# Ozonio model

Giovani Valdrighi, Vitória Guardieiro

24/09/2020

# R Markdown

This is an R Markdown presentation. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see http://rmarkdown.rstudio.com.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

# Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

# Slide with R Output

```r
summary(cars)
```
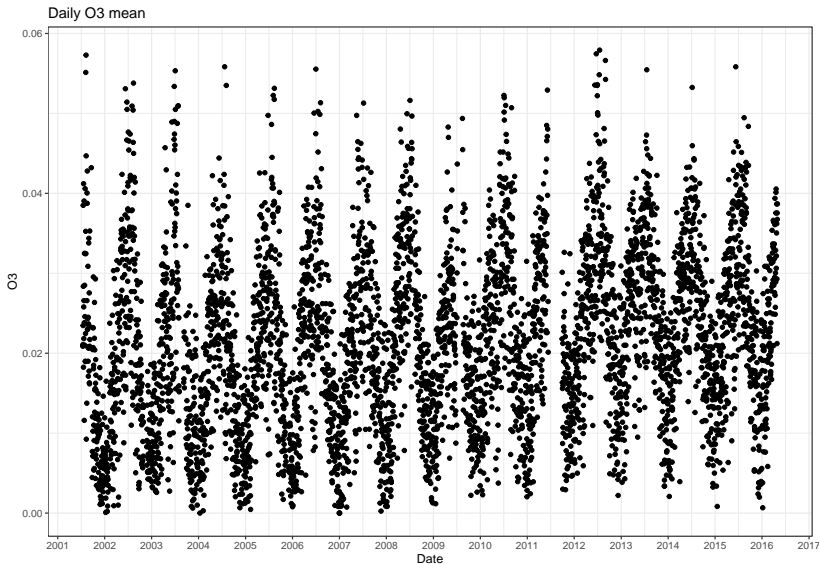
```
##      speed           dist
##  Min.   : 4.0   Min.   :  2.00
##  1st Qu.:12.0   1st Qu.: 26.00
##  Median :15.0   Median : 36.00
##  Mean   :15.4   Mean   : 42.98
##  3rd Qu.:19.0   3rd Qu.: 56.00
##  Max.   :25.0   Max.   :120.00
```
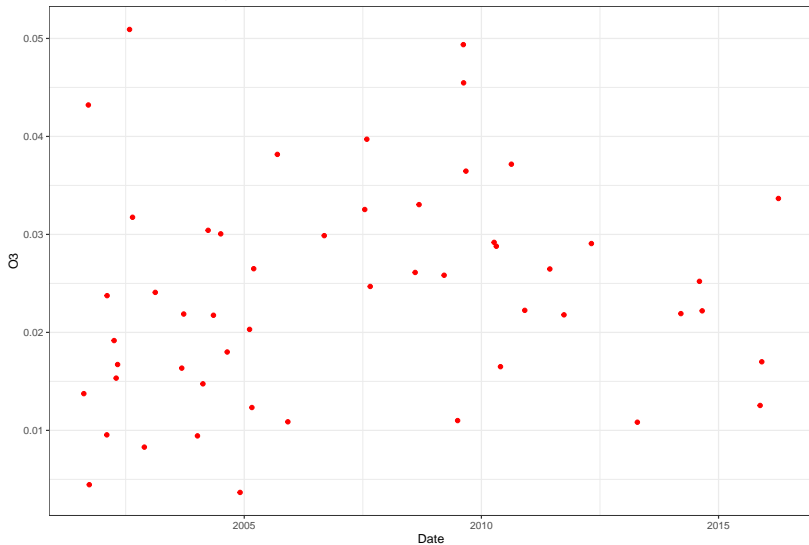
# Slide with Plot

# Data

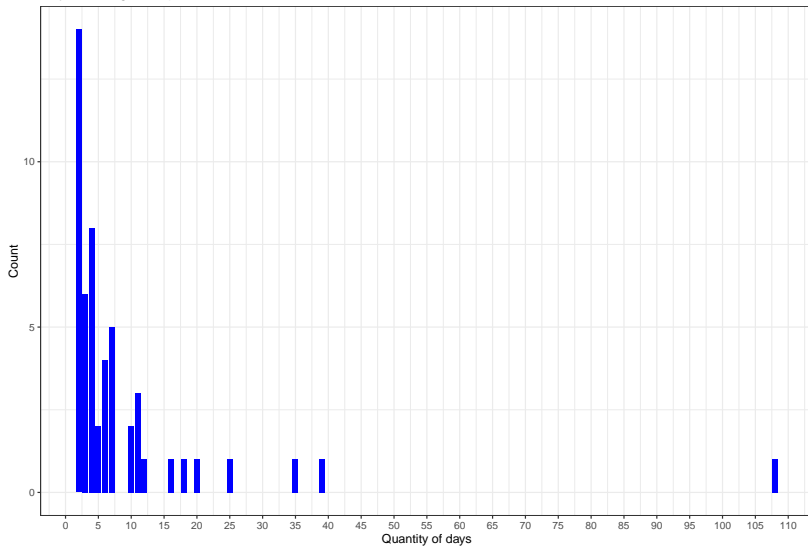▶ New York data from 15/07/2001 to 30/04/2016.



