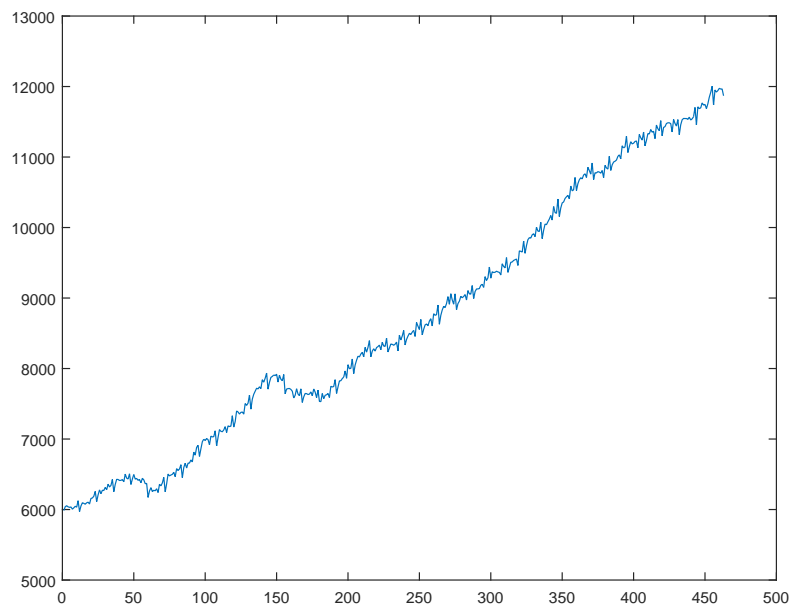


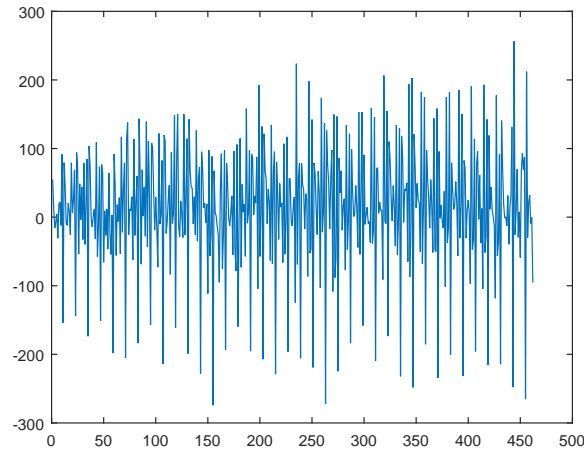
FORECASTING PROJECT: Australian labour force

Introduction

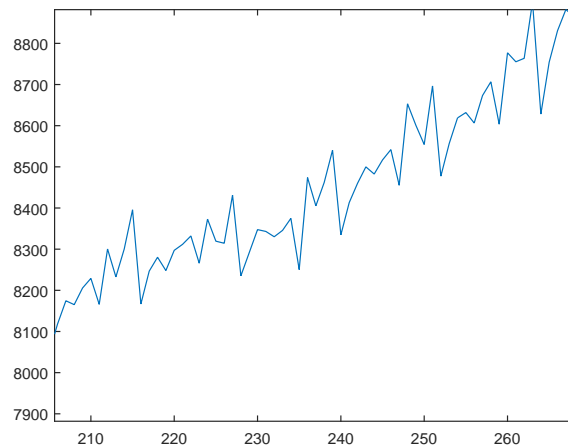
The variable I am forecasting is the Australian labour force which is the number of individuals in Australia who are able to work. The forecasts include forecast for all the months from September 2016 to August 2017 and density forecasts for the same months. The data for the month of September 2016 will be available on the website of the Australian Bureau of Statistics on October 20th. The graph below represents the evolution of the Australian labour force from February 1978 to August 2016, the y-axis represents the labour force (in thousands) and the x-axis represents the months.



