

Part 1-a

Time series analysis is all about, doing predictive analytics, predicting future values of variable or forecasting certain amount along a particular direction.

We will see :-

What is Time series

Regression vs Timeseries

Look for trends & seasonality

Smoothing

ACF

ARIMA

Lagging

Box prience test .

Filters & convolutions

frequency analysis

Multiplicative time series analysis

forecasting analysis

In Time series, Time as a dimension we are trying to forecast.

We are dealing with two variables where x axis is time & y axis can be any other thing like temp, stock price, etc.

We can't apply regression everywhere.

For example Bengaluru Petrol Price, we want to predict price of petrol.

We have to see whether we take whole previous data or just a slice of it for prediction.

Sensex daily closing ~~etc~~ → 1934 → 1994 no increase, the data can be misleading, but lastly it increased 10x.

We are only predicting data on past events, why it is happening is not our concern

Sensex Small cap → The Statistical properties

of this data is not same as sensx.

Since here is greater change, & it

Less volatility more predictable is the data.

Time series Part 1b

Take relevant slice of data in time series analysis for accurate predictions.

Time series is all about trying to separate out trend, seasonal variance from the data as in Airline passengers.

Remove yearly avg, we get or lets say components of the data are removed one by one from the original data and therefore what is left is called residue.

And what we get is quite different from original data.

Therefore Time series analysis all about revealing those pattern. For example Airline & Rainfall dataset

Time series Part 2a

Valid time series data must have.

↳ Min of 50 data points

↳ Equally spaced time intervals

↳ No missing data points.

↳ Stationary: no trend, no seasonality, and no change in frequency or amplitude of noise (i.e. its statistical distribution properties must not change over time)

Frequency and Period.

Frequency : how many times does an event occur in given unit of time

↳ Eg: 4 times in a year.

Period: Inverse of frequency.

→ Eg. Once in 3 months (annual frequency = 4)

Seasonality = periodicity (period is 1 year → seasonality otherwise periodicity)

→ Cyclic variations in repeated patterns.

Time domain : at what points in time did the event occur and at what amplitude.

Ex: SW monsoon occurs in June - Sept.

Frequency domain : at what frequencies does the event repeat.

Ex: SW monsoon as a periodicity of 1 per year

If we remove seasonal variance in atmospheric CO₂ data, we get a trend.

for non linear trend, we can use log of data for analyzing & fit regression line rather than doing Time series analysis.

Time series part 2b

property

Statistical distribution of stock data is not same for 1970 - 1990 as in 1990 to 2010; First there is an increase & then a decrease after peak. So there is non stationarity in data.

Therefore we need to split data

For Time-Series analysis

We decompose the data into:

- Trend : increase or decrease over long term
- Seasonality : periodic changes
- Noise or irregularity
- Other
- Cyclicality - (beyond 1 year)
- Other.

Time - Series operations:

1. Create Time - Series data

- a) Vector of numeric data and vector of index (date and / or time)

2. Plotting

3. Extracting, subsetting, merging, filling, padding and lagging

Main operations:

- ↳ Computing successive difference : diff
- ↳ Apply daily / weekly / monthly / quarterly / yearly
- ↳ Apply a rolling function : rollapply
- ↳ Moving averages (airpassage → annual avg).
- ↳ Smoothing a time - series
- ↳ acf.
- ↳ Test for autocorrelation: Box test.
- ↳ Partial auto correlation function: pacf
- ↳ Cross-correlation function: ccf
- ↳ Detrending
- ↳ Fitting an arima model: auto.arima
- ↳ Running diagnostics on an ARIMA model: tsdiag
- ↳ Making forecasting using ARIMA: predict
- ↳ Testing for mean reversion: adf.test

Part 3a

data(Airpassenger)

library(timetk)

library(forecast)

class(Airpassenger) \rightarrow ts.

