

Model picked for missing values and adding local trend and seasonality - Kalman Filter and Smoother for Exponential Family State Space Models Initial Claims in Texas - Covariate series considered - <https://fred.stlouisfed.org/series/TXICLAIMS> (<https://fred.stlouisfed.org/series/TXICLAIMS>)

```
#Load the required Libraries
```

```
if (!require("imputeTS")) install.packages("imputeTS")
```

```
## Loading required package: imputeTS
```

```
## Warning: package 'imputeTS' was built under R version 4.3.2
```

```
## Registered S3 method overwritten by 'quantmod':
```

```
##   method      from  
##   as.zoo.data.frame zoo
```

```
library(fredr)
```

```
## Warning: package 'fredr' was built under R version 4.3.2
```

```
library(forecast)
```

```
## Warning: package 'forecast' was built under R version 4.3.2
```

```
library(reprex)
```

```
library(zoo)
```

```
## Warning: package 'zoo' was built under R version 4.3.2
```

```
##
```

```
## Attaching package: 'zoo'
```

```
## The following object is masked from 'package:imputeTS':
```

```
##  
##   na.locf
```

```
## The following objects are masked from 'package:base':
```

```
##  
##   as.Date, as.Date.numeric
```

```
library(imputeTS)
```

```
library(forecast)
```

```
library(dlm)
```

```
## Warning: package 'dlm' was built under R version 4.3.3
```

```
library(bsts)
```

```
## Warning: package 'bsts' was built under R version 4.3.3
```

```
## Loading required package: BoomSpikeSlab
```

```
## Warning: package 'BoomSpikeSlab' was built under R version 4.3.3
```

```
## Loading required package: Boom
```

```
## Warning: package 'Boom' was built under R version 4.3.3
```

```
##  
## Attaching package: 'Boom'
```

```
## The following object is masked from 'package:stats':  
##  
##   rWishart
```

```
##  
## Attaching package: 'BoomSpikeSlab'
```

```
## The following object is masked from 'package:stats':  
##  
##   knots
```

```
## Loading required package: xts
```

```
## Warning: package 'xts' was built under R version 4.3.2
```

```
##  
## Attaching package: 'bsts'
```

```
## The following object is masked from 'package:BoomSpikeSlab':  
##  
##   SuggestBurn
```

```
library(KFAS)
```

```
## Warning: package 'KFAS' was built under R version 4.3.3
```

```
## Please cite KFAS in publications by using:
```

```
##  
## Jouni Helske (2017). KFAS: Exponential Family State Space Models in R. Journal of Statistic  
al Software, 78(10), 1-39. doi:10.18637/jss.v078.i10.
```

```
#Load the data and view the sample entries
```

```
fredr_set_key("e3dbddf1961174f1b64d469145d6eff9")  
icsa_data = fredr(series_id = "ICNSA")  
head(icsa_data)
```

```
## # A tibble: 6 × 5  
##   date      series_id value realtime_start realtime_end  
##   <date>   <chr>     <dbl> <date>         <date>  
## 1 1967-01-07 ICNSA     346000 2024-03-07     2024-03-07  
## 2 1967-01-14 ICNSA     334000 2024-03-07     2024-03-07  
## 3 1967-01-21 ICNSA     277000 2024-03-07     2024-03-07  
## 4 1967-01-28 ICNSA     252000 2024-03-07     2024-03-07  
## 5 1967-02-04 ICNSA     274000 2024-03-07     2024-03-07  
## 6 1967-02-11 ICNSA     276000 2024-03-07     2024-03-07
```

```
#Drop insignificant columns
```