There is a clear upward trend in the data, continuously increasing until August 2016, going from around 6 million at the beginning of the period to around 12 million at the end of our period. Looking at the graph, it is clear that the data is non-stationary and this will cause problems in our forecasting models: the statistical properties of our time series such as mean, variance and autocorrelation are not constant over time. This will be a problem particularly in our AR models, we then have to take the first difference (period-period change) which gives us:



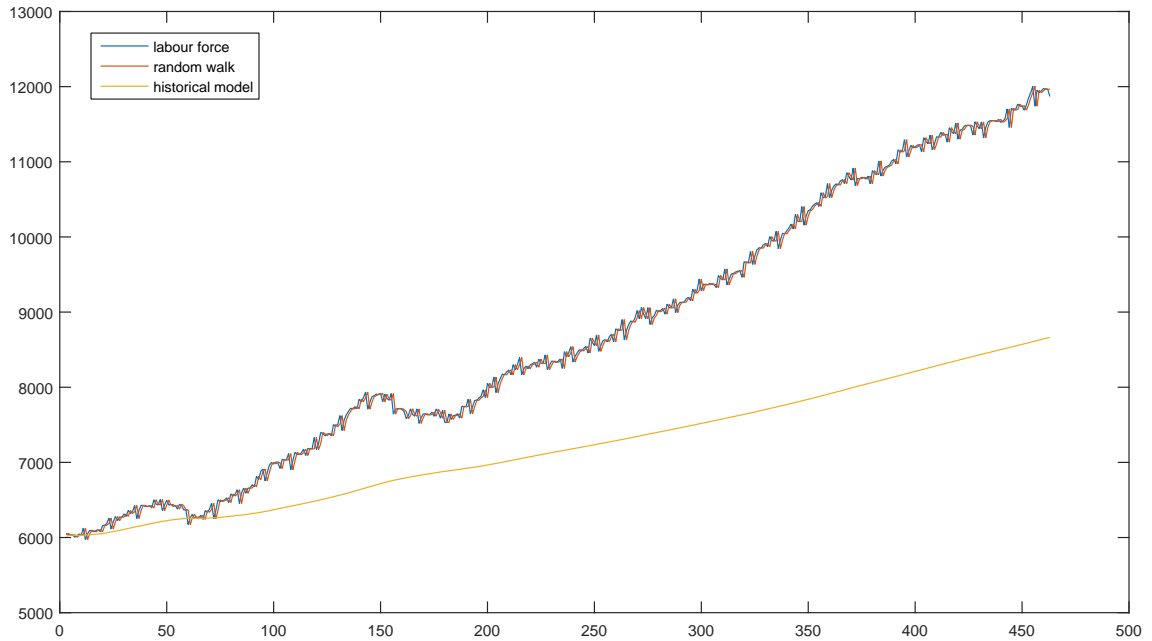
Now, we can use AR models on the first difference. Furthermore, our original data does not appear to have any monthly or seasonal effect, this is due to the scale of the graph. If we zoom in, we can indeed observe some monthly or quarterly effects:



It is clear from this graph that there are indeed seasonal effects, they appear to be quarterly effects. We will then include a model with monthly and quarterly dummies. Two predictor will also be taken into account in our model, namely the Australian GDP and the Australian population during the same period.

Benchmark model

The benchmark model will be the random walk model. The mean squared forecasting error for the random walk model is $8.9491e + 03$ whereas for the historical model, the MSFE is $3.2598e + 06$. The random walk model is a good benchmark to forecast our time series.



In yellow is the historical model, it is really far from the actual observations. The random walk model seems to forecast much better the labour force which is expected, the random walk model follows the actual observations with a lag of one period. Therefore, we will use the MSFE of the random walk model ($8.9491e + 03$) as a benchmark for our relative error measures. The code can be found in the appendix under code 1.

Seasonal effect

The excel file has been modified by adding a column with the monthly dummies and column with the quarterly dummies. Using a model with monthly and quarterly dummies:

$$y_t = \alpha t + \sum_{i=1}^4 \beta_i Q_{it} + \epsilon_t$$

and

$$y_t = \alpha t + \sum_{i=1}^{12} \theta_i M_{it} + \epsilon_t$$

where Q_{it} and M_{it} are respectively the quarterly dummies and the monthly dummies. Here, we assume that $T_0=50$, i.e. we use the first 50 months to forecast the remaining

months. For the model including an effect for each quarter, we find a MSFE of 6.9203e+04 while for the monthly model, the MSFE is 7.1473e+04. There is more a quarterly effect in our time series than a monthly effect but both models have a higher error than the error in our benchmark model. We will then use the quarterly model and try to improve it to beat our benchmark. The code can be found in the appendix under code 2.

