = c(2,2))
406
407 AIC = rep(0, 60) -> BIC
408 for (k in 1:60) {
409 fit = ar(detrended, order=k, aic=FALSE)
410 sigma2 = var(fit$resid, na.rm=TRUE)
411 BIC[k] = \log(sigma2) + (k * \log(n)/n)
412 AIC[k] = \log(sigma2) + ((n+2*k)/n) }
413 IC = cbind(AIC,BIC)
414
415 # Plotting AIC/BIC
416 data.plot <- data.frame(p = seq(1,60,1), AIC = AIC,BIC = BIC)
417
418 data_long <- melt(data.plot, id="p")
419 ggplot (data=data_long,
420 aes(x=p, y=value, colour=variable)) +
421 theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
422 geom_line() +
423 labs(title="") +ylab("AIC/BIC")
424
   lag = which.min(AIC)
425
426
427 par(mfrow=c(1,1))
428
429 sp.ar = spec.ar(detrended, log="no", order=lag, main = "Autoregressive Spectral Estimator")
430
431 frequencies = sp.ar$freq
432 spectrum = sp.ar$spec
433 rank (-spectrum)
434
435 # The frequencies corresponding to the
436 # largest spectral densities
437
438 f = NULL
439
440 f[1] = frequencies[rank(-spectrum) == 1]
441 f[2] = frequencies[rank(-spectrum) == 2]
442 f[3] = frequencies[rank(-spectrum) == 3]
443 f[4] = frequencies[rank(-spectrum) == 4]
444 f[5] = frequencies[rank(-spectrum) == 5]
445 f[6] = frequencies[rank(-spectrum) == 6]
```

```
446 f[7] = frequencies[rank(-spectrum) == 7]
447 f[8] = frequencies[rank(-spectrum) == 8]
448 f[9] = frequencies[rank(-spectrum) == 9]
449
450 # these correspond to the following periods:
451 1/f
452
453 #Removing Seasonal Components
    data = data.table(detrended = as.numeric(detrended),index1 = seq(1,12,1),index2 = seq(1,4,1),index3 = seq(1,64,1),index4 = seq
454
         (1, 6, 1))
455 seasonality1 = data[,mean(detrended),by = index1]
456 seasonality2 = data[,mean(detrended),by = index2]
457 seasonality3 = data[,mean(detrended),by = index3]
458 seasonality4 = data[,mean(detrended),by = index4]
459
460
461
    data_merged = merge(data,seasonality1, all.x = TRUE,by = 'index1', sort = 'FALSE')
462 data_merged = merge(data_merged, seasonality2, all.x = TRUE, by = 'index2', sort = 'FALSE')
463 data_merged = merge(data_merged,seasonality3, all.x = TRUE,by = 'index3', sort = 'FALSE')
464 data_merged = merge(data_merged,seasonality4, all.x = TRUE,by = 'index4', sort = 'FALSE')
465
466 data_merged = as.data.frame(data_merged)
    residual1 = data_merged$detrended - data_merged[,c(6)] # remove 1 year
467
468
    residual2 = data_merged$detrended - rowSums(data_merged[,c(6,8)]) # remove 1 year and 64 months
469 residual3 = data_merged$detrended - rowSums(data_merged[,c(6,7,9)]) # remove 1 year, 4 months, 6 months
470
471 \, \text{par}(\text{mfrow} = c(2,1))
472 Acf(residual1,100, main = 'Cycle removed: 1 year')
473 Acf(residual2,100,main = 'Cycle removed: 1 year and 5.33 years (64 months)')
474
475
    residual = data_merged$detrended - rowSums(data_merged[,c(6,7,8,9)])
476
477 adf.test(residual)
478
479 # plotting one-year cycle
480 data.plot <- data.frame(date = c("Jan","Feb","Mar","Apr","May","Jun",
481 "Jul", "Aug", "Sep", "Oct", "Nov", "Dec"), value = seasonality1$V1)
    ggplot(data.plot) + geom_point(aes(x=factor(date, levels=unique(date)),y=value)) +
482
   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
483
484 ylab("Seasonal Component (%)") +xlab("Month")
485
486
487 #ACF - PACF of residuals
488 | par(mfrow = c(2, 2))
489 acf.data = Acf(residual,length(residual),main = "Residual - All Lags")
490 pacf.data = Pacf(residual,length(residual),main = "Residual - All Lags")
491 acf.data = Acf(residual,50, main = "Residual - 50 Lags")
492 pacf.data = Pacf(residual.50.main = "Residual - 50 Lags")
493
494 # PACF seems to cut off after 14, 15 lags,
495 # ACF tails off
    # suggests AR(14) or AR(14)
496
497
498
500 #figuring out best AIC for AR model
501
502 AIC = rep(0, 40) -> BIC
```