This dataset has 12 data point per year.

actual trend can be find out by

$d \leftarrow \text{stl}(\text{Airpassenger}, s.window = 12)$

Take avg. of data points

plot(aggregate(Airpassenger, FUN = mean))

boxplot(Airpassenger)

similarly for cycles

Part 3b

In this CO₂ had linear trend. Seasonal variance is perfect.

Now we use rainfall dataset & create ts object

$fr \leftarrow \text{fseq}(\text{from} = \text{as.Date}("1871-01-01"), \text{by} = \text{"month"}$
 $\text{length.out} = 1692)$

$tsr \leftarrow \text{as.ts}(rainfall, fr)$

Rainfall data is too complex to see seasonality and trends.

Part 3c

head(t) \rightarrow oldest values

tail(t) \rightarrow newest values

- ts::first(t, "month")

- ts::last(t, "2 weeks")

ts[i] ts[j, i]

Complex

ts(as.Date("1947-06-15"))

date on particular date

dates \leftarrow seq(startdate, enddate, increment)

window(t, start = startdate, end = enddate)

merge(t1, t2)

Union of two time series with NA inserts.

merge(t1, t2, all = FALSE)

Gives intersection of two time series.

Pad a time series by merging with an empty ts object with desired ones.

lagging \rightarrow shift ts points by k units lags(t, k).
k positive, k negative.

Successive differences \rightarrow diff(n)

yearly diff = diff(n, lag=12)

rollapply(t, width, f, align = "right")

To apply at every nth point, by = n

Detrending

for linear trend use regression to fit line.

Subtract trend from data values

use decomposition function like stl

Part 4a

Smoothing

avg. of local values is called moving avg. or Rolling mean.

disadvantage: creates mountain for every spike.
Distorts shape of curve

Exponential smoothing

$$S_i = \alpha x_i + (1-\alpha) S_{i-1} \text{ prev}$$

→ previous datapoint.

$$\begin{aligned} &= \alpha x_i + (1-\alpha) (\alpha x_{i-1} + (1-\alpha) S_{i-2}) \\ &= \alpha x_i + (1-\alpha) (\alpha x_{i-1} + (1-\alpha) (\alpha x_{i-2} + (1-\alpha) S_{i-3})) \\ &= \alpha [x_i + (1-\alpha)x_{i-1} + (1-\alpha)^2 x_{i-2} + (1-\alpha)^3 S_{i-3}] \\ &= \dots \\ &= \alpha \sum_{j=0}^i (1-\alpha)^j x_{i-j} \end{aligned}$$

Single exponential smoothing works only if there is no trend,
if it is so values lag behind

Double Exponential Smoothing

$$S_i = \alpha x_i + (1-\alpha) (S_{i-1} + f_{i-1}) \quad f_i \rightarrow \text{trend}$$

$$f_i = \beta (S_i - S_{i-1}) + (1-\beta) f_{i-1}$$

→ diff b/w smooth curval and
smooth prev val

Triple smoothing

when trend & seasonality both present

for additive case:

$$S_i = \alpha (x_i - p_{i-k}) + (1-\alpha) (f_{i-1} + t_{i-1})$$

$$f_i = \beta (S_i - S_{i-1}) + (1-\beta) f_{i-1}$$

$$p_i = \gamma (x_i - S_i) + (1-\gamma) p_{i-k} \quad p_{i-k} \rightarrow \text{curr value actual}$$

$$x_{i+k} = S_i + f_i + t_i + p_{i-k} \quad \begin{array}{l} \text{→ smooth} \\ \text{periodicity value} \\ \text{→ points before} \end{array}$$

Multiplicative

$$s_i = \alpha x_i + (1-\alpha)(s_{i-1} + t_{i-1})$$

p_{i-K}

$$t_i = \beta (s_i - s_{i-1}) + (1-\beta) t_{i-1}$$

$$p_i = \gamma \frac{x_i}{s_i} + (1-\gamma)p_{i-K}$$

$$x_{ith} = (s_i + t_i p_{i-K}) p_{i-K+h}$$

p is periodic component and K is length of period.

all of this

is packaged into Holt winters Model.