Predictors

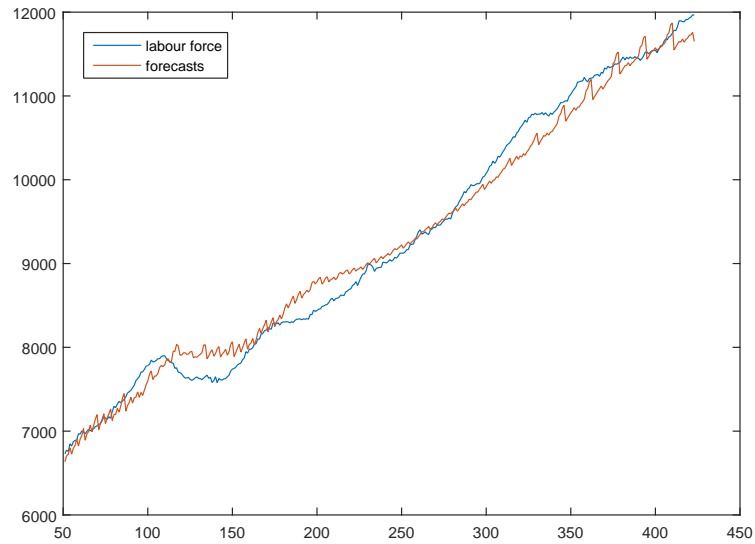
It seems relevant to include the increase of the GDP and the increase of the population during the period studied in our model. With the quarterly model, we will now introduce a predictor which is the Australian GDP for the same period. There is only quarterly data for the GDP so four consecutive months in our labour force time series will correspond to the same GDP for this quarter. The model now is:

$$y_t = \alpha t + \sum_{i=1}^4 \beta_i Q_{it} + \gamma x_{t-h} \epsilon_t$$

where x_{t-h} is the GDP at time $t - h$. (code 3 in the appendix). The MSFE becomes 4.5383e+04, lower than the previous 6.9203e+04 but the relative mean squared forecasting error is 5.0713 which means that the forecasting error is five times that of the random walk model. Now, we add another predictor (the Australian population) to improve the model.

$$y_t = \alpha t + \sum_{i=1}^4 \beta_i Q_{it} + \gamma x_{t-h} + \delta z_{t-h} + \epsilon_t$$

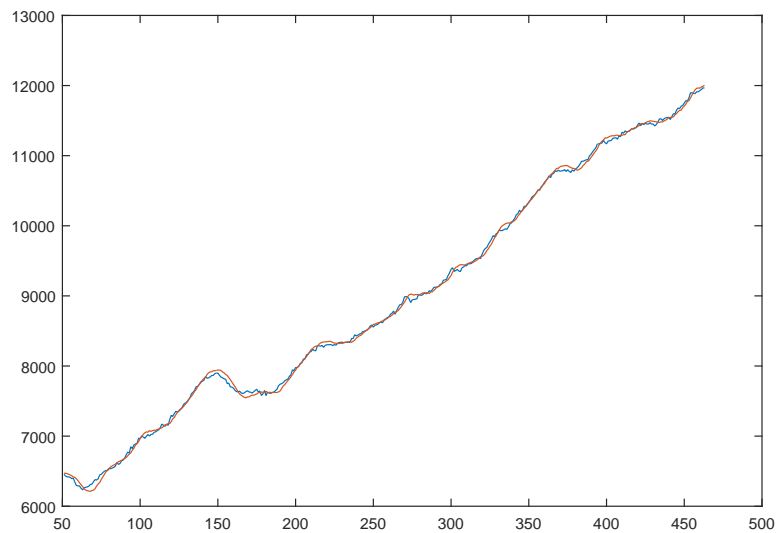
where z_{t-h} is the Australian population at time $t - h$. However, there is only data on the Australian population from June 1981 so we are cutting the first 40 months of our labour force time series and our gdp time series to get the same vector length as the population time series. Now, the MSFE becomes 3.6841e+04, it still gets lower but the relMSFE is still superior than 1 (=4.1167): our model still performs worse than a simple random walk.



