

Complete Time Series Analytics

From Zero to Hero: A Journey Through Data Patterns

Amit Dua

June 29, 2025

Course Philosophy: *"The best way to learn time series is to see patterns, understand the mathematics, and apply the knowledge to real problems."*

Contents

1	Welcome to the Amazing World of Time Series!	3
2	Module 1: The Anatomy of Time Series	6
2.1	Time Series Components - Dissecting the Data	6
2.1.1	Understanding Trend - The Long-term Story	7
2.1.2	Seasonal Patterns - The Rhythms of Time	10
2.1.3	Cyclical vs Seasonal - The Subtle Difference	14
2.1.4	Irregular Component - The Random Element	17
2.2	Moving Averages - Smoothing the Noise	18
2.2.1	Simple Moving Average - The Foundation	19
2.2.2	Weighted Moving Average - Giving Recent Data More Importance	23
2.2.3	Exponential Moving Average - The Smart Smoother	24
2.3	Autocorrelation - The Memory of Time Series	27
2.3.1	Understanding Correlation First	27
2.3.2	Autocorrelation Function (ACF) - The Mathematical Foundation	29
2.3.3	Properties of ACF for Different Processes	31
2.3.4	Partial Autocorrelation Function (PACF) - The Direct Relationship	34
2.3.5	Ljung-Box Test - Testing for Autocorrelation	38
3	Module 2: The Mathematics of Time Series Models	39
3.1	Stationarity - The Foundation of Time Series Modeling	39
3.1.1	Understanding Stationarity Mathematically	39
3.1.2	Why Stationarity Matters	42
3.1.3	Testing for Stationarity	45
3.1.4	Making Data Stationary	47
3.2	Linear Processes and Moving Average Models	51
3.2.1	General Linear Process	52
3.2.2	Moving Average Process MA(q)	53
3.2.3	MA(1) Process - Deep Dive	55

3.3	Autoregressive Models AR(p)	56
3.3.1	AR(p) Model Definition	57
3.3.2	AR(1) Process - The Building Block	59
3.3.3	AR(2) Process and Stationarity Conditions	62
3.4	ARMA Models - Combining AR and MA	64
3.4.1	ARMA(p,q) Model Definition	65
3.4.2	ARMA(1,1) - The Workhorse Model	66
3.5	Non-Stationary and Seasonal Models	68
3.5.1	ARIMA Models - Integration and Differencing	68
3.5.2	Seasonal ARIMA - SARIMA Models	71
4	Module 3: Forecasting and Model Selection	74
4.1	Forecasting in Time Series Models	74
4.1.1	Minimum Mean Square Error Forecasting	75
4.1.2	Forecasting with AR Models	76
4.1.3	Forecasting with MA and ARMA Models	78
4.1.4	Prediction Intervals	80
4.2	Durbin-Levinson Algorithm	82
4.3	Parameter Estimation Methods	85
4.3.1	Method of Moments Estimation	86
4.3.2	Least Squares Estimation	87
4.3.3	Maximum Likelihood Estimation	88
4.4	Model Selection and Comparison	90
4.4.1	Information Criteria	91
4.4.2	Model Selection Strategy	94
4.5	Residual Analysis and Diagnostic Checking	95
4.5.1	Properties of Good Residuals	96
4.5.2	Ljung-Box Test for Residuals	97
4.5.3	Comprehensive Diagnostic Strategy	100
4.6	Unit Root Tests	102
4.6.1	Phillips-Perron Test	102
4.6.2	KPSS Test Revisited	103
4.6.3	Testing for Multiple Unit Roots	105
5	Module 4: Advanced Time Series Models	106
5.1	Multivariate Time Series and VAR Models	106
5.1.1	Introduction to Multivariate Time Series	106
5.1.2	Vector Autoregression (VAR) Models	108
5.1.3	Stability and Stationarity of VAR	110
5.1.4	Impulse Response Functions	111
5.1.5	Granger Causality	114
5.2	Vector ARMA Models	116
5.3	Conditional Heteroscedastic Models	117
5.3.1	ARCH Models	118
5.3.2	GARCH Models	121
5.3.3	GARCH Model Estimation	124
5.3.4	GARCH Forecasting	125
6	Course Summary and Integration	128

1 Welcome to the Amazing World of Time Series!

Professor

Welcome, dear students! Today we embark on one of the most exciting journeys in data science. You might be wondering: "What makes time series so special?" Imagine you're a detective, and the data points are clues scattered across time. Unlike regular data where observations are independent snapshots, time series data tells a **story** - each observation is connected to what came before and influences what comes next.

Think about your daily life: Your mood today might depend on yesterday's events. Stock prices today are influenced by yesterday's news. Weather tomorrow is connected to today's atmospheric conditions. This is the essence of time series - **temporal dependence**.

Today, we'll learn to read these stories, understand their patterns, and even predict their future chapters!

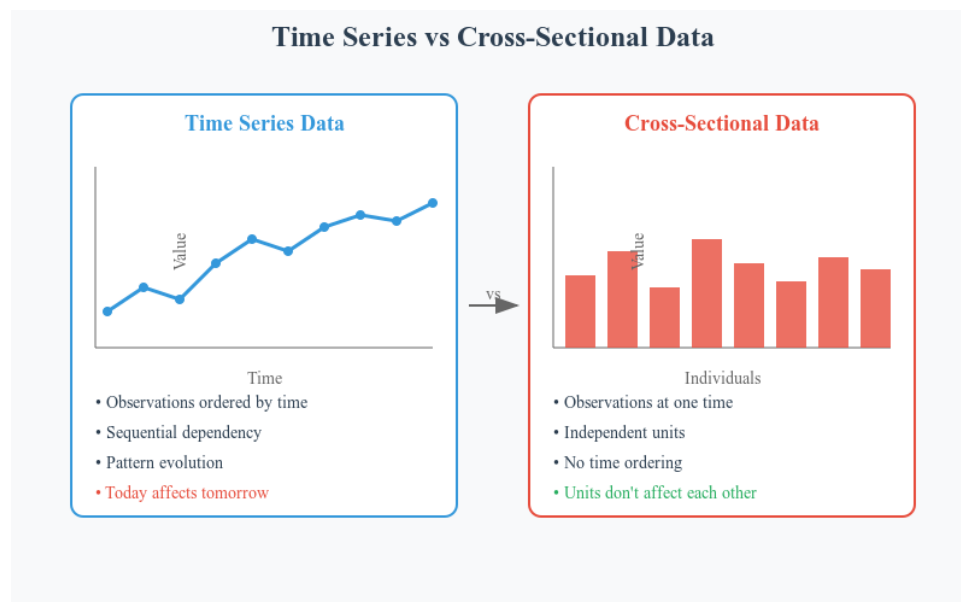


Figure 1: Time Series vs Cross-Sectional Data - The Key Difference

Paul (The Innocent)

Professor, I'm completely new to this. Can you explain what exactly is a "time series"? Is it just data with dates?

Professor

Excellent starting question, Paul! Let me explain with a simple analogy. Imagine you're keeping a diary where you write your happiness level (1-10) every day for a year. That's a time series! It's not just "data with dates" - it's data where the **order** and **timing** matter tremendously.

Mathematical Definition:

$$Y_t = \text{Value observed at time } t$$

Where t represents time (could be seconds, days, months, years).

The magic happens because:

$$Y_t \text{ depends on } Y_{t-1}, Y_{t-2}, Y_{t-3}, \dots$$

This means today's value is influenced by previous values. This is called **autocorrelation** - the correlation of a series with itself at different time lags.

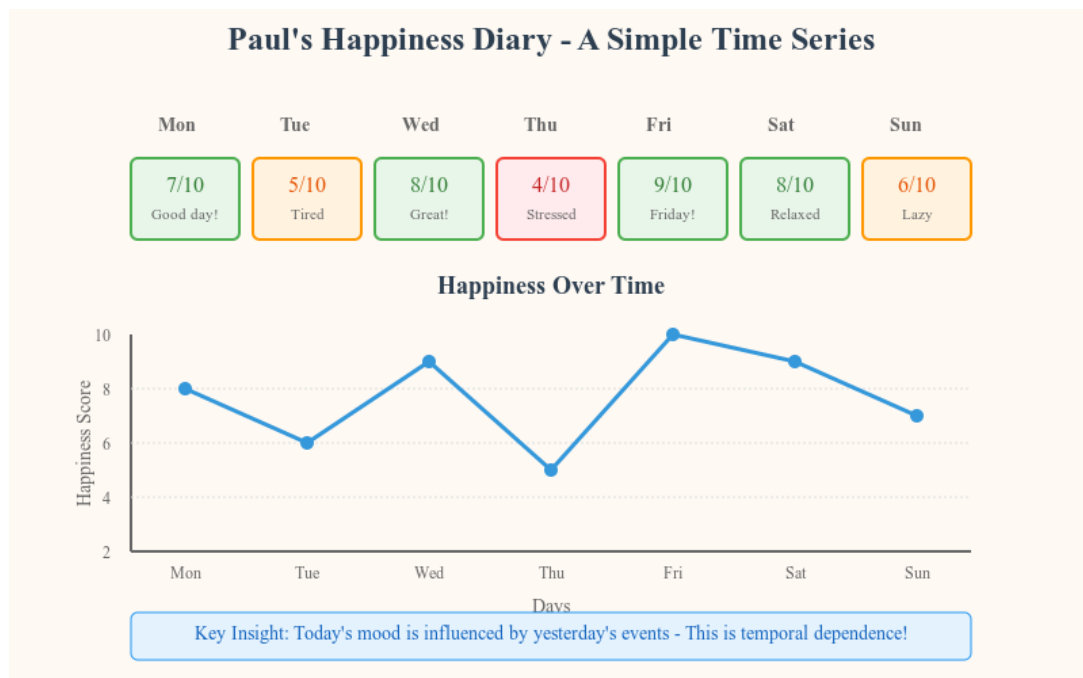


Figure 2: Paul's Happiness Diary - A Simple Time Series Example

Theory Deep Dive

Fundamental Difference: Cross-Sectional vs Time Series Data

Cross-Sectional Data:

- Heights of 100 students measured today
- Salaries of employees in a company
- Test scores of students in an exam

Key property: Observations are independent. John's height doesn't affect Mary's height.