```
503 for (k in 1:40) {
504
505 fit = ar(residual, order=k, aic=FALSE)
506 sigma2 = var(fit$resid, na.rm=TRUE)
507 BIC[k] = log(sigma2) + (k*log(n)/n)
508 AIC[k] = log(sigma2) + ((n+2*k)/n) }
509
510
511 IC = cbind(AIC,BIC)
512
513 lag = which.min(AIC)
514
515 # Plotting AIC/BIC
516 data.plot <- data.frame(p = seq(1,40,1), AIC = AIC,BIC = BIC)
517
518 data_long <- melt(data.plot, id="p")
519 ggplot (data=data_long,
520 aes(x=p, y=value, colour=variable)) +
521 theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
522 geom_line() +
523 labs(title="") +ylab("AIC/BIC")
524
525
526
    #Fitting AR model
527 model_ar = arima(residual, order = c(14, 0, 0),include.mean = FALSE,method="ML")
528
529 ts = model_ar$residuals
530
531 par(mfrow = c(2,2))
532 qqnorm(ts)
533 qqline(ts, col=2)
534 hist(ts, br=12, main = 'Histogram of Residuals')
535 shapiro.test(ts)
536 Acf(ts,lag = 50, main = 'ACF of model residuals')
537
538 ljung.box = NULL
539 for (i in 1:50){
540 ljung.box[i] = Box.test(ts,lag = i, type = c("Ljung-Box"))$p.value
541 }
542
543 # we have serial correlation from lags 12 to 16
544 plot (ljung.box)
545 abline (h = 0.05, lty="dotted", col = 'red')
546
547 # fitting an arma model to data
548 k = 9
549 model_arima = arima(residual, order = c(14, 0, k), include.mean = FALSE, method = 'ML', optim.control = list(maxit = 1000))
550 sigma2 = var(model_arima$residuals, na.rm=TRUE)
551 AIC = \log(sigma2) + ((n+2*k)/n)
552
553 ts = model_arima$residuals
554
555
556 | par(mfrow = c(2,2))
557 qqnorm(ts)
558 qqline(ts, col=2)
559 hist(ts, br=12, main = 'Histogram of Residuals')
560 shapiro.test(ts)
```

```
561 Acf(ts,lag = 50, main = 'ACF of model residuals')
562
563 ljung.box = NULL
564 for (i in 1:50) {
565 | ljung.box[i] = Box.test(ts, lag = i, type = c("Ljung-Box"))$p.value
566 }
567
568 plot (ljung.box,ylim = c(0,1))
    abline(h = 0.05, lty="dotted", col = 'red')
569
570
571
573 # AR(14) forecasts
574 regr = ar.ols(residual, order=14, demean=FALSE, intercept=FALSE)
575 fore = predict (regr, n.ahead=12)
576
577
    # Seaonality forecast
578
579 s4 = c (3,4,5,6,1,2,3,4,5,6,1,2)
580 \ s3 = seq(53, 64)
581 s2 = rep(c(1,2,3,4),3)
582 s1 = c(9,10,11,12, seq(1,8,1))
583
584
    data = data.table(index1 = s1, index2 = s2, index3 = s3, index4 = s4)
585
586 data_merged = merge(data,seasonality1, all.x = TRUE,by = 'index1', sort = 'FALSE')
587 data_merged = merge(data_merged,seasonality2, all.x = TRUE,by = 'index2', sort = 'FALSE')
588 data_merged = merge(data_merged,seasonality3, all.x = TRUE,by = 'index3', sort = 'FALSE')
589 data_merged = merge(data_merged,seasonality4, all.x = TRUE,by = 'index4', sort = 'FALSE')
590
591
    data_merged = as.data.frame(data_merged)
592 seasonal_forecast = rowSums(data_merged[,5:8])
593
594 # Trend forecast
595
596 filloutZ_test = function(x,knots){
597
598 p = length (knots)
599 n = length(x)
600 Z = matrix (nrow = n, ncol = p+1)
601 ZMax = matrix (nrow = n, ncol = p)
602
603 for (i in 1:n){
604
605 for (j in 1:p){
606
607 ZMax[i,j] = max(0,x[i]-knots[j])
608 }
609
610 }
611
612 Z = as.matrix(cbind(1,ZMax))
613
614 return (Z)
615 }
616
617
618
```