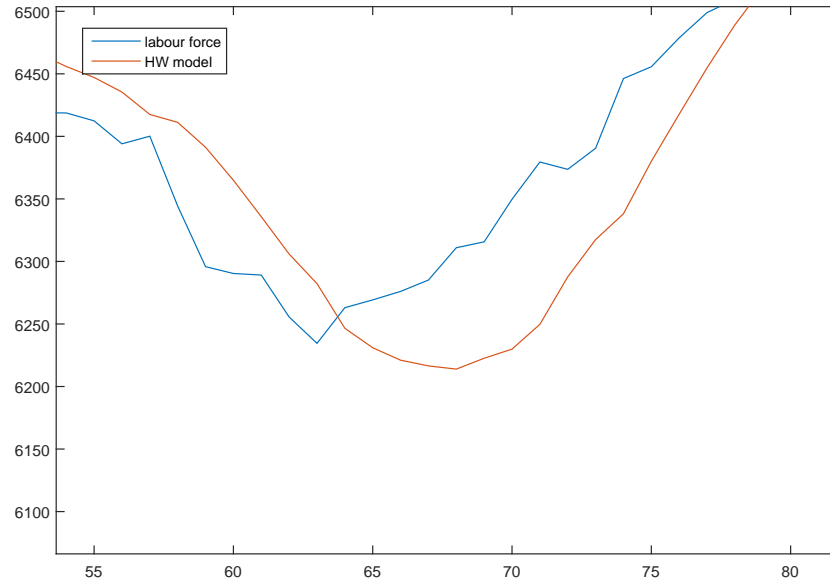
(code 4 in the appendix)

Holt-Winters Smoothing

We now introduce a more sophisticated model. With the Holt-Winters Smoothing, we find a MSFE of 1.9810×10^3 which is lower than our benchmark's, the relMSFE is 0.2214, the error in this model is 0.2214 times the error in the random walk model. We have improved the precision of our model.



From the graph, we can see that the Holt-Winters model (in orange) forecasts well the time series (in blue) but if we zoom in, we still have some significant differences:



(code 5 from the appendix)

AR model

From the data of the labour force, we can see that the values are closely related to past value, there is not a large difference between the present value and close past values. Past values might have an effect on current values. It might then be relevant to use an autoregressive model (AR) to capture this effect. The form of an AR(p) model with p lags is:

$$y_t = c + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t$$

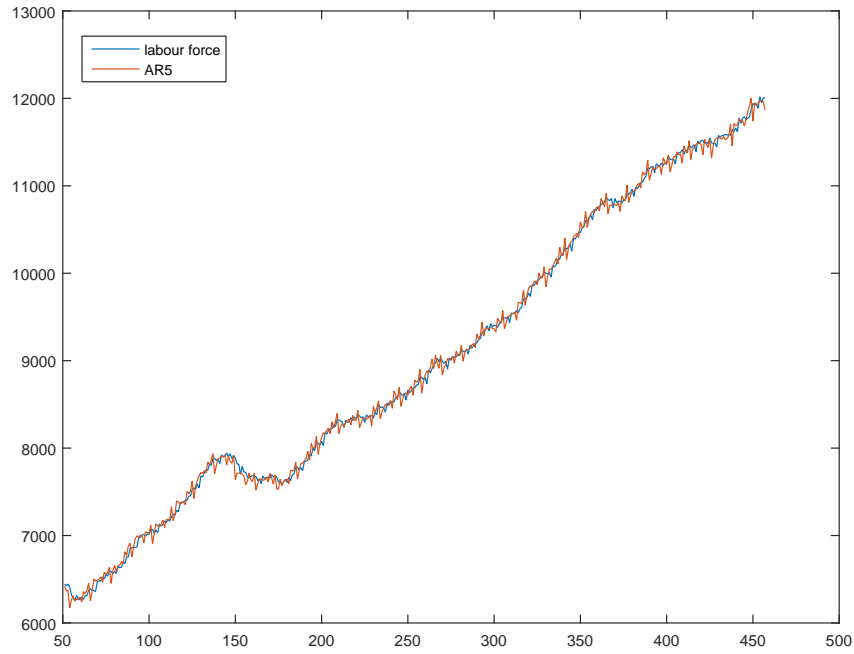
where ϵ_t is an uncorrelated innovation process with mean zero.