Daily O3 mean

# Missing data

- There are 52 time skips in the data, in a total of 473 days.
- The biggest skips is 108 days in 2011.
- The majority of skips are of 1 or 2 days.
- Around 9.5% missing data.
- The missing observations are distributed along the time without a clear pattern.
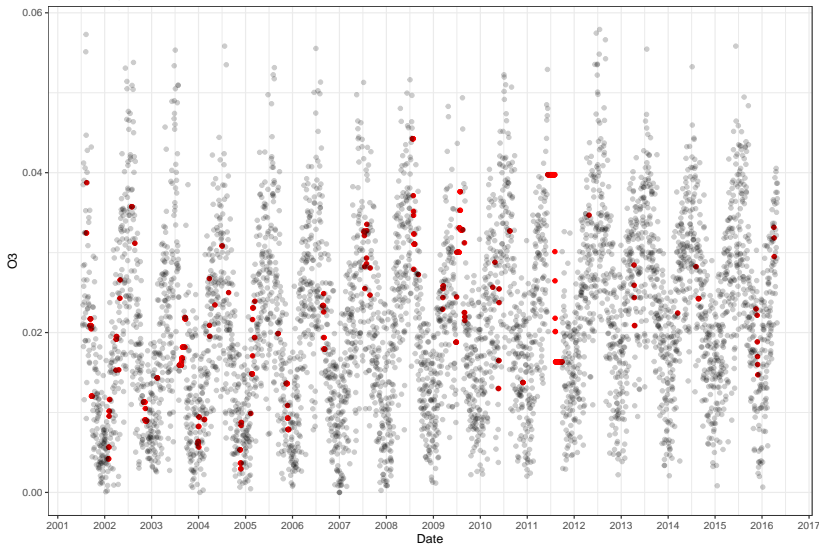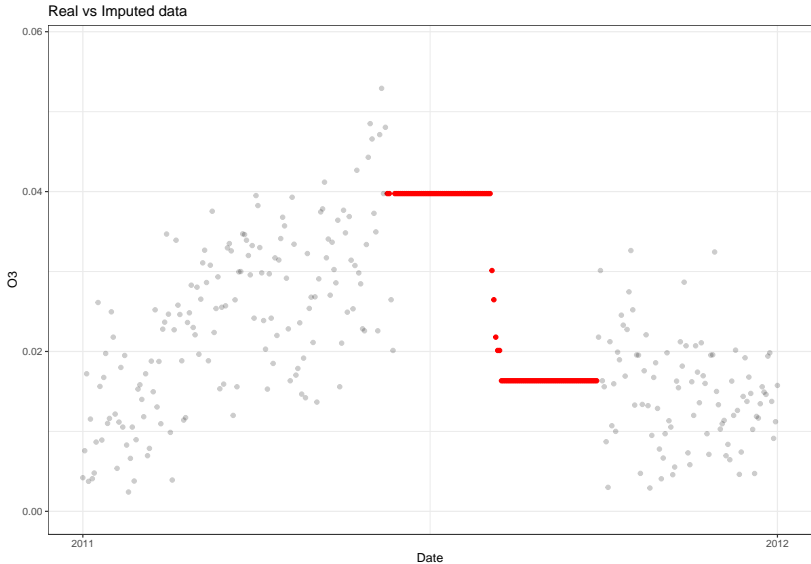
Observations after data skips

Days missing in sequence

# Imputation method

- It was used the kNN method to imputate values on missing observations.
- The kNN method needs the parameter k, the number of closest points considered.
- Starting with $k = 7$.