```
619
620 time_test = x[length(x)]+seq(1,12,1)/12
621 Z_test = filloutZ_test(time_test,knots)
622 Y.hat = Z_test % *% alpha.hat
623
624 # Adding up all components
625 forecast = as.numeric(Y.hat) + as.numeric(seasonal_forecast)+as.numeric(fore$pred)
626
    fore$pred = ts(forecast, start = c(2000,12), end = c(2001,11), frequency = 12)
627
   forese = ts(as.numeric(fore se), start = c(2000, 12), end = c(2001, 11), frequency = 12)
628
629
630 window.plot = data.frame(Beggining = as.Date("1999-01-01"), End = as.Date("2001-11-01"))
631 window.index = employment.df$date >= window.plot$Beggining &
632 employment.df $date <= window.plot $End
633
    emp.plot = ts(employment.df[window.index,2], start = c(1995,1),end = c(2001,11),frequency = 12)
634
635 par(mfrow = c(1,1))
636 dataplot = window(emp.plot, c(1995,1), c(2001, 11))
637 ts.plot(dataplot, fore$pred, col=1:2,
638 ylab="Employment/Population", ylim = c(62,66))
639 lines(fore$pred, type="p", col=2)
640 lines(fore$pred+2*fore$se, lty="dashed", col=4)
    lines(fore$pred-2*fore$se, lty="dashed", col=4)
641
642
643
644
    645
646 fore = predict (model arima.n.ahead = 12)
647
648 # Adding up all components
649 forecast = as.numeric(Y.hat) + as.numeric(seasonal_forecast)+as.numeric(fore$pred)
650 fore$pred = ts(forecast, start = c(2000,12), end = c(2001,11), frequency = 12)
651 fore$se = ts(as.numeric(fore$se), start = c(2000,12), end = c(2001,11), frequency = 12)
652
653
654 window.plot = data.frame(Beggining = as.Date("1999-01-01"), End = as.Date("2001-11-01"))
   window.index = employment.df$date >= window.plot$Beggining &
655
    employment.df $date <= window.plot $End</pre>
656
    emp.plot = ts(employment.df[window.index,2], start = c(1995,1),end = c(2001,11),frequency = 12)
657
658
659 par(mfrow = c(1,1))
660 dataplot = window(emp.plot, c(1995,1), c(2001, 11))
661 ts.plot(dataplot, fore$pred, col=1:2,
662 ylab="Employment/Population", ylim = c(62,66))
663 lines(fore$pred, type="p", col=2)
664 lines(fore$pred+2*fore$se, lty="dashed", col=4)
665 lines(fore$pred-2*fore$se, lty="dashed", col=4)
666
667
668
669
    670
    671
   trend.season.forecast = as.numeric(Y.hat) + as.numeric(seasonal_forecast)
672
673
674 window.plot = data.frame(Beggining = as.Date("2000-12-01"), End = as.Date("2001-11-01"))
675 window.index = employment.df$date >= window.plot$Beggining &
676 employment.df $date <= window.plot $End
```

```
677
678
test = ts(employment.df[window.index,2]-trend.season.forecast,
679
start = c(2000,12),end = c(2001,11),frequency = 12)
680
681
# One-step-forecasts
682
fit <- Arima(residual,order = c(14,0,0),include.mean = FALSE)
683
fit_test <- Arima(c(residual,test), model=fit)
684
onestep <- fitted(fit_test)[309:320] + trend.season.forecast
685
onestep = ts(onestep, start = c(2000,12), end = c(2001,11), frequency = 12)
686
687
# plot
688
ts.plot(dataplot, onestep, col=1:2,
689
ylab="Employment/Population",ylim = c(62,66))</pre>
```