We now have to determine the best AR model in terms of lags.

Table 1: AR process

	AR1	AR2	AR3	AR4	AR5	AR6
MSFE	9.4032e+03	6.6201e+03	5.4193e+03	5.4143e+03	5.3204e+03	5.3380e+03
relMSFE	1.0507	0.7398	0.6056	0.6050	0.5945	0.5965

We see that after one lag, we get a model that beats our benchmark model, the relMSFE is lower than 1. The best model is the AR5 process which means that after 5 lags, the past starts losing its predictive power: the error for the AR6 process is higher than for the AR5 process. (code 6 in the appendix).



This model still performs worse than the Holt-Winters model for one-step ahead forecasts. Now, we have to remember that our time series is not stationary so it is more relevant to study the first difference of our time series than the actual observations.

Table 2: AR process on first diff

	AR1_firstdiff	AR2_firstdiff	AR3_firstdiff	AR4_firstdiff
MSFE	6.6110e+03	6.2173e+03	5.9776e+03	6.0108e+03
relMSFE	0.7387	0.6947	0.6680	0.6717

Now, the best model is the AR3 model (code 7 in the appendix). For one lag, the AR model already performs better than the benchmark. Although in this case, the error is larger than for the best AR model with our original time series.

Forecasts

We now have two competing models which are the Holt-winters model and the AR3 model on the first difference of our time series. We have to remember that the Holt-Winters model loses its precision as the number of steps ahead forecasted, h , increases. The Holt-Winters model might not be very accurate in forecasting the labour force in the next twelve months, we then have to compare the errors of the two models as the number of steps ahead forecasted increases.

Table 3: Comparison between AR3 and HW

relMSFE	AR3_firstdiff	Holt-Winters
h=1	0.6680	0.2214
h=2	0.9152	0.3197
h=3	0.9275	0.4466
h=4	0.9646	0.6043
h=5	0.8756	0.7932
h=6	0.9423	1.0168
h=7	0.9057	1.2743
h=8	0.9038	1.5669
h=9	0.6836	1.8924
h=10	0.2118	2.2581
h=11	0.2108	2.6577
h=12	0.2137	3.0947

We observe that the Holt-Winters model performs better than the AR3 process until h reaches 6. For $h = 6$, the Holt-Winters model performs worse than the benchmark (i.e. $relMSFE > 1$). Therefore, to forecast the Australian labour force, we will use the Holt-Winters model to forecast the first 5 months from September 2016 to January 2017 and then use the AR3 process to forecast the remaining months until August 2018. The point forecasts are then:

Table 4: Point forecasts for the coming year

	Point forecasts
Sep 2016	1.2011e+04 (HW model)
Oct 2016	1.2018e+04 (HW model)
Nov 2016	1.2047e+04 (HW model)
Dec 2016	1.2057e+04 (HW model)
Jan 2017	1.2071e+04 (HW model)
Feb 2017	-2.0295 (AR3 on first diff)
March 2017	-15.0879 (AR3 on first diff)
April 2017	16.4162 (AR3 on first diff)
May 2017	-15.4250 (AR3 on first diff)
June 2017	-3.1912 (AR3 on first diff)
July 2017	5.1112 (AR3 on first diff)
August 2017	-83.0737 (AR3 on first diff)

which gives us the following forecasts:

Table 5: Points forecasts for the following 12 months starting in Sept 2016

	Point forecasts (in thousands)
Sep 2016	12011.189
Oct 2016	12018.175
Nov 2016	12047.355
Dec 2016	12057.371
Jan 2017	12070.916
Feb 2017	12068.887
March 2017	12053.799
April 2017	12070.215
May 2017	12054.790
June 2017	12051.599
July 2017	12056.710
August 2017	11973.636