Time Series Data:

- Daily stock prices of MRF for 5 years
- Monthly rainfall in Mumbai
- Weekly COVID cases in India

Key property: Observations are dependent. Today's stock price influences tomorrow's price.

Mathematical Representation:

Cross-sectional: $Y_i \perp Y_j$ (independent)

Time series: $Y_t \not\perp Y_{t-k}$ (dependent)

Brain Teaser

Quick Understanding Check!

Which of these are time series? Mark True (T) or False (F):

1. Daily temperature in Delhi for 2023 ____
2. Heights of students in your class ____
3. Monthly sales of iPhones in India ____
4. Exam scores of 500 students ____
5. Hourly heart rate during exercise ____
6. Age of employees in a company ____

Answers: T, F, T, F, T, F

Think: What makes the "True" ones special? They all have the element of **time dependency**!

2 Module 1: The Anatomy of Time Series

2.1 Time Series Components - Dissecting the Data

Professor

Just like a doctor studies human anatomy to understand how the body works, we need to understand the "anatomy" of time series. Every time series is made up of fundamental components, just like the human body has organs, muscles, and bones.

The classical decomposition says:

$$Y_t = T_t + S_t + C_t + I_t$$

Where each component tells a different part of the story:

- T_t = **Trend** (The backbone - long-term direction)
- S_t = **Seasonal** (The heartbeat - regular patterns)
- C_t = **Cyclical** (The breath - irregular long-term cycles)
- I_t = **Irregular** (The reflexes - random movements)

Let's explore each component in detail!

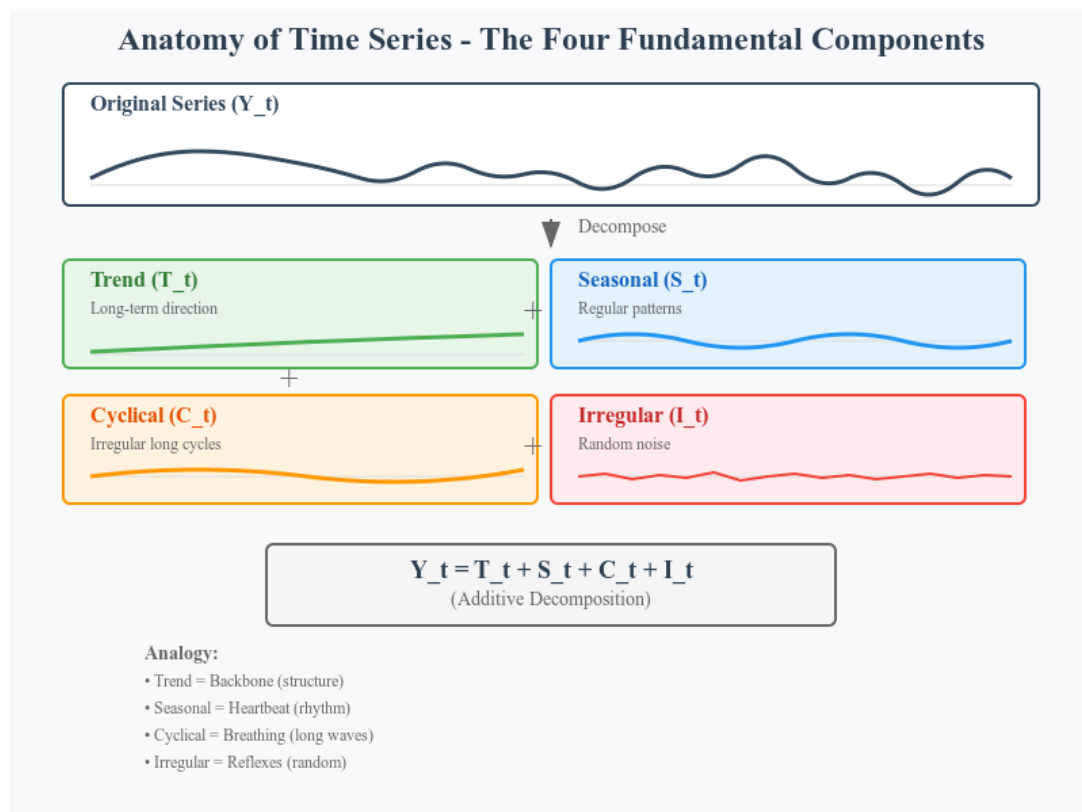


Figure 3: Anatomy of Time Series - The Four Fundamental Components

2.1.1 Understanding Trend - The Long-term Story

Theory Deep Dive

Trend Component (T_t)

The trend represents the long-term movement in the data. It answers the question: "Where is this series going over time?"

Mathematical Definition:

$$T_t = \alpha + \beta t + \epsilon_t$$

Where:

- α = Starting level (intercept)
- β = Rate of change (slope)
- t = Time
- ϵ_t = Small random variations

Types of Trends:

1. **Linear Trend:** $T_t = \alpha + \beta t$ (straight line)
2. **Exponential Trend:** $T_t = \alpha e^{\beta t}$ (curved growth)
3. **Polynomial Trend:** $T_t = \alpha + \beta_1 t + \beta_2 t^2$ (curved)
4. **No Trend:** $T_t = \alpha$ (horizontal line)

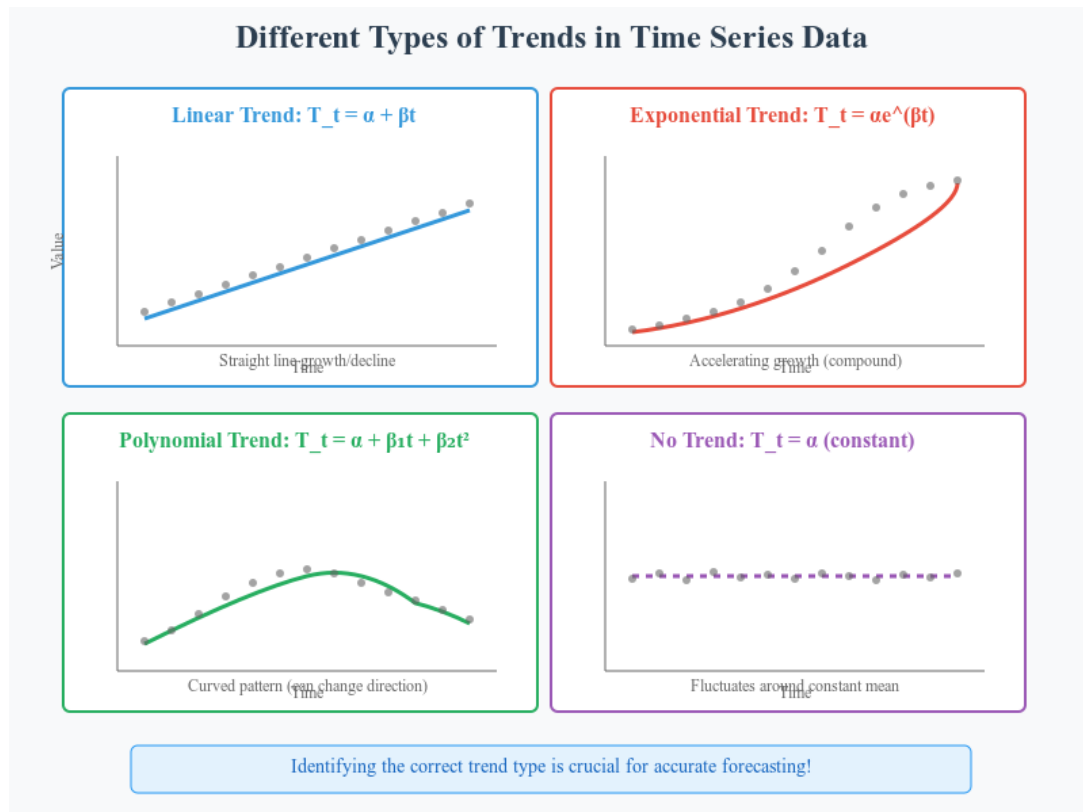


Figure 4: Different Types of Trends in Time Series Data

Rohan (The Visual Learner)

Professor, I can see the different trend patterns in the diagram, but how do I know which type of trend my data has? And why does it matter for analysis?

