

Stats 326 A3

```
Richard Choi
2022-05-15

library(fpp3)

## Warning: package 'fpp3' was built under R version 4.1.3

## -- Attaching packages ----- fpp3 8.4.8 --

## v tidbits 3.1.6 v tidbits 3.1.1
## v dplyr 1.2.8 v tidbits 8.4.8
## v lubridate 1.2.6 v feasts 0.2.2
## v lubridate 1.2.8 v feasts 0.2.1
## v ggplot2 3.3.5

## Warning: package 'tidbits' was built under R version 4.1.3

## Warning: package 'dplyr' was built under R version 4.1.3

## Warning: package 'tidyr' was built under R version 4.1.3

## Warning: package 'lubridate' was built under R version 4.1.3

## Warning: package 'ggplot2' was built under R version 4.1.3

## Warning: package 'tidbitsdata' was built under R version 4.1.3

## Warning: package 'feasts' was built under R version 4.1.3

## Warning: package 'fabletools' was built under R version 4.1.3

## Warning: package 'fable' was built under R version 4.1.3

## -- Conflicts ----- fpp3_conflicts --
## x lubridate::date() masks base::date()
## x dplyr::filter() masks stats::filter()
## x tidbits::intersect() masks base::intersect()
## x tidbits::interval() masks lubridate::interval()
## x dplyr::lag() masks stats::lag()
## x tidbits::setdiff() masks base::setdiff()
## x tidbits::union() masks base::union()

library(tidyverse)

## Warning: package 'tidyverse' was built under R version 4.1.3

## -- Attaching packages ----- tidyverse 1.3.1 --

## v purrr 0.3.4 v strings 1.4.8
## v purrr 0.3.4 v forecasts 0.5.1

## Warning: package 'rread' was built under R version 4.1.3

## Warning: package 'purrr' was built under R version 4.1.3

## Warning: package 'strings' was built under R version 4.1.3

## Warning: package 'forecast' was built under R version 4.1.3

## -- Conflicts ----- tidyverse_conflicts() --
## x lubridate::as_datetime() masks base::as_datetime()
## x lubridate::date() masks base::date()
## x dplyr::filter() masks stats::filter()
## x tidbits::intersect() masks lubridate::intersect()
## x tidbits::interval() masks lubridate::interval()
## x dplyr::lag() masks stats::lag()
## x tidbits::setdiff() masks lubridate::setdiff()
## x tidbits::union() masks lubridate::union()

library(lubridate)
library(ggally)

## Warning: package 'Ggally' was built under R version 4.1.3

## Registered 53 method overwritten by 'Ggally':
## method from
## + gg ggplot2
```

Question 1

Plot the data and comment on what you can observe

```
productivity.df = read_csv("productivity.csv")

## Rows: 44 Columns: 2
## -- Column specification -----
## delimiter: ","
## db1 (2): Year, Productivity
##
## I use 'spec()' to retrieve the full column specification for this data.
## I specify the column types or set 'show_col_types = FALSE' to quiet this message.

productivity.tibble = productivity.df %>%
  as_tibble(index="year")

productivity.tibble %>% autoplot() = kable("year") +
  ggtitle("The labour productivity index for primary industries in New Zealand 1978 - 2021")

## Plot variable not specified, automatically selected `vars = Productivity`

The labour productivity index for primary industries in New Zealand 1978 - 2021

## We can see that there is a strong
```

positive linear trend in the plot. This means that the New Zealand productivity in primary industries is being more efficient as the time progresses.

Question 2

Create a training set that contains data from the years 1978–2016.

```
productivityTraining = productivity.df %>%
  filter(year < 2017) %>%
  as_tibble(index="year")
```