The density forecasts are (code 8 in the appendix):

Table 6: Density forecasts for the following 12 months starting in September 2016

	90% confidence interval	
Sep 2016	-61.3510	177.6884
Oct 2016	-155.3112	124.1457
Nov 2016	-164.1685	116.8099
Dec 2016	-116.5341	169.6654
Jan 2017	-184.1954	88.1459
Feb 2017	-142.6665	138.6076
March 2017	-152.7846	122.6087
April 2017	-121.0130	153.8454
May 2017	-134.8624	104.0125
June 2017	-68.9598	62.5774
July 2017	-60.4403	70.6626
August 2017	-148.6281	-17.5193

For the density forecasts, the AR3 process on the first difference of the time series is used so the values represent the difference from the previous months' level: it gives us an interval of variation of the level of labour force from one month to the next.

Appendix

a) Code 1: random walk

```

1 data=xlsread('labor','Data1');
2 data=data(6:end,3);
3 T=length(data);
4
5 %one-step ahead forecasts historical mean and random walk
6 h=1; T0=2;
7 datahat_h=zeros(T-h-T0+1,1);
8 datahat_r=zeros(T-h-T0+1,1);
9 data_obs=data(T0+h:end);
10 for t=T0:T-h
11     datat=data(1:t);
12     datahat_r(t-T0+h,:)=datat(end);
13     datahat_h(t-T0+h,:)=mean(datat);
14 end
15 MSFERW=mean((data_obs-datahat_r).^2)
16 MSFE_h=mean((data_obs-datahat_h).^2)
17
18 plot(1:T,data);
19 print -depsc myfig.eps
20
21 plot(T0+h:T,data_obs,T0+h:T,datahat_r,T0+h:T,datahat_h);

```

b) Code 2: seasonal effect

```

1 data=xlsread('data');
2 y=data(:,1);
3 T=length(y);
4 t=(1:T)';
5
6 Q=data(:,4);
7 R=data(:,2);
8 T0=50;
9
10 %% construct 4 dummy variables
11 D1=(Q==1); D2=(Q==2);
12 D3=(Q==3); D4=(Q==4);
13 %% construct 12 dummy variables
14 B1=(R==1); B2=(R==2); B3=(R==3);
15 B4=(R==4); B5=(R==5); B6=(R==6);
16 B7=(R==7); B8=(R==8); B9=(R==9);

```

```

17 B10=(R==10); B11=(R==11); B12=(R==12);
18
19 h=1;
20 syhat=zeros(T-h-T0+1, 1);
21 ytp=y(T0+h:end);%observed y_{t+h}
22 for t=T0:T-h
23     yt=y(1:t);
24     D1t=D1(1:t); D2t=D2(1:t);
25     D3t=D3(1:t); D4t=D4(1:t);
26     Xt1=[(1:t)' D1t D2t D3t D4t];
27     beta=(Xt1'*Xt1)\(Xt1'*yt);
28     yhat=[t+h D1(t+h) D2(t+h) D3(t+h) D4(t+h)]*beta;
29     syhat(t-T0+1)=yhat;
30 end
31 MSFE=mean((ytp-syhat).^2)
32
33 h=1;
34 syhat1=zeros(T-h-T0+1, 1);
35 ytp1=y(T0+h:end);%observed y_{t+h}
36 for t=T0:T-h
37     y1t=y(1:t);
38     B1t=B1(1:t); B2t=B2(1:t); B3t=B3(1:t);
39     B4t=B4(1:t); B5t=B5(1:t); B6t=B6(1:t);
40     B7t=B7(1:t); B8t=B8(1:t); B9t=B9(1:t);
41     B10t=B10(1:t); B11t=B11(1:t); B12t=B12(1:t);
42     Xt2=[(1:t)' B1t B2t B3t B4t B5t B6t B7t B8t B9t B10t B11t
           B12t];
43     beta1=(Xt2'*Xt2)\(Xt2'*y1t);
44     yhat1=[t+h B1(t+h) B2(t+h) B3(t+h) B4(t+h) B5(t+h) B6(t+h)
            B7(t+h) B8(t+h) B9(t+h) B10(t+h) B11(t+h) B12(t+h)]*
            beta1;
45     syhat1(t-T0+1)=yhat1;
46 end
47 MSFE1=mean((ytp-syhat1).^2)