Professor

Great visual question, Rohan! Identifying the trend type is crucial because it determines how we should model and forecast the data.

How to Identify Trend Type:

1. **Visual Inspection:** Plot the data and look at the shape
2. **Linear Regression:** If R^2 is high, likely linear trend
3. **Log Transformation:** If $\log(\text{data})$ shows linear trend, original has exponential trend
4. **Polynomial Fitting:** Try different polynomial degrees

Why It Matters:

- **Forecasting:** Wrong trend type = Wrong predictions
- **Stationarity:** Trends make data non-stationary
- **Model Selection:** Different trends need different models

Let me show you a real example with MRF stock prices!

Real-World Example

Real Example: MRF Stock Price Trend (2019-2024)

MRF (Madras Rubber Factory) stock shows a clear exponential trend:

- 2019: 60,000 per share
- 2024: 1,20,000 per share
- Pattern: Accelerating growth (not linear!)

If we assume linear trend: Forecast for 2025 = 1,32,000
If we assume exponential trend: Forecast for 2025 = 1,45,000

The difference of 13,000 per share could mean millions in investment decisions!

Python Code Reference

Python Code Reference: 01.time_series_decomposition.py

Key functions to understand:

```
# Load stock data
data = yf.download('MRF.NS', start='2019-01-01')

# Decompose into components
decomposition = seasonal_decompose(data['Close'],
model='multiplicative')

# Plot all components
decomposition.plot()
```

Run this code to see the actual decomposition of MRF stock!

2.1.2 Seasonal Patterns - The Rhythms of Time

Theory Deep Dive

Seasonality Component (S_t)

Seasonality represents regular, predictable patterns that repeat over fixed periods.

Mathematical Definition:

$$S_t = S_{t-s} \text{ for all } t$$

Where s is the seasonal period (e.g., 12 months, 4 quarters, 7 days).

Key Properties:

1. **Fixed Period:** Pattern repeats every s time units
2. **Predictable:** Same pattern occurs at same time each cycle
3. **Additive or Multiplicative:**
 - Additive: $Y_t = T_t + S_t + I_t$ (constant seasonal effect)
 - Multiplicative: $Y_t = T_t \times S_t \times I_t$ (proportional seasonal effect)

Common Seasonal Patterns:

- Daily: Hour-of-day patterns (traffic, website visits)
- Weekly: Day-of-week patterns (restaurant sales)
- Monthly: Month-of-year patterns (ice cream sales)
- Quarterly: Quarter-of-year patterns (retail sales)

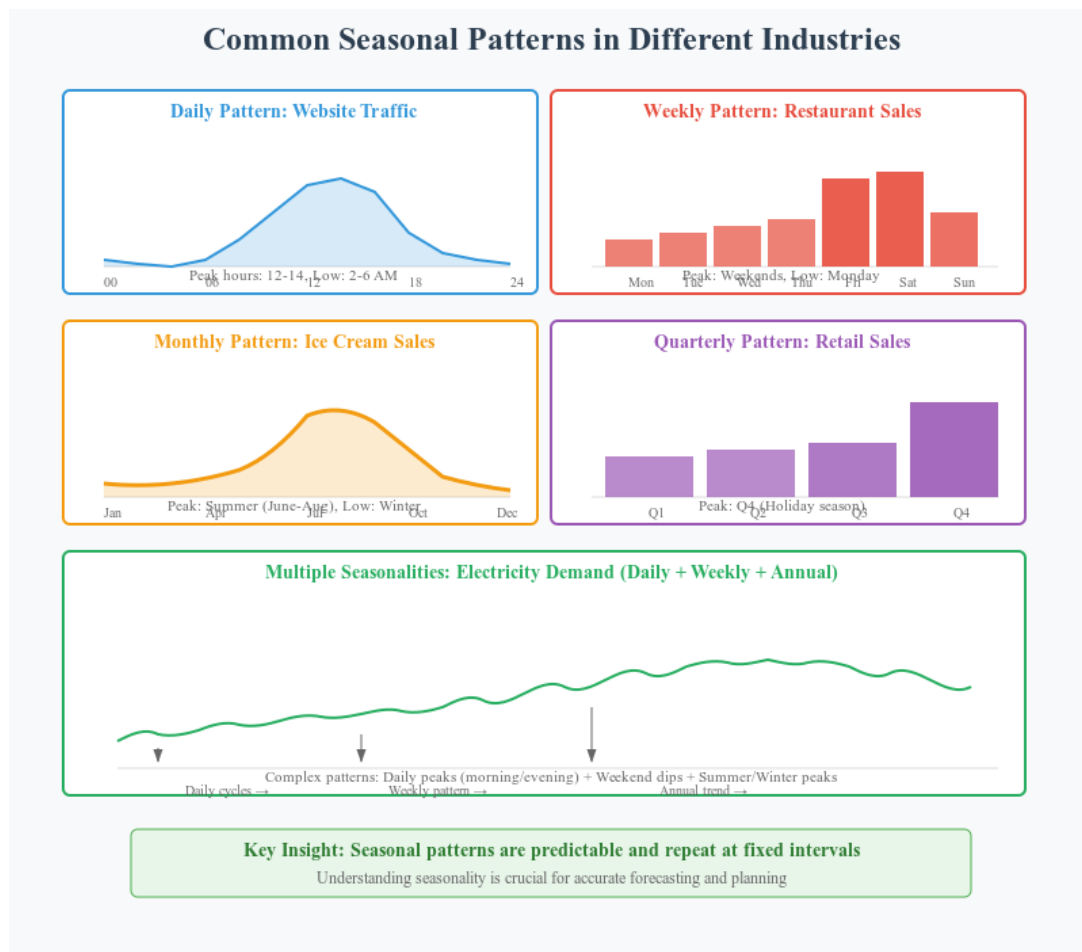


Figure 5: Common Seasonal Patterns in Different Industries

Sumit (The Inquisitive)

Professor, I understand that seasonality repeats, but how do we mathematically distinguish between additive and multiplicative seasonality? And when would we use each model?

Professor

Excellent analytical question, Sumit! The choice between additive and multiplicative seasonality depends on how the seasonal effect changes with the trend level.


Additive Model: $Y_t = T_t + S_t + I_t$

- Seasonal effect is **constant** over time
- Example: Temperature always varies by $\pm 10^\circ\text{C}$ seasonally, regardless of climate change trend
- Use when: Seasonal fluctuations don't change with trend level

Multiplicative Model: $Y_t = T_t \times S_t \times I_t$

- Seasonal effect is **proportional** to trend level
- Example: Sales increase 20% every December, whether baseline is 1 lakh or 10 lakh
- Use when: Seasonal fluctuations grow/shrink with trend

Mathematical Test: Plot the data. If seasonal amplitude increases with trend \rightarrow Multiplicative If seasonal amplitude stays constant \rightarrow Additive



fig_additive_vs_multiplicative.png

Figure 6: Additive vs Multiplicative Seasonality - Visual Comparison

Brain Teaser

Seasonality Detective Challenge!

Look at these scenarios and identify the type of seasonality:

1. Ice cream sales: Summer sales are always 50,000 higher than winter, regardless of store size

Ⓐ Additive
Ⓑ Multiplicative
2. E-commerce website traffic: 30% spike every festival season, proportional to current user base

Ⓐ Additive
Ⓑ Multiplicative
3. Daily temperature: Varies by $\pm 15^{\circ}\text{C}$ seasonally in Delhi, independent of global warming trend

Ⓐ Additive
Ⓑ Multiplicative

Answers: A, B, A

Key Insight: When seasonal effect is a fixed amount \rightarrow Additive. When it's a percentage \rightarrow Multiplicative.

2.1.3 Cyclical vs Seasonal - The Subtle Difference

Professor

Many students confuse cyclical and seasonal patterns. Let me clarify this important distinction!

Seasonal Patterns:

- **Fixed duration:** Always same length (12 months, 7 days)
- **Regular timing:** Occurs at same time each cycle
- **Predictable:** Can forecast exactly when pattern occurs

Cyclical Patterns:

- **Variable duration:** Can last 2-10 years (business cycles)
- **Irregular timing:** No fixed schedule
- **Unpredictable:** Cannot forecast exact timing