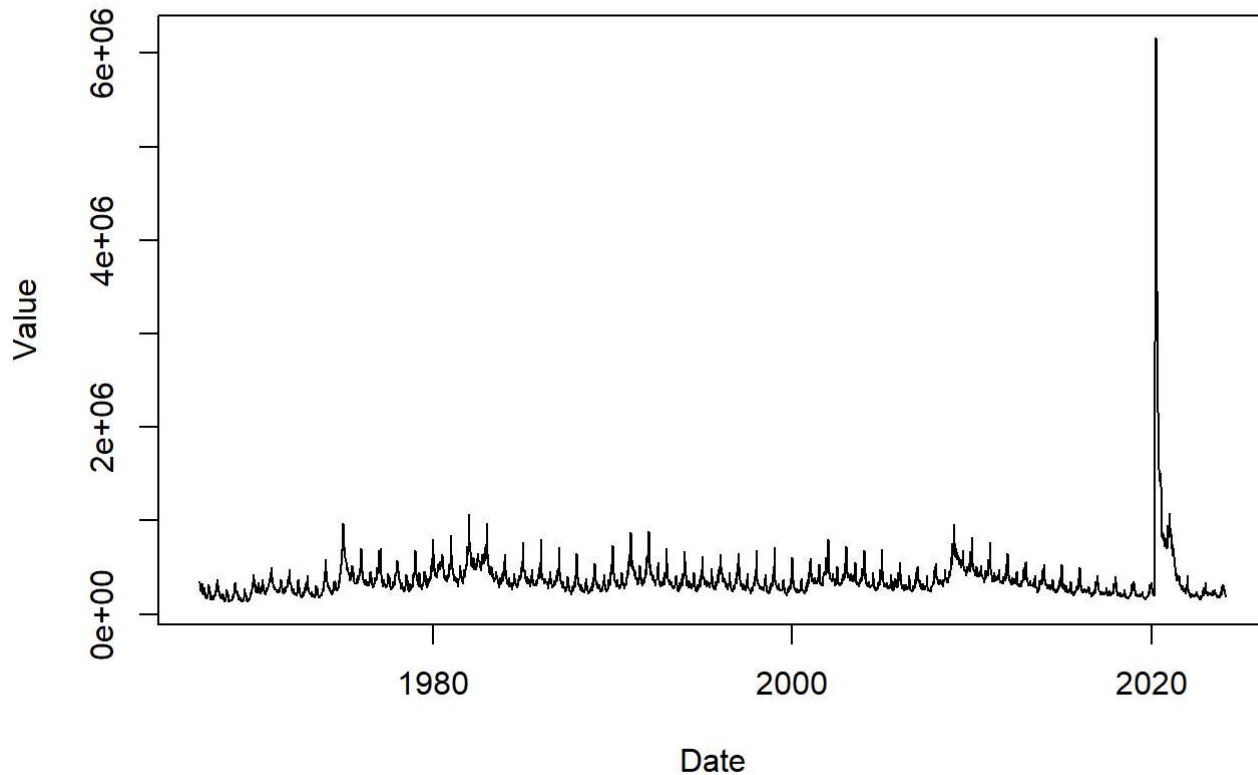
```
icsa_data <- icsa_data[, c("date", "value")]  
head(icsa_data)
```

```
## # A tibble: 6 × 2  
##   date      value  
##   <date>     <dbl>  
## 1 1967-01-07 346000  
## 2 1967-01-14 334000  
## 3 1967-01-21 277000  
## 4 1967-01-28 252000  
## 5 1967-02-04 274000  
## 6 1967-02-11 276000
```

```
#Time Series plot of ICNSA data
```

```
plot(icsa_data$date, icsa_data$value, type = "l", xlab = "Date", ylab = "Value", main = "Time  
Series Plot")
```

## Time Series Plot



## Exploratory Data Analysis

```
#Displays number of rows and columns in the dataset
dim(icnsa_data)
```

```
## [1] 2982    2
```

```
#Datatype of each column of the Iris Dataset
str(icnsa_data)
```

```
## tibble [2,982 × 2] (S3: tbl_df/tbl/data.frame)
## $ date : Date[1:2982], format: "1967-01-07" "1967-01-14" ...
## $ value: num [1:2982] 346000 334000 277000 252000 274000 276000 247000 248000 326000 240000
...