```

c) Code 3: GDP predictor

```

1 data=xlsread('data');
2 y=data(:,1);
3 g=data(:,3);
4 T=length(y);
5 t=(1:T)';
6

```

```

7 Q=data (:,4) ;
8 R=data (:,2) ;
9 T0=50;
10
11 %% construct 4 dummy vairables
12 D1=(Q==1); D2=(Q==2);
13 D3=(Q==3); D4=(Q==4);
14 %% construct 12 dummy variables
15 B1=(R==1); B2=(R==2); B3=(R==3);
16 B4=(R==4); B5=(R==5); B6=(R==6);
17 B7=(R==7); B8=(R==8); B9=(R==9);
18 B10=(R==10); B11=(R==11); B12=(R==12);
19
20 h=1;
21 syhat=zeros (T-h-T0+1, 1) ;
22 ytp=y (T0+h:end) ;%observed y_{t+h}
23 for t=T0:T-h
24     yt=y (1:t) ;
25     yt=yt (h+1:t , : ) ;
26     D1t=D1 (1:t) ; D2t=D2 (1:t) ;
27     D3t=D3 (1:t) ; D4t=D4 (1:t) ;
28     gdp=lagmatrix (g (1:t) ,h) ;
29     Xt=[lagmatrix ((1:t) ,h) D1t D2t D3t D4t gdp] ;
30     Xt=Xt (h+1:t , : ) ;
31     beta=(Xt' *Xt) \ (Xt' *yt) ;
32     yhat=[t+h D1(t+h) D2(t+h) D3(t+h) D4(t+h) gdp(t)] *beta ;
33     syhat (t-T0+1)=yhat ;
34 end
35 MSFE1=mean ((ytp-syhat) .^2)

```

d) Code 4: GDP and population predictors

```

1 data=xlsread ( 'data' ) ;
2 y=data (41:end,1) ;
3 p=data (41:end,5) ;
4 g=data (41:end,3) ;
5 T=length (y) ;
6 t=(1:T) ' ;
7
8 Q=data (:,4) ;
9 R=data (:,2) ;
10 T0=50;
11

```

```

12 %% construct 4 dummy variables
13 D1=(Q==1); D2=(Q==2);
14 D3=(Q==3); D4=(Q==4);
15 %% construct 12 dummy variables
16 B1=(R==1); B2=(R==2); B3=(R==3);
17 B4=(R==4); B5=(R==5); B6=(R==6);
18 B7=(R==7); B8=(R==8); B9=(R==9);
19 B10=(R==10); B11=(R==11); B12=(R==12);
20
21 h=1;
22 syhat=zeros(T-h-T0+1, 1);
23 ytp=y(T0+h:end);%observed y_{t+h}
24 for t=T0:T-h
25     yt=y(1:t);
26     yt=yt(h+1:t, :);
27     D1t=D1(1:t); D2t=D2(1:t);
28     D3t=D3(1:t); D4t=D4(1:t);
29     gdp=lagmatrix(g(1:t),h);
30     p1=lagmatrix(p(1:t),h);
31     Xt=[lagmatrix((1:t),h) D1t D2t D3t D4t gdp p1];
32     Xt=Xt(h+1:t, :);
33     beta=(Xt'*Xt)\(Xt'*yt);
34     yhat=[t+h D1(t+h) D2(t+h) D3(t+h) D4(t+h) gdp(t) p1(t)]*
           beta;
35     syhat(t-T0+1)=yhat;
36 end
37
38 MSFE1=mean((ytp-syhat).^2)
39 MSFERW= 8.949053318210024e+03;
40 RelMSE=MSFE1/MSFERW
41
42 plot(T0+h:T,ytp,T0+h:T,syhat)

```

e) Code 5: Holt-Winters smoothing

```

1 data=xlsread('data');
2 y=data(:,1);
3
4 T=length(data);
5 t=(1:T)';
6
7 T0 = 50; h = 1; s=12;
8 syhat = zeros(T-h-T0+1,1);