Method

↳ exponential weighting

↳ model level: single

↳ Trend level: double

↳ Seasonality level: triple

Holt Winters (ts , alpha = 0.2, beta = 0.2, gamma = 0.2,

seasonal = c("additive", "multiplicative")

Single exponential: beta = NULL, gamma = NULL

Double exponential: gamma = NULL

Specify above values to good fit b/w smooth curve

& actual ~~obs~~ curve

Time Series Modelling

↳ A classic Time Series Model

$$Y = T + S + C + I$$

$Y = \text{data}$

$T = \text{trend}$.

$S = \text{seasonal}$ (regular pattern occur with time)

(\rightarrow long term patterns like business cycle)

$I = \text{irregular}$.

$\Rightarrow \rightarrow$ additive model

Multiplication can be used if seasonality is multiplicative

Part 4b

data ("Airpassenger")

since in seasonal plot, seasonality was amplifying over time we choose multiplication

$h_1 \leftarrow \text{HoltWinters}(\text{Airpassenger}, \text{alpha}=0.5, \text{beta}=0.3, \text{gamma}=0.2, \text{season}=\text{"multiplicative"})$

Try for various values of alpha, beta & gamma

got good fit since of single value exponential smoothing

for prediction

$\text{predict}(h_1, \text{n.ahead}=12)$

if you have data points say 100, then
you shouldn't predict more than 50 for future case.

Part 5a

Auto correlation

$$c(k) = \frac{\sum_{i=1}^{N-k} (x_i - \bar{x})(x_{i+k} - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2}$$

$$\text{with } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

diff with \bar{x} & correlating with k values ahead

repeat for K times, and it's actually a K
then autocorrelation length.

Box-T胥 → Box-T胥 (x) → return p value

Time Series Regression Model

Linear Regression: $Y_i = \beta X_i + \varepsilon_i$

linear trend + white noise

self-blown & Autoregression (AR): $X_t = \phi X_{t-1} + \varepsilon_t$

Moving Average (MA): $\varepsilon_t = u_t + \theta u_{t-1}$

Combine two - ARIMA $X_t = \phi X_{t-1} + u_t + \theta u_{t-1}$

Differencing: $X_t - X_{t-1}$

Integrated → ARIMA: initial diff. steps non-stationary

Assumptions

Noise terms are

↳ Independent

↳ Identically distributed

↳ Normally distributed

↳ Having zero mean

ARIMA

ARIMA → Auto-Regressive Moving Average

ARIMA → Generalization of ARIMA

Auto-regressive Integrated Moving Average.

Arima if non stationary, initial diff. steps are
integrated to reduce non-stationarity

arima (p, d, q)

p = order of AR

d = diff. order

(# of time lags)
(# moving avg.) (coeff)

d = degree of differencing (#times past values are subtracted)

for ex)

P → ~~for~~ 2, find two values of k
two different periodic components

q → Should I do monthly, annual
aug. etc.

arima is popular, because it can be automated

Part 5b

Non seasonal ARIMA

ARIMA(p, d, q)

p → order of auto-regression

d → degree of diff.

q → order of moving avg.

Seasonal arima.

↳ ARIMA(p, d, q) (P, D, Q) m

$m \rightarrow$ no. of periods in each season

↳ $P, D, Q \rightarrow$ for seasonal component

ARIMA(0, 0, 0) $X_t = \epsilon_t$ - White Noise

AR(1) = ARIMA(1, 0, 0)

Random walk = $I(1) = \text{ARIMA}(0, 1, 0)$

$$- X_t = X_{t-1} + \epsilon_t$$

MA(1) = ARIMA(0, 0, 1)

ARIMA(0, 1, 1) → basic exponential smoothing

ARIMA(0, 2, 2) → double exponential smoothing

ARIMA Method

Box Jenkins approach

To find right values of p, d, q for ARIMA modelling?