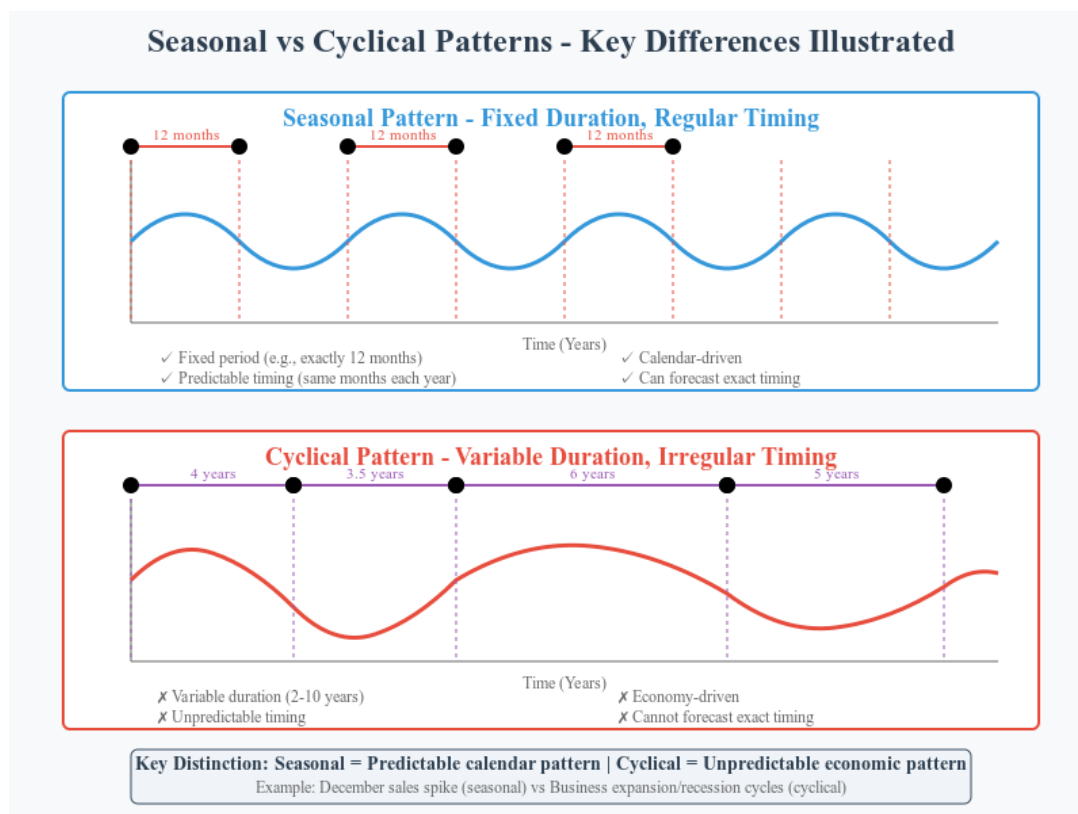


Figure 7: Seasonal vs Cyclical Patterns - Key Differences Illustrated

Neha (The Skeptic)

Professor, I'm skeptical about this distinction. In real data, how can we be sure what we're seeing is truly cyclical and not just irregular seasonal patterns? Couldn't what we call "business cycles" just be long-term seasonal effects we don't understand yet?

Professor

Brilliant skepticism, Neha! You've touched on a deep philosophical question in time series analysis. The truth is, the distinction isn't always clear-cut in practice.

Practical Identification Criteria:

Evidence for Seasonality:

- Pattern repeats at exact intervals (365 days, 52 weeks)
- Strong correlation with calendar events
- Pattern persists across different economic conditions

Evidence for Cycles:

- Duration varies (2-year recession, 8-year expansion)
- Amplitude varies (mild vs severe recessions)
- Correlation with economic indicators, not calendar

Gray Areas: Some patterns blur the lines. For example:

- El Niño cycles: 2-7 year irregular "seasonal" weather patterns
- Real estate cycles: 18-year cycles with seasonal overlay

The key is: If you can predict the *timing* precisely, it's seasonal. If you can only predict the *existence* of the pattern, it's cyclical.

2.1.4 Irregular Component - The Random Element

Theory Deep Dive

Irregular Component (I_t)

The irregular (or random) component represents the "leftover" variation after removing trend, seasonal, and cyclical patterns.

Mathematical Definition:

$$I_t = Y_t - T_t - S_t - C_t \text{ (additive model)}$$

$$I_t = \frac{Y_t}{T_t \times S_t \times C_t} \text{ (multiplicative model)}$$

Properties of Good Irregular Component:

1. **Zero Mean:** $E[I_t] = 0$
2. **Constant Variance:** $Var[I_t] = \sigma^2$
3. **No Autocorrelation:** $Cov[I_t, I_{t-k}] = 0$ for all $k \neq 0$
4. **Normal Distribution:** $I_t \sim N(0, \sigma^2)$ (ideally)

Sources of Irregular Variation:

- Measurement errors
- Unexpected events (natural disasters, policy changes)
- Market sentiment and psychology
- Random economic shocks

Important Warning

Decomposition Quality Check

If your irregular component shows patterns, your decomposition is incomplete! Common issues:

- **Trending residuals:** Missed a trend component
- **Seasonal residuals:** Wrong seasonal period or missing multiple seasonalities
- **Autocorrelated residuals:** Missing cyclical or other systematic patterns

Always check residuals before accepting your decomposition!

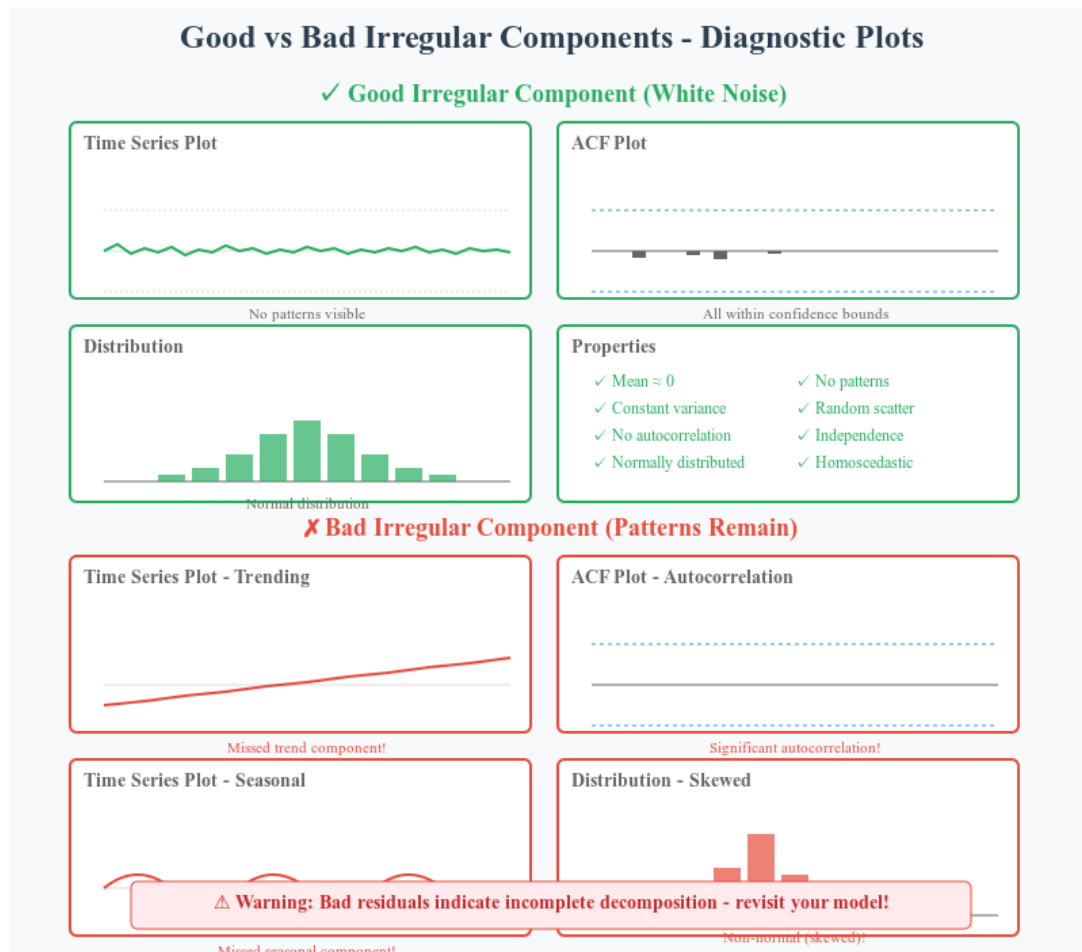


Figure 8: Good vs Bad Irregular Components - Diagnostic Plots

2.2 Moving Averages - Smoothing the Noise

Professor

Now that we understand time series components, let's learn our first tool for analysis: moving averages. Think of moving averages as a "smoothing brush" that helps us see the underlying patterns by reducing noise.

Imagine you're trying to see the shape of a mountain through thick fog. Moving averages are like waiting for the fog to clear in patches - they help reveal the underlying structure.

2.2.1 Simple Moving Average - The Foundation

Theory Deep Dive

Simple Moving Average (SMA)

The simple moving average is the foundation of time series smoothing.