Fit Holt's linear trend model and Holt's damped linear trend model to the training data.

```
holtFit = productivityTraining %>%
  model(holt = ETS(Productivity ~ error("A") + trend("M") + season("M")))

dampedFit = productivityTraining %>%
  model(damped = ETS(Productivity ~ error("A") + trend("M") + season("M")))

report(holtFit)

## Series: Productivity
## Model: ETS(A,Ad,M)
## Smoothing parameters:
## alpha = 0.4889842
## beta = 0.889288943
##
## Initial states:
## i[0] h[0]
## 889.7265 48.07191
##
## sigma^2 = 1844.32
##
## AIC AICC BIC
## 533.7385 535.5036 543.4903

report(dampedFit)

## Series: Productivity
## Model: ETS(A,Ad,M)
## Smoothing parameters:
## alpha = 0.5398882
## beta = 0.889288943
## phi = 0.98
##
## Initial states:
## i[0] h[0]
## 879.492 48.03462
##
## sigma^2 = 20988.5
##
## AIC AICC BIC
## 537.6485 548.2768 547.4289
```

Interpret the estimates for the model parameters (α, β, ϕ) of Holt's damped linear trend model.

α represents smoothing parameter for the level, $\alpha = 0.5398882$ shows that the level reacts moderately to each new observation. $\beta = 0.889288943$ shows that the level reacts moderately to each new observation. $\phi = 0.98$ shows that the level reacts moderately to each new observation. This means that the slope of the model change over time by small degree.

ϕ represents the dampness in the Holt's damped linear trend model. ϕ value is close to 1 which means that the forecast trend is not much damped and the forecast is trended.

Compare AICc. Which model has a better fit to the training data?

Comparing AICc: Holt model seems better than Holt's damped linear trend model at fitting to the training data.

Question 3

Based on the models fitted in part 2, do the following: [9 Marks]

Forecast the next 5 years into the future.

```
holtFc = holtFit %>%
  forecast(h=5)

dampedFc = dampedFit %>%
  forecast(h=5)
```

Create a plot where you overlay the point forecasts on the original data.

```
productivity.tibble = productivity.df %>%
  as_tibble(index="year")

holtFc %>% autoplot(productivity.tibble, level=NULL) + ggtitle("5 year Holt forecast of productivity (2017 - 2022)")

5 year Holt forecast of productivity (2017 - 2022)

dampedFc %>% autoplot(productivity.tibble, level=NULL) + ggtitle("5 year Holt Damped forecast of productivity (2017 - 2022)")

5 year Holt Damped forecast of productivity (2017 - 2022)
```

Compute appropriate measures of forecast accuracy. Comment on which model provides better forecasts.

```
accuracy(holtFc, productivity.tibble)

## # A tibble: 1 x 10
##   model_type ME RMSE MAE MPE MAPE RMSE RMSE ACF1
##   <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Holt Test -194. 288. 194. -0.43 5.43 1.92 1.31 -0.115

accuracy(dampedFc, productivity.tibble)

## # A tibble: 1 x 10
##   model_type ME RMSE MAE MPE MAPE RMSE RMSE ACF1
##   <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Damped Test -0.8 99.4 10.8 -2.37 2.37 0.658 0.68 -0.057
```

The root mean squared error (RMSE) is 99.45 for Damped model compared to Holt's model RMSE of 208.13. This means that the Holt's model is a much more accurate model for forecast.

Report 95% prediction interval for the year 2022 in the model with the better forecasts. Interpret this in plain English.

```
dampedFit %>%
  holo(h=5) %>%
  holo(h=5) %>%
  filter(year == "2022")

## # A tibble: 1 x 5 [Y]
## # Key:   model [Y]
##   model Year Productivity_mean 200%
##   <chr> <dbl> <dbl> <dbl> <dbl>
## 1 Damped 2022 N(3758, 58049) 3758. [3285.731, 4259.176]95
```

We estimate that the year 2022 on average will have the productivity index between 3285.73 and 4259.18.

Using the data you have available, discuss how you could you reduce the forecast uncertainty?

Using the data available, we could perform cross validation to reduce the forecast uncertainty. Cross validation such as evaluation on a rolling forecasting origin can help to choose a good forecasting model which can reduce the forecast uncertainty.