```

```

9  ytph = y(T0+h:end); % observed y {t+h}
10 alpha = .2; beta = .2; gamma=.2;% smoothing parameters
11 St=zeros(T-h,1);
12 Lt = mean(y(1:s)); bt = 0; St(1:12)=y(1:s)/Lt
13 for t = s+1:T-h
14 newLt = alpha*(y(t)-St(t-s)) + (1-alpha)*(Lt+bt);
15 newbt = beta*(newLt-Lt) + (1-beta)*bt;
16 St(t)=gamma*(y(t)-newLt) + (1-gamma)*St(t-s);
17 yhat = newLt + h*newbt + St(t+h-s);
18 Lt = newLt; bt = newbt; % update Lt and bt
19 if t>= T0 % store the forecasts for t >= T0
20 syhat(t-T0+1,:) = yhat;
21 end
22 end
23
24 MSFE1 = mean((ytph-syhat).^2)
25 MSFERW= 8.949053318210024e+03;
26 RelMSE=MSFE1/MSFERW
27
28 plot(T0+h:T,ytph,T0+h:T,syhat)

```

f) Code 6: AR5 process

```

1  data=xlsread('labor','Data1');
2  data=data(6:end,3);
3
4  m=6;
5  data0=data(1:m);
6  data=data(m+1:end);
7  T=length(data);
8  T0=50; h=1;
9  data_AR5=zeros(T-h-T0+1,1);
10
11 for t=T0:T-h
12     datat=data(h:t);
13     yt=[[data0(m);data(1:t-h)] [data0(m-1:end);data(1:t-h-1)]
14         [data0(m-2:end);data(1:t-h-2)]...
15         [data0(m-3:end);data(1:t-h-3)] [data0(m-4:end);data(1:
16             t-h-4)]];
17     betahat5 =(yt'*yt)\(yt'*datat);
18     data_AR5(t-T0+1,:)= [data(t) data(t-1) data(t-2) data(t-3)
19         data(t-4)]*betahat5;
20 end

```

```

18 data_obs=data (T0+h:end) ;
19
20 MSFE_AR5=mean(( data_obs-data_AR5).^2)
21 MSFERW= 8.949053318210024e+03;
22 RelMSE=MSFE_AR5/MSFERW
23
24 plot (T0+h:T, data_AR5 ,T0+h:T, data_obs) ;

```

g) Code 7: AR3 process on first diff

```

1 clear all
2 clc
3
4 data=xlsread('labor','Data1');
5 data=data(6:end,3);
6 T=length(data);
7
8 Newdata=data(2:end)-data(1:end-1);
9
10 m=4;
11 Newdata0=Newdata(1:m);
12 Newdata= Newdata(m+1:end);
13 T=length(Newdata);
14 T0=50;h=1;
15 Newdata_AR=zeros(T-h-T0+1,1);
16 for t=T0:T-h
17     Newdatat=Newdata(h:t);
18     yt=[[Newdata0(m);Newdata(1:t-h)] [Newdata0(m-1:end);
19         Newdata(1:t-h-1)]...
20         [Newdata0(m-2:end);Newdata(1:t-h-2)]];
21     betahat1=(yt'*yt)\(yt'*Newdatat);
22     Newdata_AR(t-T0+1,:)= [Newdata(t) Newdata(t-1) Newdata(t-2)
23         ]*betahat1;
24 end
25
26 dataobs=Newdata(T0+h:end);
27 MSFE_AR=mean(( dataobs-Newdata_AR).^2);
28 MSFERW= 8.949053318210024e+03;
29 RelMSE=MSFE_AR/MSFERW

```

h) Code 8: Density forecasts

```

1 sigmahat=(1/(length(data)-1))*sum((dataobs-Newdata_AR).^2); %
    where dataobs are the actual observations, Newdata_AR are

```



```
the forecasts using an AR3 process on the first difference
2
3 IntForecast1=[Newdata_AR1-1.645*sqrt(sigmahat) Newdata_AR1
+1.645*sqrt(sigmahat) ] %where Newdata_AR1 is the point
forecast for a particular month
```