

Time-Series:

- ① Time-series Data: Data points collection, often regular, time intervals.
- Cross-sectional Data: Comparing different groups at a specific moment.
- key components of time series: Trend, Seasonality, Cyclical variation, Irregular variation
- A series is stationary: mean, variance, autocorrelation - remain constant over time
- How do you test for stationary: statistical test: ADF (Augmented Dickey-Fuller) KPSS (Kwiatkowski-Phillips-Schmidt-Shin)

⇒ ADF:

$$\Delta y_t = \alpha + \beta t + \gamma y_{t-1} + \sum_{i=1}^p \delta_i y_{t-i} + \epsilon_t$$

$$\Delta y_t = y_t - y_{t-1}$$

t = time trend

α = lag or drift

$\gamma = 0$ (test static)

- P-value < 0.05 → reject H_0 → series is stationary
- P-value > 0.05 → fail to reject H_0 non-stationary

⇒ KPSS:

$$y_t = \delta_t + \beta t + \epsilon_t$$

δ_t = random walk

β_t = deterministic trend

ϵ_t = stationary error term

- P-value < 0.05 → non-stationary
- P-value > 0.05 → stationary

- When the ADF test fails to reject the unit root null, but the KPSS test rejects its stationarity null. In this case, the series is often trend-stationary.

- A time series is stationary if its mean, variance, autocovariance are constant over time.

a) Differencing: Remove trends & unit roots.

$$y'_t = y_t - y_{t-1}$$

$$y''_t = y'_t - y'_{t-1}$$

b) Log transformation:

$$y'_t = \log(y_t)$$

→ stabilizes variance

c) Square Root/Cube Root:

d) Box-Cox transformation:

$$y'_t \xrightarrow{\lambda \neq 0} \frac{y_t^\lambda - 1}{\lambda}$$

$$\xrightarrow{\lambda = 0} \log(y_t)$$

e) Seasonal Differencing:

$$y'_t = y_t - y_{t-s}$$

f) Standardization/Normalization:

g) Detrending by model fitting:

• stationarity: How a time series is changing will remain the same in the future.

• A unit root is a stochastic trend called a "random walk with drift". Randomness can't be predicted.

Important: • unit root present = not stationary
• unit root absent = stationary

→ To test for stationarity with a unit root test.

assume: Null hypothesis (H_0): no unit root
Alternative hypothesis (H_1): unit root

② White Noise \doteq A sequence of random values with, $(\mu=0) + (\sigma^2 = \text{constant}) + (\text{No autocorrelation})$

In time series modeling, pure randomness.

• Residuals = actual data - model's fitted values.

After fitting ARIMA/ETS/Regressions

If [residuals ~ white noise], then model has captured all systematic patterns. (trend, seasonality, autocorrel)

Trend stationary \rightarrow Differenced stationary series \doteq

Series that becomes stationary once you remove a deterministic trend. (straight line)

$$Y_t = a + bt + \epsilon_t$$

Series that becomes stationary only after differencing.

$$Y_t = Y_{t-1} + \epsilon_t$$

$$\Delta Y_t = Y_t - Y_{t-1} = \epsilon_t$$

Autoregressive model

AR vs MA vs ARMA \doteq

moving avg. \rightarrow autoregressive moving avg.

• Current value of the series depends on the past [value + noise]. (like momentum)

$$Y_t = c + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \epsilon_t$$

Y_t : current value

Y_{t-1}, Y_{t-2} : past value

ϵ_t : white noise

• yesterday's return influences today's return

influences today's return

MA \doteq

• Current value = past shocks (errors)

$$Y_t = c + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

$\epsilon_{t-1}, \epsilon_{t-2}$: past shocks

• Today's return may depend on yesterday's partly.

ARMA \doteq

• Current value = (AR + MA)

$$Y_t = c + \phi_1 Y_{t-1} + \dots + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots$$

Most important

if series is stationary \rightarrow ARMA
if non-stationary \rightarrow ARIMA

ARIMA \doteq (Autoregressive Integrated moving avg.)