Real vs Imputed data

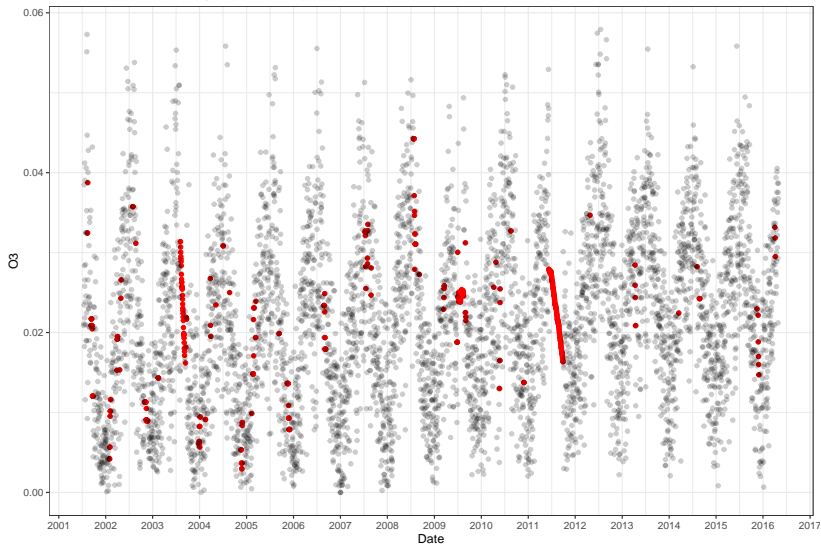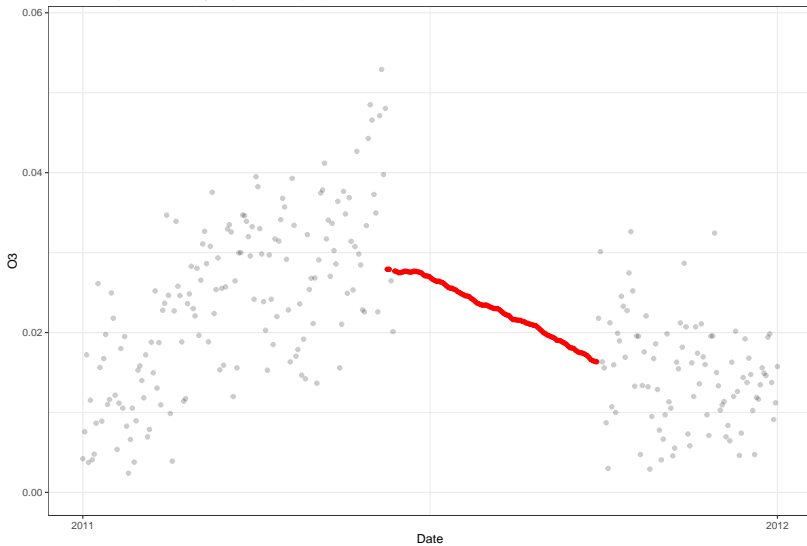▶ Method create a bad behavior where the size of the skips is bigger than 7 days.



Real vs Imputed data

- ▶ To deal with this, the parameter k used for imputation will be different if the size of the skip is minor tem 30 days, between 30 days and 100 days, or bigger than 100 days.
- ▶ k = 7, k = 45, k = 120, respectively.
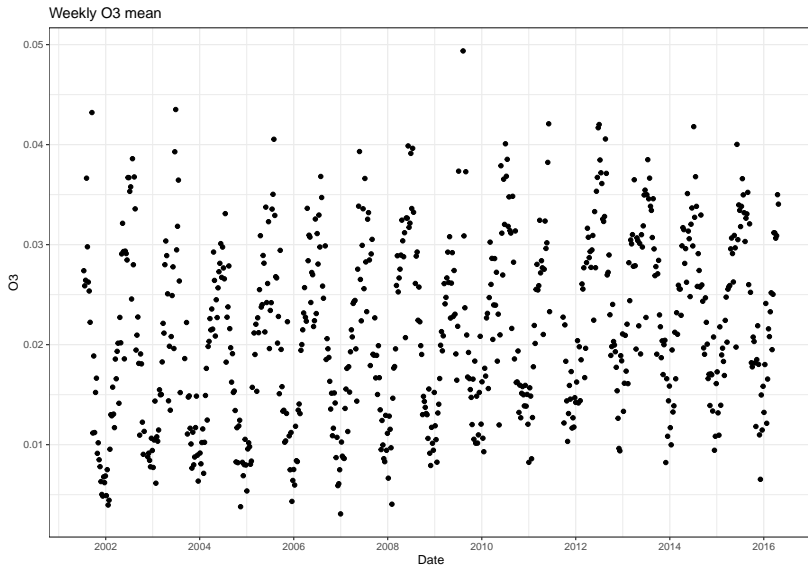- ▶ We will aggregate closest points by weighted by distance mean.

Real vs Imputed data (by separeted imput.)