```

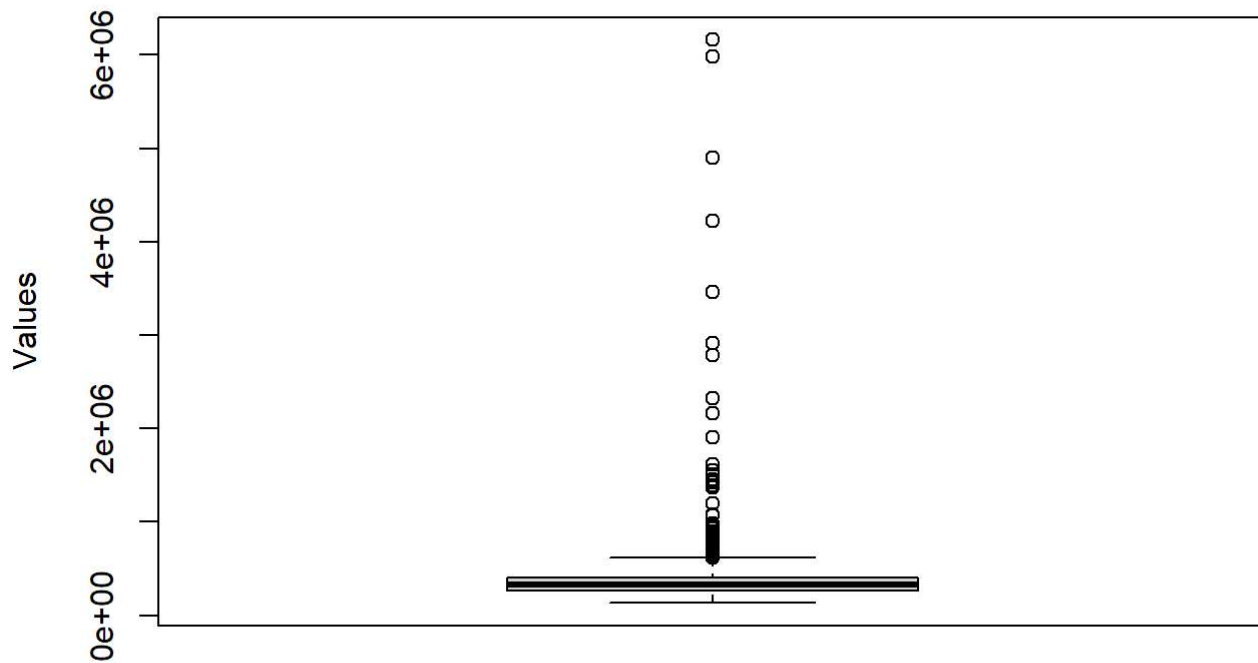
```
#Count number of null values in each column - No null values found in any column
```

```
Icnsa_NullValues<-as.data.frame(sapply(icnsa_data, function(x) sum(is.na(x))))
colnames(Icnsa_NullValues) <- c('Null Value Count')
Icnsa_NullValues
```

```
##      Null Value Count
## date              0
## value             0
```

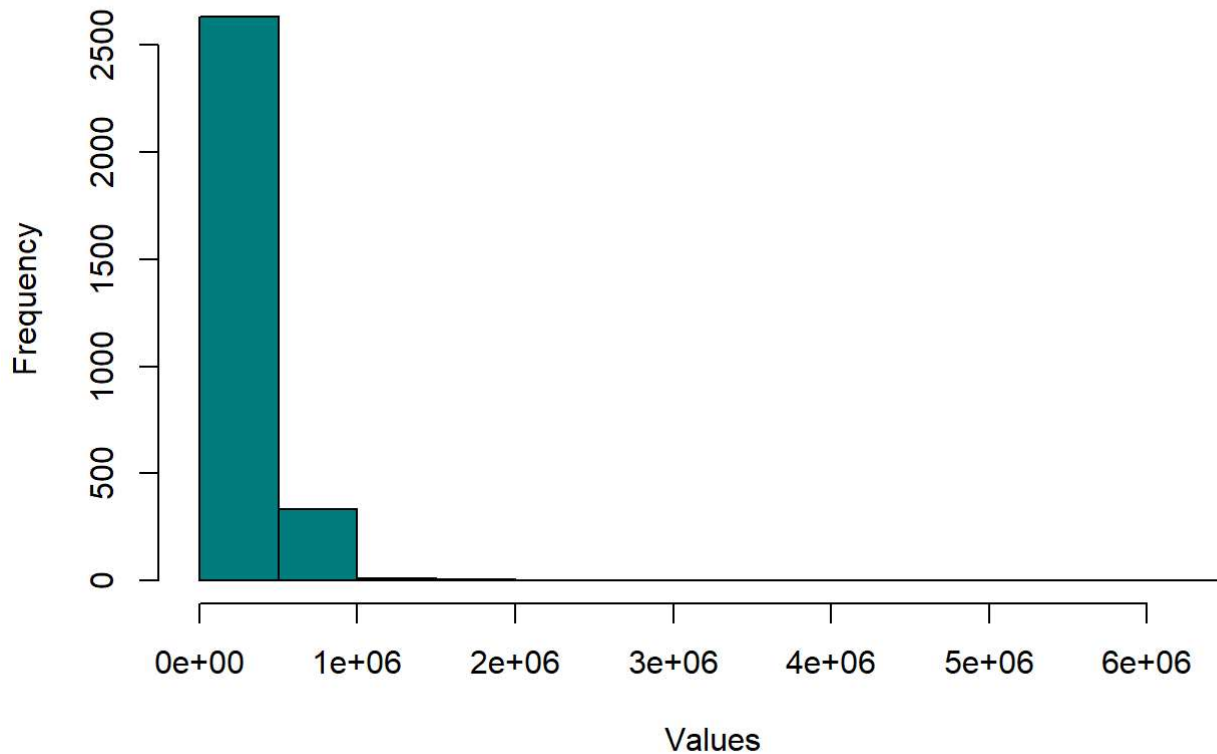
```
#Box-plot to check the outliers
boxplot(icsa_data['value'], main = "Boxplot with Outliers", ylab = "Values")
```