Mathematical Definition:

$$SMA_t^{(k)} = \frac{1}{k} \sum_{i=0}^{k-1} Y_{t-i}$$

Where:

- k = window size (number of periods to average)
- Y_{t-i} = observation i periods ago
- $SMA_t^{(k)}$ = moving average at time t with window k

Why Each Part Matters:

- $\frac{1}{k}$: Ensures we get an average (not just a sum)
- \sum : Combines multiple observations for stability
- Y_{t-i} : Uses recent history to predict current
- Window k : Controls smoothness vs responsiveness trade-off

Properties:

1. **Lag**: MA always lags behind actual data by $(k - 1)/2$ periods
2. **Smoothing**: Larger k = smoother but less responsive
3. **Equal Weights**: All observations in window have same weight $1/k$

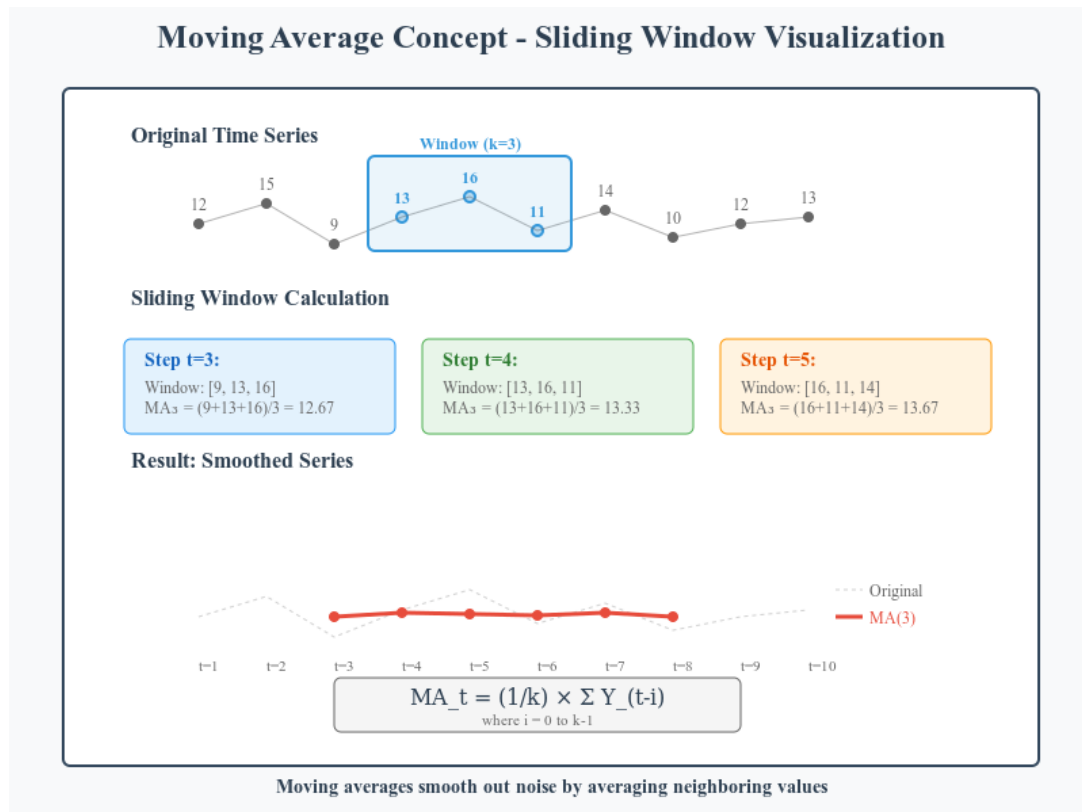


Figure 9: Moving Average Concept - Sliding Window Visualization

Paul (The Innocent)

Professor, I understand the formula, but I'm confused about the window size. How do I choose the right value of k ? If I make it too big or too small, what happens?

Professor

Excellent practical question, Paul! Choosing the window size is like tuning a radio - too high or too low and you miss the signal.

Window Size Effects:

Small Window ($k = 3-5$):

- Responsive to changes
- Captures short-term patterns
- More noise remains
- Less smoothing

Large Window ($k = 20-50$):

- Very smooth
- Removes most noise
- Slow to respond to changes
- Misses short-term patterns

Choosing Window Size:

1. **Data frequency:** Daily data \rightarrow larger windows (20-50), Monthly data \rightarrow smaller windows (3-12)
2. **Purpose:** Trend identification \rightarrow larger, Signal detection \rightarrow smaller
3. **Seasonality:** Match seasonal period (12 for monthly data, 252 for daily stock data)

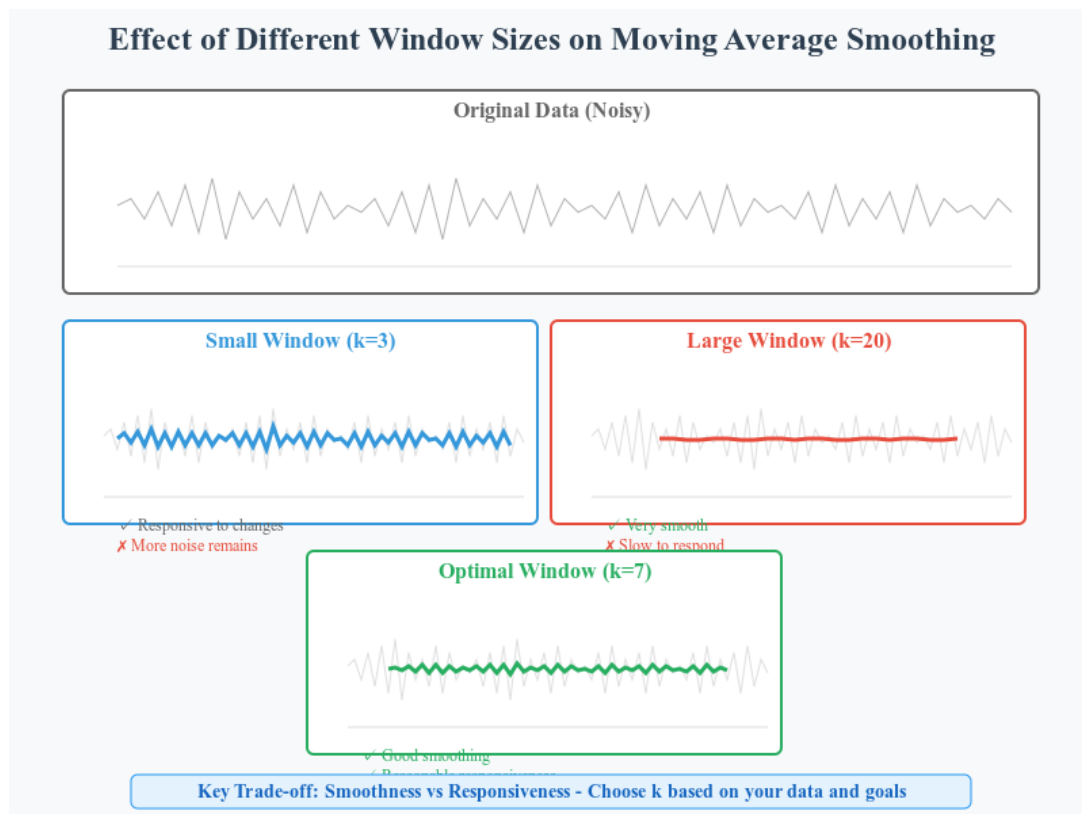


Figure 10: Effect of Different Window Sizes on Moving Average Smoothing

Python Code Reference

Python Code Reference: 01.time.series.decomposition.py

Moving average calculation:

```
# Simple moving averages
data['MA_5'] = data['Close'].rolling(window=5).mean()
data['MA_20'] = data['Close'].rolling(window=20).mean()
data['MA_50'] = data['Close'].rolling(window=50).mean()

# Plot comparisons
plt.plot(data.index, data['Close'], label='Original')
plt.plot(data.index, data['MA_20'], label='MA(20)')
```

This shows how different window sizes affect smoothing!

2.2.2 Weighted Moving Average - Giving Recent Data More Importance

Theory Deep Dive

Weighted Moving Average (WMA)

Sometimes we want recent observations to have more influence than older ones.

Mathematical Definition:

$$WMA_t = \frac{\sum_{i=0}^{k-1} w_i Y_{t-i}}{\sum_{i=0}^{k-1} w_i}$$

Where w_i are weights assigned to each observation.

Common Weight Schemes:

1. **Linear:** $w_i = k - i$ (recent observations get higher weights)
2. **Exponential:** $w_i = \alpha^i$ where $0 < \alpha < 1$
3. **Custom:** Assign weights based on domain knowledge

Example - Linear WMA with k=4:

$$WMA_t = \frac{4Y_t + 3Y_{t-1} + 2Y_{t-2} + 1Y_{t-3}}{4 + 3 + 2 + 1} = \frac{4Y_t + 3Y_{t-1} + 2Y_{t-2} + Y_{t-3}}{10}$$

2.2.3 Exponential Moving Average - The Smart Smoother

Theory Deep Dive

Exponential Moving Average (EMA)

EMA is the most elegant moving average - it uses *all* historical data but gives exponentially decreasing weights to older observations.

Mathematical Definition:

$$EMA_t = \alpha Y_t + (1 - \alpha)EMA_{t-1}$$

Where $0 < \alpha < 1$ is the smoothing parameter.

Recursive Form Expansion:

$$EMA_t = \alpha Y_t + (1 - \alpha)EMA_{t-1} \quad (1)$$

$$= \alpha Y_t + \alpha(1 - \alpha)Y_{t-1} + (1 - \alpha)^2 EMA_{t-2} \quad (2)$$

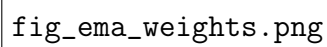
$$= \alpha \sum_{i=0}^{\infty} (1 - \alpha)^i Y_{t-i} \quad (3)$$

Weight Pattern:

- Current observation: weight = α
- Previous observation: weight = $\alpha(1 - \alpha)$
- Two periods ago: weight = $\alpha(1 - \alpha)^2$
- And so on...