• AR(p) + MA(q) + ARMA(p,q) \rightarrow p = no. of autoregressive terms

q = no. of moving avg. terms

d = no. of differences

handle non-stationarity

• ARIMA(1,0,0) \rightarrow AR(1)

" (0,1,0) \rightarrow MA(1)

" (1,1,0) \rightarrow ARMA(1,1)

" (0,1,1) \rightarrow Random walk

" (0,1,1) \rightarrow Exponential smoothing

" (1,1,1) \rightarrow SARIMA

$$Y_t = \phi_1 Y_{t-1} + \epsilon_t + \theta_1 \epsilon_{t-1}$$

SARIMA \doteq (seasonality)

• SARIMA(p,d,q)(P,D,B,s)

seasonal period

- ③ ACF & PACF \div (Partial Autocorrelation Function) :
- Correlation of y_t with $y_{t-1}, y_{t-2} \dots$, it's tell you how past values in general affect the current value.
 - Correlation of y_t with y_{t-2} , after removing y_{t-1} .
 - Before plotting ACF/PACF, you must make the series stationary.
 - PACF helps guess P , ACF helps guess Q .
- \Rightarrow Confidence band \div (statistical threshold for significance)
- The spikes in ACF & PACF plots tells you whether a correlation at a given lag is statistically significant or just noise.

- Bands \div A threshold dashed blue lines.
- How to calculate \div For a large sample size N : std. error of autocorrelation $\approx \frac{1}{\sqrt{N}}$
- All spikes within the band, series behave like white noise.

MAE (mean absolute error) $\rightarrow mae = \frac{\sum_{t=1}^n |y_t - \hat{y}_t|}{n}$

RMSE (root mean squared error) $\rightarrow rmse = \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2}$

MAPE (mean absolute percentage error) $\rightarrow mape = \frac{100}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right|$

AIC (akaike information criterion) $\rightarrow AIC = -2 \ln(L) + 2K$

BIC (bayesian information criterion) \rightarrow Likelihood of the model L , no. of parameters K . Lower AIC \rightarrow better model.

- These are used to measure How well the forecasted values match actuals value.
- How well the model structure fit, penalizing complexity.

- How to handle the missing values in time series data?
- Random missing
Systematic missing
- i) Forward Fill $\div y_t = y_{t-1}$
 - ii) Backward Fill $\div y_t = y_{t+1}$
 - iii) linear interpolation \div Missing points b/w known values.
 - iv) Time-based Interpolation \div Spline/polynomial (use curve line)
 - v) Model-based \div
 - vi) static (mean, median, mode) \div

- How to Backtest a Forecasting model.
- Train-test Split
 -

⑥ GARCH (generalized autoregressive condition Heteroskedasticity) :

- Designed for variance (volatility). where ARIMA designed for mean (Returns)
- today volatility = past volatility + past forecast errors.
- GARCH (P, q)

$$h_t = \alpha_0 + \sum_{i=1}^p \alpha_i e_{t-i}^2 + \sum_{j=1}^q \beta_j h_{t-j}$$

$: h_t = \sigma_t^2$

conditional variance
 (volatility at time t)
 yesterday's & before
 volatility
 yesterday's & before
 squared error

⑦ Exponential Smoothing Models :

- Recent data gets more weight, older data get exponentially less weight.
- Simple ES : $\hat{y}_{t+1} = \alpha y_t + (1-\alpha) \hat{y}_t$ → data has no trend
- Holt's linear trend : level: $l_t = \alpha y_t + (1-\alpha)(l_{t-1} + T_{t-1})$
 Trend: $T_t = \beta(l_t - l_{t-1}) + (1-\beta)T_{t-1}$
 Seasonal: $S_t = \gamma(Y_t - l_t) + (1-\gamma)S_{t-m}$

- Use Forecast risk (VaR, option pricing)
- position size based on expected volatility
- Identify High-risk & low-risk regimes

⑧ Machine Learning Approaches :

- Random Forests & Logistic Regression :
- Support vector Machine (SVM)
- Recurrent Neural Network (RNN) + Reinforcement learning
- Long short-term memory (LSTM)
- Gated Recurrent units (GRU)
- XGBoost classifier

- How to deal with large datasets or high-dimensional feature spaces.
- Optimization Algorithms : Stochastic gradient descent (SGD) (Time consuming + power computation)
 - Bayesian Optimization
 - Parallel Computing : Use Amdahl's law
 - Dimensionality Reduction : PCA (Principal component analysis)
- $$S = \frac{1}{(1-P)+P} \frac{N}{\text{Number of processes}}$$
- speedup program,