Problem 2

Manually determine an appropriate non-seasonal ARIMA model for the productivity training data (from Problem 1), by doing the following: [9 Marks]

Conduct a KPSS unit root test on the training data. Keep differencing the data until it is stationary. What is the order of differencing d?

```
productivityTraining %>% features(Productivity, unitroot_kpss)

## # A tibble: 1 x 2
##   kpss.stat kpss.pvalue
##   <dbl> <dbl>
## 1 1.06 0.01

productivityTraining %>%
  mutate(diff_productivity = difference(Productivity)) %>%
  features(diff_productivity, unitroot_kpss)

## # A tibble: 1 x 2
##   kpss.stat kpss.pvalue
##   <dbl> <dbl>
## 1 0.8528 0.1

productivityTraining %>% features(Productivity, unitroot_diffs)

## # A tibble: 1 x 1
##   ndiffs
##   <int>
## 1 1
```

The KPSS test has a null hypothesis that the data is stationary and non-seasonal. We have a p-value = 0.01 which means that we have a strong evidence that the data is not stationary and seasonal. Therefore, differencing was applied so that p-value is 0.1 which means that we have a strong evidence that the data is stationary. The order of differencing is 1.

Plot the ACF and PACF plots for the differenced data and comment on what you observe.

```
productivityTraining %>% ACF(Productivity %>% difference(1)) %>% autoplot() +
  labs(y = "ACF", title = "ACF of the differenced series")

ACF of the differenced series

productivityTraining %>% PACF(Productivity %>% difference(1)) %>% autoplot() +
  labs(y = "PACF", title = "PACF of the differenced series")

PACF of the differenced series

Based on ACF and PACF plots, we
```

can determine the working model.

We see a dampening sinusoidal pattern in the ACF plot, and a significant spike at lag 2 but none beyond lag 2.

Based on what you have learned on lectures, what ARIMA model would you suggest fitting to the training data?

Therefore, $p = 2$ and $d = 1$ so we should fit working model ARIMA(2,1,0) to the training data.

Write the equation of this model using backshift notation.

$(1 - \phi_1 B - \phi_2 B^2)(1 - B)Y_t = c$

Fit the following ARIMA models, compare them using information criteria, and write down the equation of the best model using backshift notation: [6 Marks]

```
fit = productivityTraining %>%
  model(ARIMA(Productivity ~ pdc(2,1,0)))

fit %>% report()

## Series: Productivity
## Model: ARIMA(0,1,1) w/ drift
##
## Coefficients:
##      e[0] constant
##      -0.4889 68.8464
## s.e. 0.2321 11.4906
##
## sigma^2 estimated as 184873: log likelihood=-248.09
## AIC=486.19 AICC=487.39 BIC=487.78
```

Your suggested model from part 1.

An automatic model using the stepwise algorithm.

```
autoStepfit <- productivityTraining %>%
  model(stepwise = AREMA(Productivity))

autoStepfit %>% report()

## Series: Productivity
## Model: ARIMA(0,1,1) w/ drift
##
## Coefficients:
##      e[0] constant
##      -0.4889 68.8464
## s.e. 0.2321 11.4906
##
## sigma^2 estimated as 184873: log likelihood=-248.09
## AIC=486.19 AICC=487.39 BIC=487.78
```

An automatic model without using the stepwise algorithm.

```
autoFit <- productivityTraining %>%
  model(stepsize = AREMA(Productivity, stepsize = FALSE))

autoFit %>% report()

## Series: Productivity
## Model: ARIMA(0,1,1) w/ drift
##
## Coefficients:
##      e[0] constant
##      -0.4889 68.8464
## s.e. 0.2321 11.4906
##
## sigma^2 estimated as 184873: log likelihood=-248.09
## AIC=486.19 AICC=487.39 BIC=487.78
```

The automatic model using stepwise and without stepwise have both fitted ARIMA(0,1,1). They both have the same information criteria and AICc value of 486.19 which is lower than the manual fit (AICc=487.39). The best model using backshift notation $(1 - B)Y_t = c + (1 + \phi_1 B)Y_t$.

Using the best model found in part 2, do the following: [5 Marks]

Conduct a diagnostic check on the residuals. Discuss whether or not you have any concerns about the model assumptions.

Based on the AICc results.

```
autoFit %>% gg_tsresiduals()

## Residuals:
##      e[0] constant
##      -0.4889 68.8464
## s.e. 0.2321 11.4906
##
## sigma^2 estimated as 184873: log likelihood=-248.09
## AIC=486.19 AICC=487.39 BIC=487.78
```

than the actual data.