### Boxplot with Outliers



```
#Histogram to check the distribution of the data
hist(icsa_data$value, main = "Data Distribution ", xlab = "Values", ylab = "Frequency", col="#008080")
```

## Data Distribution



#Since the data doesn't have normal distribution, we have to use IQR to find the outliers in the data instead of Z-score

```
#IQR for calculating outliers range in the data

data <- icnsa_data$value

Q1 <- quantile(data, 0.25)
Q3 <- quantile(data, 0.75)

# Calculate IQR
IQR <- Q3 - Q1

#Threshold Multiplier
threshold <- 1.5

# Calculate lower and upper bounds for outliers
lower_bound <- Q1 - threshold * IQR
upper_bound <- Q3 + threshold * IQR

# Identify outliers
outliers <- data[data < lower_bound | data > upper_bound]

print(outliers)
```

```
## [1] 693000 637000 677000 813000 681000 969000 850000 729000 699000
## [10] 691000 708000 685000 704000 680000 804000 648000 643000 628000
## [19] 839000 638000 726000 657000 644000 702000 694300 1073500 761700
## [28] 771200 692300 671000 631400 647000 638100 650400 709400 638200
## [37] 653600 745100 976600 773600 650600 621600 637900 770000 803900
## [46] 710493 654620 730995 637449 649471 651775 872742 691092 625772
## [55] 652046 882118 687914 704930 676076 637910 654473 682016 713805
## [64] 647045 637343 799246 620929 620004 724111 677897 693776 760481
## [73] 629867 719691 717000 731958 956791 763987 620143 682176 710152
## [82] 619951 645827 652635 623279 677038 673097 651215 825891 659173
## [91] 773499 646219 2914268 5981838 6161268 4898119 4221704 3466665 2790860
## [100] 2325889 2162822 1902006 1610909 1555824 1456353 1446176 1425969 1390315
## [109] 1511632 1372247 1200962 982808 830266 878555 813744 826436 857006
## [118] 784765 814441 731391 722387 820590 757745 731660 736618 719716
## [127] 708717 827447 710901 946091 928178 854113 823512 898597 1082696
## [136] 930184 833704 843702 860209 827835 707268 749251 714661 750287
## [145] 630621 701073 668346 620420
```

```
print(lower_bound)
```

```
##      25%
## 54563.62
```

```
print(upper_bound)
```

```
##      75%
## 617740.6
```

```
#Print the dates for which outliers exist and filter the data post 2020 to get covid period
```

```
class(icnsa_data$date)
```

```
## [1] "Date"
```

```
outlier_dates <- c()
for (value in outliers) {
  outlier_rows <- which(icnsa_data$value == value)
  outlier_date <- icnsa_data$date[outlier_rows]
  if (as.numeric(format(outlier_date, "%Y")) >= 2020){
    print(outlier_date)
  }
}
```

```
## [1] "2020-03-21"  
## [1] "2020-03-28"  
## [1] "2020-04-04"  
## [1] "2020-04-11"  
## [1] "2020-04-18"  
## [1] "2020-04-25"  
## [1] "2020-05-02"  
## [1] "2020-05-09"  
## [1] "2020-05-16"  
## [1] "2020-05-23"  
## [1] "2020-05-30"  
## [1] "2020-06-06"  
## [1] "2020-06-13"  
## [1] "2020-06-20"  
## [1] "2020-06-27"  
## [1] "2020-07-04"  
## [1] "2020-07-11"  
## [1] "2020-07-18"  
## [1] "2020-07-25"  
## [1] "2020-08-01"  
## [1] "2020-08-08"  
## [1] "2020-08-15"  
## [1] "2020-08-22"  
## [1] "2020-08-29"  
## [1] "2020-09-05"  
## [1] "2020-09-12"  
## [1] "2020-09-19"  
## [1] "2020-09-26"  
## [1] "2020-10-03"  
## [1] "2020-10-10"  
## [1] "2020-10-17"  
## [1] "2020-10-24"  
## [1] "2020-10-31"  
## [1] "2020-11-07"  
## [1] "2020-11-14"  
## [1] "2020-11-21"  
## [1] "2020-11-28"  
## [1] "2020-12-05"  
## [1] "2020-12-12"  
## [1] "2020-12-19"  
## [1] "2020-12-26"  
## [1] "2021-01-02"  
## [1] "2021-01-09"  
## [1] "2021-01-16"  
## [1] "2021-01-23"  
## [1] "2021-01-30"  
## [1] "2021-02-06"  
## [1] "2021-02-13"  
## [1] "2021-02-20"  
## [1] "2021-02-27"  
## [1] "2021-03-06"  
## [1] "2021-03-13"
```



```
## [1] "2021-03-20"
## [1] "2021-03-27"
## [1] "2021-04-03"
## [1] "2021-04-10"
```

