

# Project: The Forecasting Tourism 2010 Competition

EM1415

Marco Solari, 875475

## Table of contents

Setup and Data Loading	2
Setup	2
Loading Data	2
Creating tsibble	3
Assignment	3
Full Plot	3
Creating Validation Set	4
Naïve Forecasts	5
Choosing Measures	6
Computing MAPE	6
Computing MASE	6
MAPE Pairs	6
Ensemble Methods	6
a. Write the exact formula used for generating the first method, in the form $F_{t+k} = \dots (k = 1, 2, 3, 4)$ ,	6
b. What is the rationale behind multiplying the naive forecasts by a constant?	6
c. What should be the dependent variable and the predictors in a linear regression model for this data? Explain..	6
d. Fit the linear regression model to the first five series and compute forecast errors for the validation period.	6
e. Before choosing a linear regression, the winner described the following process:	7
f. If we were to consider exponential smoothing, what particular type(s) of exponential smoothing are reasonable candidates?	7
g. The winner concludes with possible improvements one being "an investigation into how to come up with a blending ensemble method that doesn't use much manual twerking would also be of benefit". Can you suggest methods or an approach that would lead to easier automation of the ensemble step?	7
h. The competition focused on minimizing the average MAPE of the next four values across all 518 series. How does this goal differ from goals encountered in practice when considering tourism demand? Which steps in the forecasting process would likely be different in a real-life tourism forecasting scenario?	7

## Setup and Data Loading

### Setup

```
knitr::opts_chunk$set(  
  echo = T,  
  dev = "cairo_pdf"  
)  
  
libraries_list <- c(  
  "tidyverse",  
  "fpp3"  
)  
  
lapply(  
  X = libraries_list,  
  FUN = require,  
  character.only = TRUE  
)
```

```
[[1]]  
[1] TRUE
```

```
[[2]]  
[1] TRUE
```

```
theme_set(  
  ggthemes::theme_tufte(  
    base_size = 16,  
    base_family = "Atkinson Hyperlegible"  
  )  
)
```

### Loading Data

```
data_main <- readr::read_csv(  
  "Data/tourism_data.csv",  
  show_col_types = F  
)
```

```
data_main %>% dim
```

```
[1] 43 518
```

```
data_main %>% is.na() %>% sum
```

```
[1] 11668
```

We are missing 52.38% of the observations.

## Creating **tsibble**

```
tourism_full <- data_main %>%  
  mutate(  
    Year = 1965:2007  
  ) %>%  
  as_tsibble(  
    index = Year  
  )
```

## Assignment

### Full Plot

*Plot all the series (an advanced data visualization tool is recommended) - what type of components are visible? Are the series similar or different? Check for problems such as missing values and possible errors.*

```
tmelt <- reshape2::melt(tourism_full, id="Year")
```

```
tmelt %>%  
  ggplot(  
    aes(  
      x = Year,  
      y = value,  
      colour = variable,  
      group = variable  
    )  
  ) +  
  geom_line(  
    alpha = .8  
  ) +  
  scale_y_log10() +  
  scale_color_viridis_d(  
    option = "rocket"  
  ) +  
  labs(  
    title = "Tourism Time Series",  
    y = "Value"  
  ) +  
  theme(  
    legend.position = "none"  
  )
```

Warning: Removed 11668 rows containing missing values (`geom\_line()`).