Choosing Alpha:

- $\alpha = 0.1$: Heavy smoothing, slow response
- $\alpha = 0.3$: Moderate smoothing
- $\alpha = 0.9$: Light smoothing, fast response



fig_ema_weights.png

Figure 11: Exponential Weight Decay in EMA - Why Recent Data Matters More

Rohan (The Visual Learner)

Professor, I can see from the diagram how the weights decrease exponentially, but why is this better than simple moving average? And how do I choose the right alpha value?

Professor

Great visual question, Rohan! EMA has several advantages over SMA:

Why EMA is Better:

1. **Uses all data:** SMA throws away older data, EMA uses everything
2. **No sudden jumps:** When old data "falls off" SMA window, values can jump. EMA is smooth
3. **More responsive:** Reacts faster to recent changes
4. **Memory efficient:** Only need to store previous EMA value, not entire window

Choosing Alpha:

- **Start with:** $\alpha = \frac{2}{k+1}$ where k is equivalent SMA window
- **High volatility data:** Lower alpha (more smoothing)
- **Trending data:** Higher alpha (more responsiveness)
- **Forecasting:** Optimize alpha to minimize forecast errors

Rule of Thumb: EMA with $\alpha = 0.2$ SMA with window = 9
EMA with $\alpha = 0.1$ SMA with window = 19

Brain Teaser

Moving Average Mastery Challenge!

Given this stock price data: [100, 102, 98, 105, 103, 107, 101]

Calculate by hand:

1. 3-period SMA for the last value: _____
2. 3-period WMA with weights [3,2,1] for the last value: _____
3. If EMA yesterday was 104 and $\alpha = 0.3$, what's today's EMA if today's price is 101? _____

Solutions:

1. $SMA = (107+103+101)/3 = 103.67$
2. $WMA = (3 \times 101 + 2 \times 107 + 1 \times 103)/(3+2+1) = 517/6 = 86.17$
3. $EMA = 0.3 \times 101 + 0.7 \times 104 = 30.3 + 72.8 = 103.1$

2.3 Autocorrelation - The Memory of Time Series

Professor

We now come to one of the most important concepts in time series: autocorrelation. This is what makes time series different from regular data - the ability of a series to be correlated with itself at different time lags.

Think of autocorrelation as the "memory" of time series. Just like you might have similar moods on consecutive days, time series often have values that are similar to their recent past. Understanding this memory is key to prediction and modeling.

2.3.1 Understanding Correlation First

Theory Deep Dive

Correlation Refresher

Before understanding autocorrelation, let's refresh correlation between two different variables X and Y.

Correlation Coefficient:

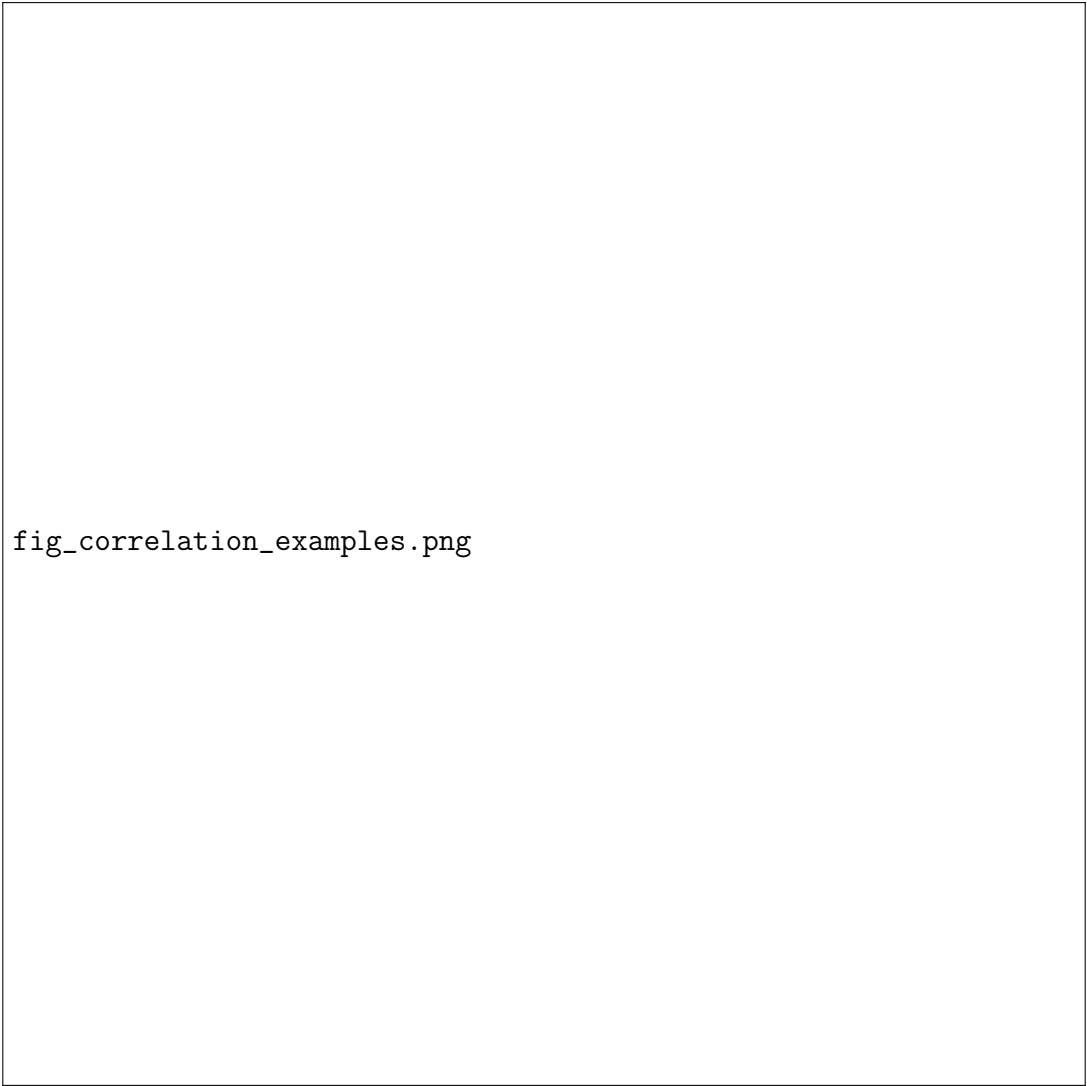
$$\rho_{X,Y} = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

Properties:

- Range: $-1 \leq \rho_{X,Y} \leq 1$
- $\rho = 1$: Perfect positive correlation
- $\rho = 0$: No linear correlation
- $\rho = -1$: Perfect negative correlation

Sample Correlation:

$$r_{X,Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$



`fig_correlation_examples.png`

Figure 12: Correlation Examples - From Perfect Negative to Perfect Positive

2.3.2 Autocorrelation Function (ACF) - The Mathematical Foundation

Theory Deep Dive

Autocorrelation Function (ACF)

Autocorrelation is simply the correlation of a time series with a lagged version of itself.

Mathematical Definition:

$$\rho_k = \frac{\gamma_k}{\gamma_0} = \frac{Cov(Y_t, Y_{t-k})}{Var(Y_t)}$$

Where:

- ρ_k = autocorrelation at lag k
- γ_k = autocovariance at lag k
- γ_0 = variance of the series (autocovariance at lag 0)

Autocovariance Definition:

$$\gamma_k = Cov(Y_t, Y_{t-k}) = E[(Y_t - \mu)(Y_{t-k} - \mu)]$$

Sample Autocorrelation:

$$r_k = \frac{c_k}{c_0} = \frac{\frac{1}{n} \sum_{t=k+1}^n (Y_t - \bar{Y})(Y_{t-k} - \bar{Y})}{\frac{1}{n} \sum_{t=1}^n (Y_t - \bar{Y})^2}$$

Key Properties:

1. $\rho_0 = 1$ (perfect correlation with itself)
2. $-1 \leq \rho_k \leq 1$ for all k
3. $\rho_k = \rho_{-k}$ (symmetry property)

Sumit (The Inquisitive)

Professor, I understand the mathematical definition, but could you explain the intuitive meaning? What does it really mean when we say "autocorrelation at lag 2 is 0.7"?

Excellent question, Sumit! Let me break this down intuitively.

Autocorrelation at lag k answers: "How similar is today's value to the value k periods ago?"

Practical Interpretation:

- $\rho_1 = 0.7$: Today's value is strongly similar to yesterday's value
- $\rho_2 = 0.3$: Today's value is moderately similar to the value 2 days ago
- $\rho_5 = 0.1$: Today's value is weakly similar to the value 5 days ago
- $\rho_{10} = 0.0$: Today's value has no linear relationship with 10 days ago

Real Example - Stock Returns: If daily stock returns have $\rho_1 = 0.15$, it means: "If stock went up today, there's a 15% higher tendency for it to go up tomorrow too"

This is evidence of **momentum** - a key concept in finance!