#Observation: 1. Therefore, "2020-03-21" to "2021-04-10" is the covid period from the above data

##3. Use the state-space missing value methodology discussed in lecture to impute new values for the Covid period. Compare this imputation method to that of Hw 2 when you used a spline to impute values

```
#Define Covid period

covid_start <- as.Date("2020-03-21")
covid_end <- as.Date("2021-04-10")

# Filter the relevant COVID period
covid_data <- icnsa_data[icnsa_data$date >= covid_start & icnsa_data$date <= covid_end, ]

# Define a function to create a BSTS model
create_bsts_model <- function(data) {
  model <- AddLocalLinearTrend(list(), data$value) #Adds Local Linear trend to the time series such that model can adapt to the changes in trend over time
  model <- AddSeasonal(model, data$value, nseasons = 12) #Adds seasonal component to the time series to account for seasonal variations in the data
  return(model)
}

#Model Training
bsts_model <- create_bsts_model(icnsa_data)
filtered_state <- bstsmcmc(icnsa_data$value, state.specification = bstsmcmc_model, niter=1000)
```

```
## ===== Iteration 0 Thu Mar 7 08:07:26 2024
## =====
## ===== Iteration 100 Thu Mar 7 08:07:29 2024
## =====
## ===== Iteration 200 Thu Mar 7 08:07:32 2024
## =====
## ===== Iteration 300 Thu Mar 7 08:07:36 2024
## =====
## ===== Iteration 400 Thu Mar 7 08:07:42 2024
## =====
## ===== Iteration 500 Thu Mar 7 08:07:52 2024
## =====
## ===== Iteration 600 Thu Mar 7 08:08:01 2024
## =====
## ===== Iteration 700 Thu Mar 7 08:08:07 2024
## =====
## ===== Iteration 800 Thu Mar 7 08:08:10 2024
## =====
## ===== Iteration 900 Thu Mar 7 08:08:13 2024
## =====
```