## Tourism Time Series

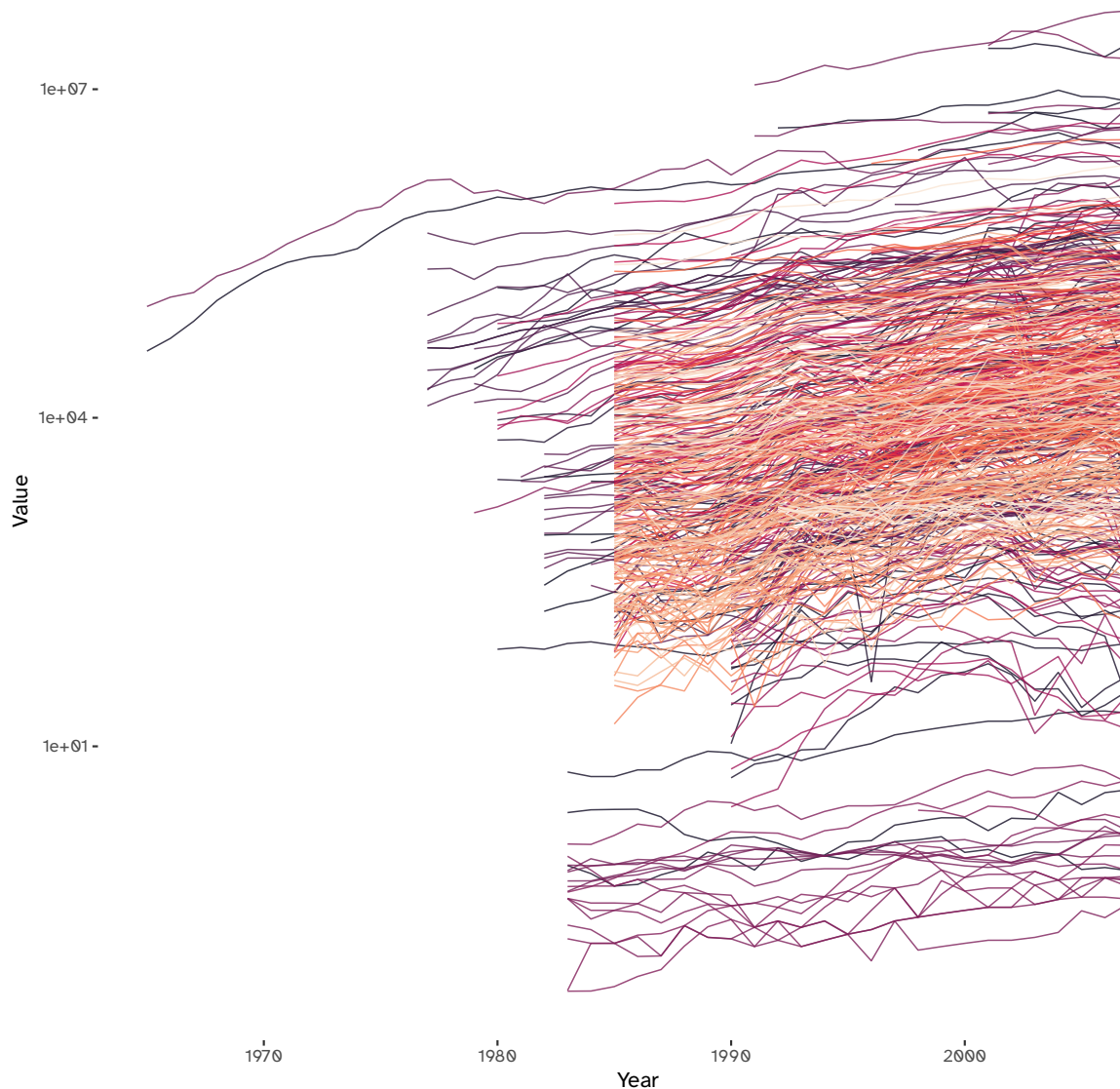


Figure 1: Printing a legend for 518 different series is not possible. However, color has been used only to differentiate the series and does not contain further information. Plotting the y-axis variable on the log scale was made necessary by the huge variation in the series values.

A check for NAs has already been made while loading data and it showed the presence of a large number of missing values. This is mostly to be attributed to the different starting time of the series. From the visualization we can definitely group different starting points: this suggests that another approach to visualisation might be successfully conducted by grouping series by their starting date.

Plotting all 518 series does not allow to spot details, such as the presence of seasonal patterns. However, a general upward trend is clear; moreover, we can spot some probable outliers, that should be further investigated, and some clues about the presence of cyclicity in some of the series.

## Creating Validation Set

*Partition the series into training and validation, so that the last 4 years are in the validation period for each series. What is the logic of such a partitioning? What is the disadvantage?*

```

tourism_train <- tourism_full %>%
  filter(Year < 2004)
tourism_validation <- tourism_full %>%
  filter(Year >= 2004)

```

The logic behind partitioning the series into a *training* and *validation* set is to *estimate the forecasting error*: we can train a model or apply a filter to the train set and use it to assess its performance with out-of-sample data. The main disadvantage with this approach is that we are not using all the information available to train our model; moreover, we are not computing *true forecasts*, therefore the accuracy measures from the residuals will be smaller.