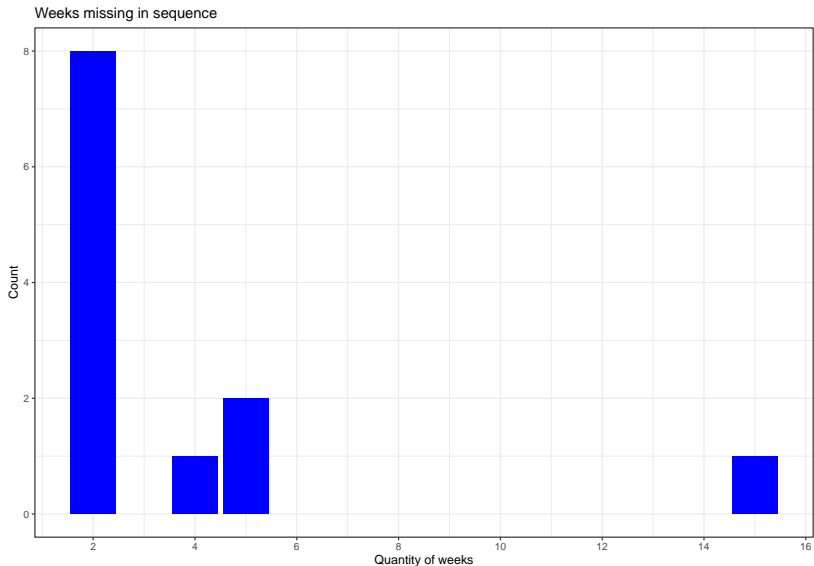Real vs Imputed data (by separated imput.)
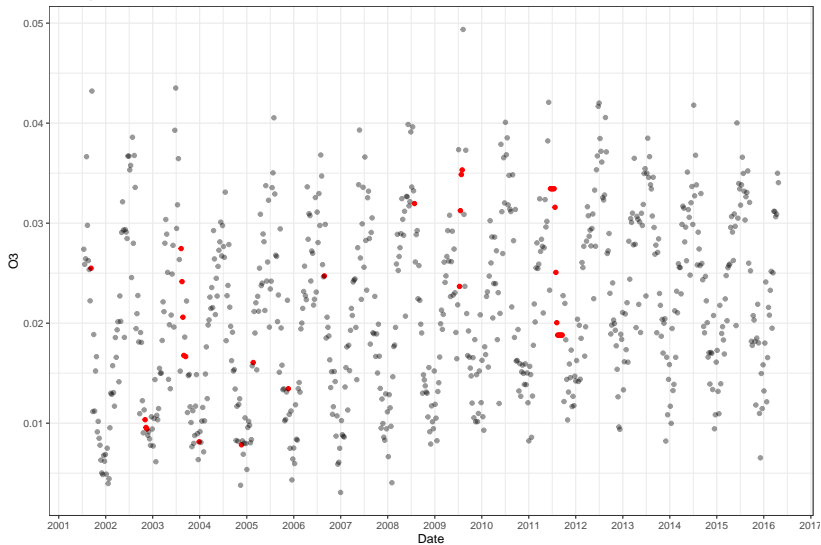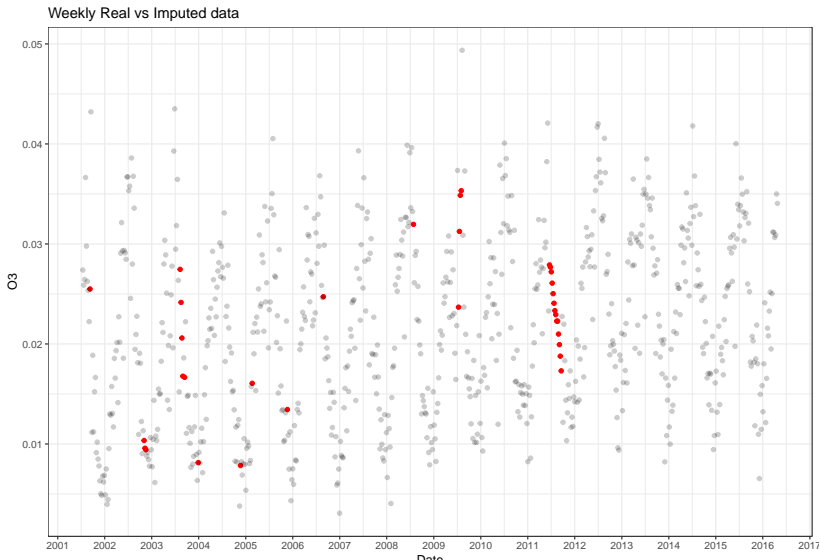
# Weekly data



Weekly O3 mean

- ► If the data is grouped by week, ignoring the missing values when aggregating, it'll have 33 missing observations.
- ► Around 4.3% missing data.



Weeks missing in sequence

Weekly Real vs Imputed data

- ▶ It has the same problem when the sequence of missing data is to big.
- ▶ Again, if there is more than 5 missing weeks, it will be used k = 16, if it's less, it'll be k = 4.



Weekly Real vs Imputed data

Weekly Real vs Imputed data