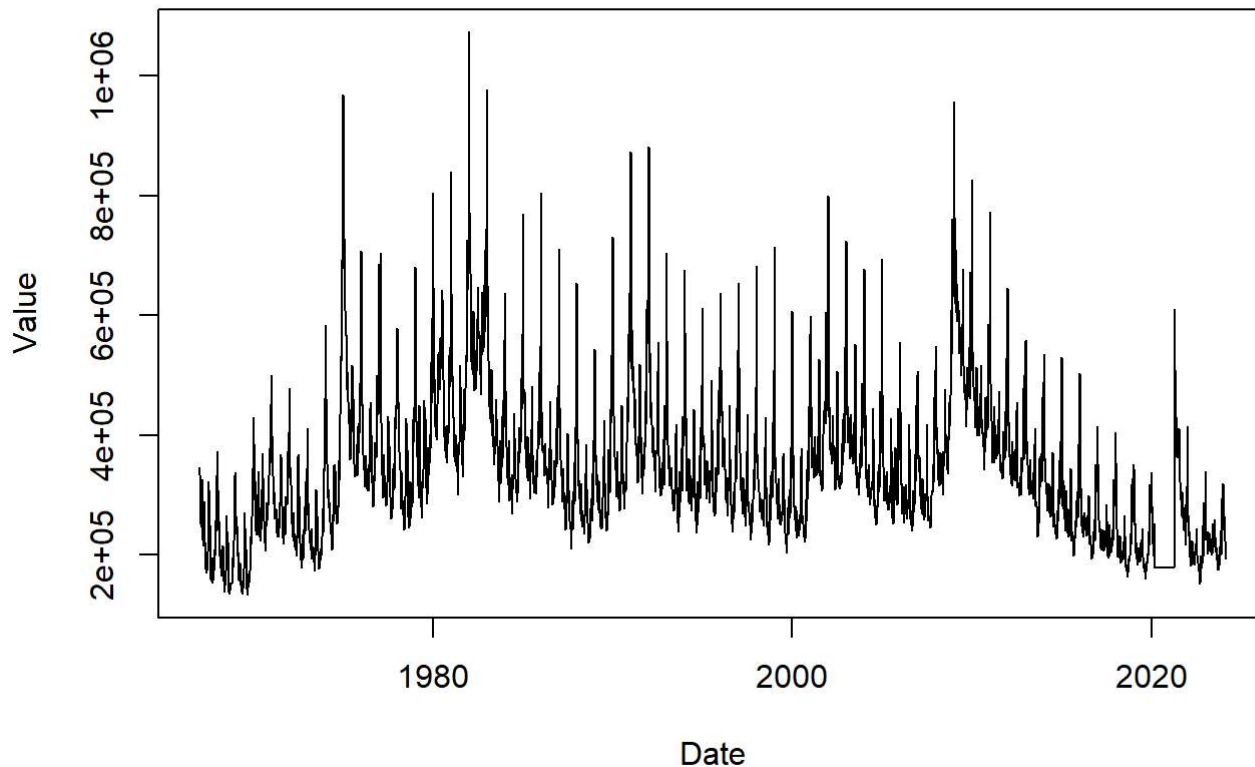
```
predict(filtered_state)$mean
```

```
## [1] 178595.1
```

```
imputed_values <- predict(filtered_state, range=c(covid_start, covid_end))  
imputed_values_covid <- imputed_values$mean  
covid_data$value <- imputed_values_covid  
icnsa_data[icnsa_data$date >= covid_start & icnsa_data$date <= covid_end, ] <- covid_data
```

```
#Time Series plot of ICNSA data post outlier treatment  
plot(icnsa_data$date, icnsa_data$value, type = "l", xlab = "Date", ylab = "Value", main = "Time  
Series Plot")
```

## Time Series Plot



**#Observation:**

1. Mean of the predicted values was used to impute the covid period values of ICNSA data
2. SSM can capture patterns in data and adds local trend and seasonality to the data but may give some uncertainty in estimates

```
##SSM for point forecast
```

```
library(forecast)
ts_data <- ts(icnsa_data$value, frequency = 12)
ets_model <- ets(ts_data)
forecast_ets <- forecast(ets_model, h = 1)
print("ETS Forecast:")
```

```
## [1] "ETS Forecast:"
```

```
print(forecast_ets)
```

```
##           Point Forecast   Lo 80   Hi 80   Lo 95   Hi 95
## Jul 249           190148.5 158453.8 221843.1 141675.6 238621.3
```