## Naïve Forecasts

*Generate naïve forecasts for all series for the validation period. For each series, create forecasts with horizons of 1, 2, 3, and 4 years ahead ( $F_{t+1}$ ,  $F_{t+2}$ ,  $F_{t+3}$ , and  $F_{t+4}$ ).*

We know that *naïve forecasts* consist in the last observation,  $\forall h$ .

$$y_{T+h} | T = y_T$$

It follows that we can produce the forecasts with the following code:

```

naive_forecast <- tourism_train %>%
  filter(
    Year == 2003
  ) %>% as_tibble()

```

`naive_forecast` will contain a `tsibble` with  $y_t$  for all the series.

```
naive_forecast %>% dim
```

```
[1] 1 519
```

To obtain  $F_{t+1}$ ,  $F_{t+2}$ ,  $F_{t+3}$ , and  $F_{t+4}$ :

```

merge(
  x = 2003 + 1:4,
  y = naive_forecast
) %>%
mutate(
  Year = x
) %>%
select(
  -x
) %>%
as_tibble(
  index = Year
)

```

```
# A tsibble: 4 x 519 [1Y]
```

	Y1	Y2	Y3	Y4	Y5	Y6	Y7	Y8	Y9	Y10	Y11	Y12
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	30217.	351175	174960	78519	14280	2340	6442	7181	7781	2613	1462	604
2	30217.	351175	174960	78519	14280	2340	6442	7181	7781	2613	1462	604
3	30217.	351175	174960	78519	14280	2340	6442	7181	7781	2613	1462	604
4	30217.	351175	174960	78519	14280	2340	6442	7181	7781	2613	1462	604

```
# i 507 more variables: Y13 <dbl>, Y14 <dbl>, Y15 <dbl>, Y16 <dbl>, Y17 <dbl>,
#   Y18 <dbl>, Y19 <dbl>, Y20 <dbl>, Y21 <dbl>, Y22 <dbl>, Y23 <dbl>,
#   Y24 <dbl>, Y25 <dbl>, Y26 <dbl>, Y27 <dbl>, Y28 <dbl>, Y29 <dbl>,
#   Y30 <dbl>, Y31 <dbl>, Y32 <dbl>, Y33 <dbl>, Y34 <dbl>, Y35 <dbl>,
#   Y36 <dbl>, Y37 <dbl>, Y38 <dbl>, Y39 <dbl>, Y40 <dbl>, Y41 <dbl>,
#   Y42 <dbl>, Y43 <dbl>, Y44 <dbl>, Y45 <dbl>, Y46 <dbl>, Y47 <dbl>,
#   Y48 <dbl>, Y49 <dbl>, Y50 <dbl>, Y51 <dbl>, Y52 <dbl>, Y53 <dbl>, ...
```

## Choosing Measures

*Which measures are suitable if we plan to combine the results for the 518 series? Consider MAE, Average error, MAPE and RMSE.*

## Computing MAPE

*For each series, compute MAPE of the naive forecasts once for the training period and once for the validation period.*

## Computing MASE

*The performance measure used in the competition is Mean Absolute Scaled Error (MASE). Explain the advantage of MASE and compute the training and validation MASE for the naive forecasts.*

## MAPE Pairs

*Create a scatter plot of the MAPE pairs, with the training MAPE on the x-axis and the validation MAPE on the y-axis. Create a similar scatter plot for the MASE pairs. Now examine both plots. What do we learn? How does performance differ between the training and validation periods? How does performance range across series?*

## Ensemble Methods

*The competition winner, Lee Baker, used an ensemble of three methods:*

- Naive forecasts multiplied by a constant trend<sup>1</sup>.
- Linear regression.
- Exponentially-weighted linear regression.

<sup>1</sup> Global/local trend: "globally tourism has grown "at a rate of 6% annually."

a. Write the exact formula used for generating the first method, in the form  $F_{t+k} = \dots (k = 1, 2, 3, 4)$ ,

b. What is the rationale behind multiplying the naive forecasts by a constant?<sup>2</sup>

<sup>2</sup> Hint: think empirical and domain knowledge.

c. What should be the dependent variable and the predictors in a linear regression model for this data? Explain..

d. Fit the linear regression model to the first five series and compute forecast errors for the validation period.

```
train_subset <- tourism_train %>%
  select(
    Y1,
    Y2,
    Y3,
    Y4,
```

Y5  
)

e. *Before choosing a linear regression, the winner described the following process:*

"I examined fitting a polynomial line to the data and using the line to predict future values. I tried using first through fifth order polynomials to find that the lowest MASE was obtained using a first order polynomial (simple regression line). This best fit line was used to predict future values. I also kept the  $R^2$  value of the fit for use in blending the results of the prediction."

*What are two flaws in this approach?*

f. *If we were to consider exponential smoothing, what particular type(s) of exponential smoothing are reasonable candidates?*

g. *The winner concludes with possible improvements one being "an investigation into how to come up with a blending ensemble method that doesn't use much manual twerking would also be of benefit". Can you suggest methods or an approach that would lead to easier automation of the ensemble step?*

h. *The competition focused on minimizing the average MAPE of the next four values across all 518 series. How does this goal differ from goals encountered in practice when considering tourism demand? Which steps in the forecasting process would likely be different in a real-life tourism forecasting scenario?*