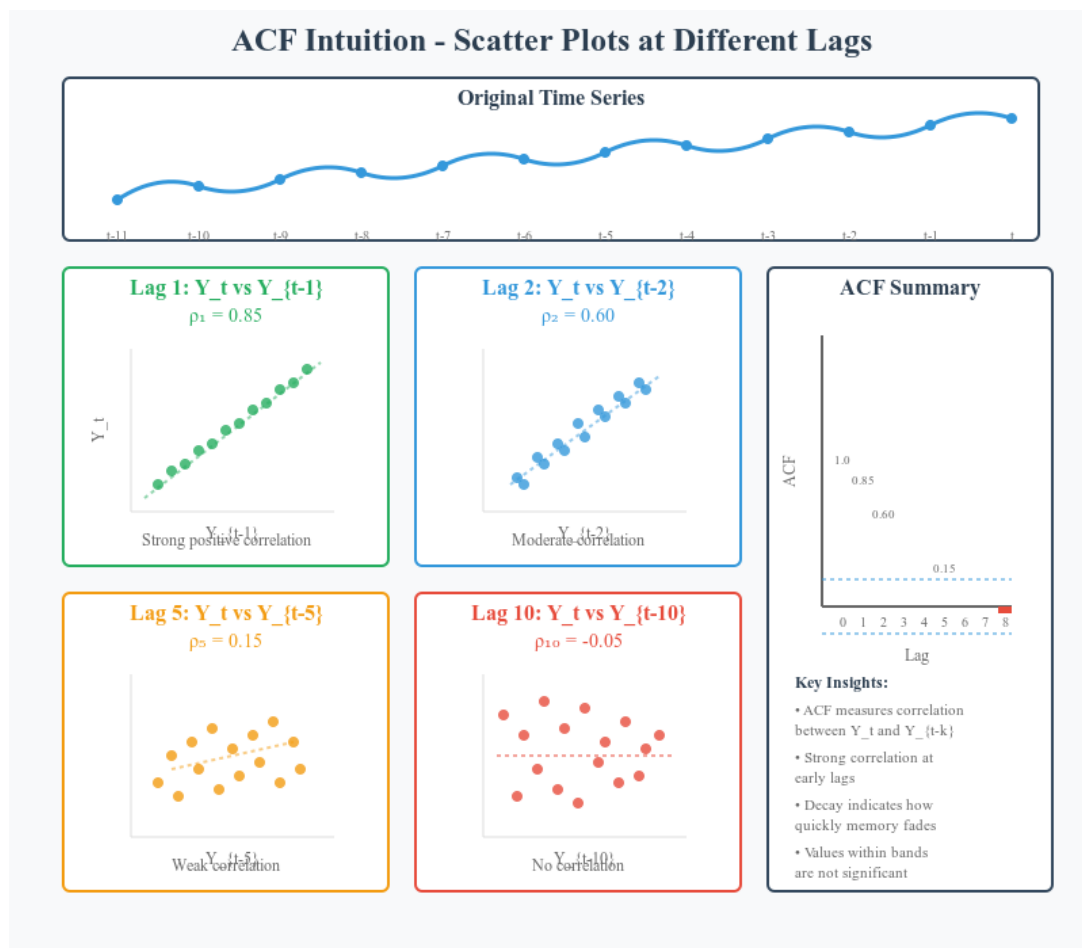


Figure 13: ACF Intuition - Scatter Plots at Different Lags

2.3.3 Properties of ACF for Different Processes

Theory Deep Dive

ACF Patterns for Common Time Series Models

Different time series processes have characteristic ACF patterns:

1. White Noise Process:

$$Y_t = \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2)$$

ACF: $\rho_k = 0$ for all $k > 0$

2. Random Walk:

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

ACF: $\rho_k \approx 1$ for all k (very slow decay)

3. AR(1) Process:

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

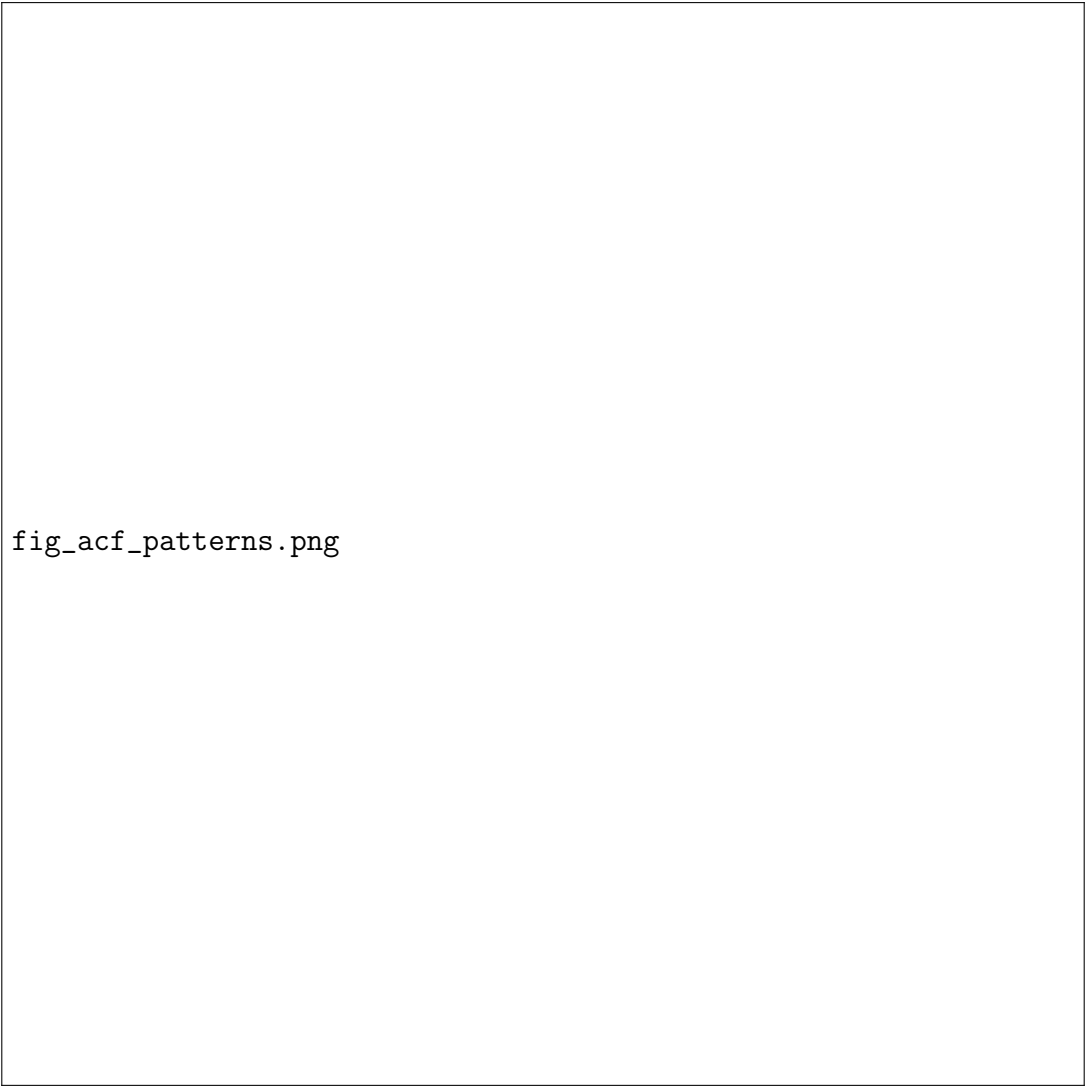
ACF: $\rho_k = \phi^k$ (exponential decay)

4. MA(1) Process:

$$Y_t = \epsilon_t + \theta \epsilon_{t-1}$$

ACF: $\rho_1 = \frac{\theta}{1+\theta^2}$, $\rho_k = 0$ for $k > 1$

5. Seasonal Process: ACF: Peaks at seasonal lags (12, 24, 36... for monthly data)



fig_acf_patterns.png

Figure 14: Characteristic ACF Patterns for Different Time Series Models

Neha (The Skeptic)

Professor, I'm skeptical about interpreting these ACF patterns. How do we know that high autocorrelation at lag 1 isn't just a coincidence? What if we're seeing patterns where none exist?

Professor

Outstanding skepticism, Neha! This is exactly why we need statistical significance testing for autocorrelations.

Statistical Significance of ACF:

For a white noise process (no autocorrelation), the sample autocorrelations are approximately:

$$r_k \sim N\left(0, \frac{1}{n}\right)$$

95% Confidence Bounds:

$$\pm 1.96 \times \frac{1}{\sqrt{n}}$$

Interpretation:

- If $|r_k| > \frac{1.96}{\sqrt{n}}$: Autocorrelation is statistically significant
- If $|r_k| \leq \frac{1.96}{\sqrt{n}}$: Could be just random noise

Example: With $n=1000$ observations:

- Confidence bounds = $\pm 1.96/1000 = \pm 0.00196$
- If $r_1 = 0.15$: Significant ($0.15 \gg 0.00196$)
- If $r_5 = 0.03$: Not significant ($0.03 \gg 0.00196$)

This is why ACF plots always show confidence bands!

Python Code Reference

Python Code Reference: 03_autocorrelation_analysis.py

ACF calculation and plotting:

```
from statsmodels.graphics.tsaplots import plot_acf
from statsmodels.tsa.stattools import acf

# Calculate ACF values
acf_values = acf(data, nlags=20, fft=False)

# Plot with confidence intervals
plot_acf(data, lags=20, alpha=0.05)
plt.show()
```

The $\alpha=0.05$ parameter adds 95% confidence bands!