- ↳ Model identification (i.e., its order)
- ↳ Parameter estimation → fit model to data.
- ↳ Diagnostic checking to validate model

Pacf → partial autocorrelation function. at lag k is corr. between all data points that are exactly k steps apart.
 → after accounting for their correlation with the data b/w those k steps.

Pacf helps to find degree of AR in ARIMA.

R code

library (forecast)

auto.arima(x)

$m \leftarrow \text{arima}(n, \text{order} = c(p, d, q))$

fdiag(m)

predict(m, n.ahead=10)

Testing for Mean Reversion

- ↳ If a time series data is moving away from its mean, then it has a tendency to come back to its mean.

Test

Augmented Dickey-Fuller test

adf.test(x)

Part-6a

Box-test (CO₂) give p-value

acf (Airpassenger) shows auto-correlat.

acf (CO₂)

acf (fri) both the ^{auto} -ve correlation

$m_1 \leftarrow$ auto.arima (Air Passengers)

since ma going from -ve to the so
we should not consider

Here we checked whether p, q, d values
given ~~Part~~ by auto.arima, really fit the
model or not?

Part 6-b

$\emptyset m_2 \leftarrow$ auto.arima (CO₂)

ARIMA (1,1,1) (1,1,2) [12]

$m_3 \rightarrow$ auto.arima (rainfall)

ARIMA (0,0,0) (3,1,0) [,2]

since -ve to the so is not
a good model.

Taking your own model

$m_4 \leftarrow$ arima (701)

$m_4 \leftarrow$ arima (fri, c(2,0,1))

Part 7

Periodogram

Lomb-Scargle Periodogram

↳ Converting from time domain to freq. domain

↳ ~~Yearly~~ Seeing various long term cycles.

↳ Doing moving average

↳ plotting periodicity

Heteroskedasticity

Non constant residual variance

Volatility clusters in financial data

ARCH and GARCH methods

ARCH : Auto regressive conditional heteroskedasticity

GARCH : Generalized ARCH

Applications in Financial Modelling & Economics

Graph Analytics

(graphs & graph theory)

$$G = (V, E)$$

V = {vertex or node}

F

Directed & undirected graphs

Paths & Cycles

Connected components & cliques

Graph algs shortest path etc

Diameter of graph

• Graph centrality

• Degree centrality

• Closeness centrality

• Eigenvector centrality

• Between centrality

Centre & Centrality

Centrality: denotes importance of a node

Centre of graph is point that is at least distance from all nodes.

Degree centrality: Number of connections that a node has

$$cd(v_i) = d_i$$

Between-ness Centrality

↳ Bridges or nodes that span boundaries

↳ proportion of all pairs shortest paths passing through node

Higher network importance = Higher Betweenness

It helps to find key participants in network
those who connect people.

Closeness Centrality: measure of shortest path in a graph.

↳ Inversely proportional to distances in graph

$$cc(v_i) = \frac{1}{\sum_{j} d_{ij}}$$

Diameter of graph: Length of longest shortest path

Eigen vector centrality
↳ finding influential nodes in a network

Ramayana

- ↳ Characters are nodes
- ↳ Relations b/w characters as edges

from text to social graph

Input text + of Ramayana

Step 1 → program to identify characters

Step 2 → co-occurrence of potential relations

Step 3 → build social graph of Ramayana characters

Step 4 → apply graph analysis to graph

Importance of Relation

no. of times of their co-occurrence, is strength

or Importance of their relationship

Centrality Findings

Centre of graph: Rama & Laxman

Centrality Scores:

Rama = 87.5% 56 relations

Sita = 46.88% 36 relations

Lakshma = 43.75% 28 relations

Role of hanuman

Show Hanuman as highest Betweenness centrality

- ↳ Third highest Eigenvector Centrality
- ↳ Fourth in Degree Centrality
- ↳ Fourth in closeness Centrality

Identifying sub-stores

↳ Subgraph or clique detection algorithm to find sub stores