2.3.4 Partial Autocorrelation Function (PACF) - The Direct Relationship

Professor

While ACF shows the total correlation between Y_t and Y_{t-k} , it includes both direct and indirect effects. The Partial Autocorrelation Function (PACF) shows only the **direct** relationship after removing the influence of intermediate lags.

Think of it this way: If you want to know how your mood today relates to your mood 3 days ago, ACF includes the influence of yesterday and the day before yesterday. PACF removes these influences to show only the direct 3-day relationship.

Theory Deep Dive

Partial Autocorrelation Function (PACF)

PACF measures the correlation between Y_t and Y_{t-k} after removing the linear dependence on the intermediate variables $Y_{t-1}, Y_{t-2}, \dots, Y_{t-k+1}$.

Mathematical Definition: The PACF at lag k is the coefficient ϕ_{kk} in the regression:

$$Y_t = \phi_{k1}Y_{t-1} + \phi_{k2}Y_{t-2} + \dots + \phi_{kk}Y_{t-k} + \epsilon_t$$

Recursive Calculation (Durbin-Levinson Algorithm):

$$\phi_{kk} = \frac{\rho_k - \sum_{j=1}^{k-1} \phi_{k-1,j} \rho_{k-j}}{1 - \sum_{j=1}^{k-1} \phi_{k-1,j} \rho_j}$$

Key Properties:

1. $\phi_{11} = \rho_1$ (PACF at lag 1 equals ACF at lag 1)
2. PACF helps identify the order of autoregressive processes
3. PACF at lag k is the partial correlation coefficient

PACF Patterns:

- **AR(p):** PACF cuts off after lag p
- **MA(q):** PACF decays gradually
- **White Noise:** PACF = 0 for all lags $k \neq 0$

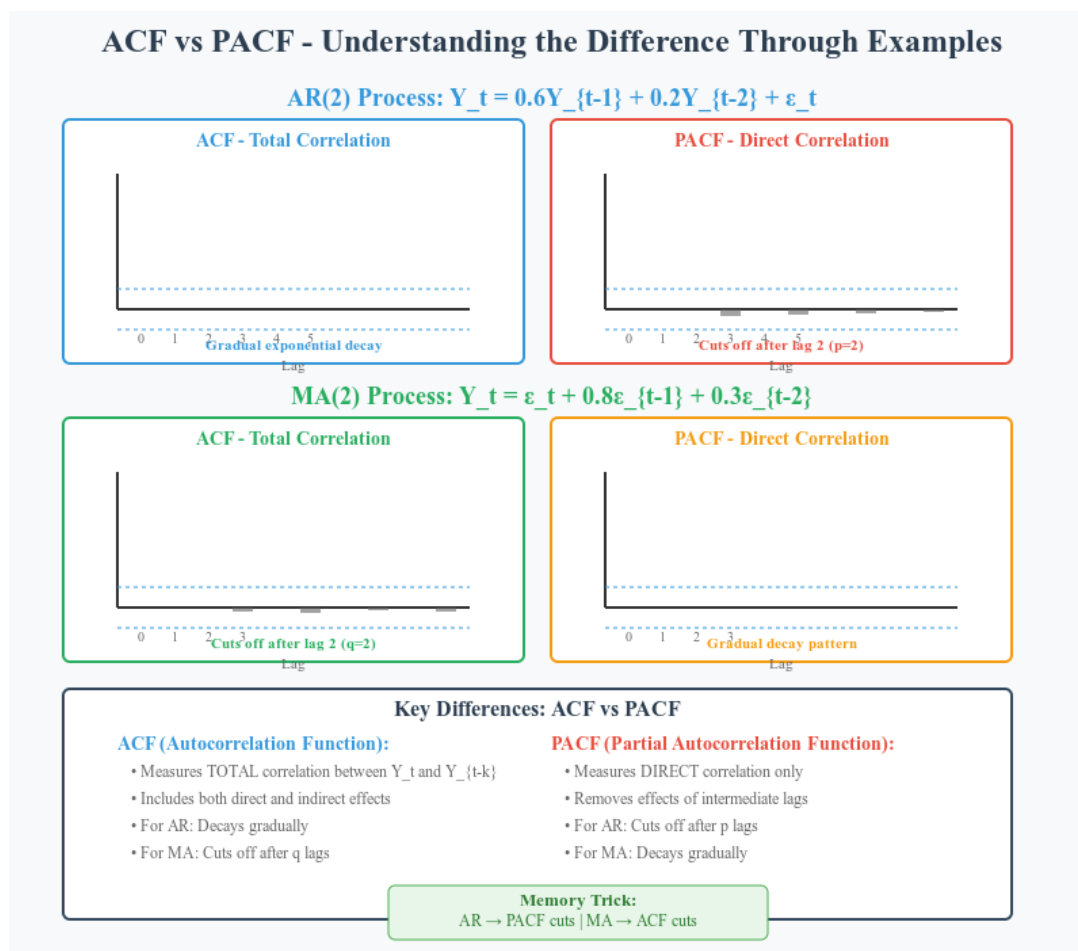


Figure 15: ACF vs PACF - Understanding the Difference Through Examples

Paul (The Innocent)

Professor, I'm struggling to understand the difference between ACF and PACF. Could you give me a simple, real-world analogy to help me remember?

Professor

Perfect question, Paul! Let me give you an analogy that will make this crystal clear.

The Gossip Chain Analogy:

Imagine a gossip chain in your office:

- Monday: You tell a secret to Person A
- Tuesday: Person A tells Person B
- Wednesday: Person B tells Person C
- Thursday: Person C tells Person D

ACF Perspective: "How similar is the story on Thursday to the original Monday story?"

- This measures the **total** correlation
- Includes all the distortions along the chain
- Story similarity decreases with each step

PACF Perspective: "How much does Person C directly contribute to Person D's version, ignoring what they already knew from A and B?"

- This measures the **direct** contribution only
- Removes the influence of intermediate people
- Shows only the new information added at each step

Time Series Translation:

- ACF: Total influence of past on present
- PACF: Direct influence after removing intermediate effects

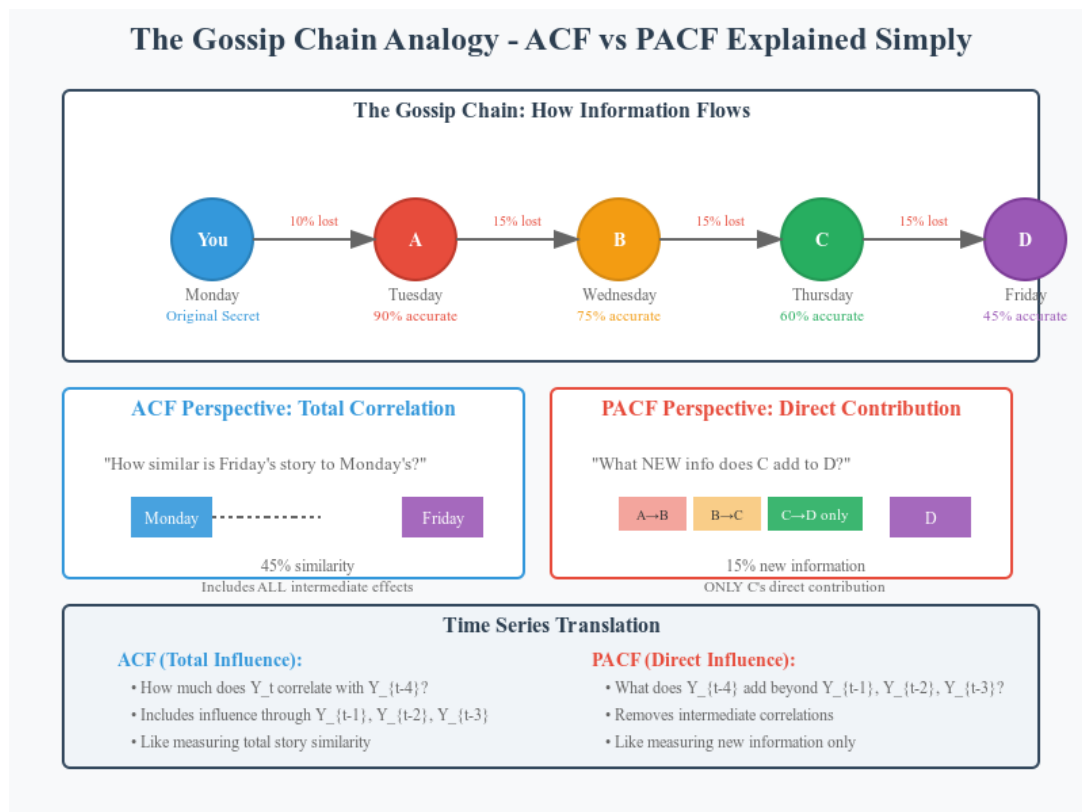


Figure 16: The Gossip Chain Analogy - ACF vs PACF Explained Simply

2.3.5 Ljung-Box Test - Testing for Autocorrelation

Theory Deep Dive

Ljung-Box Test for Serial Correlation

The Ljung-Box test formally tests whether a time series is white noise (no autocorrelation).

Hypotheses:

$$H_0 : \rho_1 = \rho_2 = \dots = \rho_h = 0 \text{ (white noise)}$$

$$H_1 : \text{At least one } \rho_k \neq 0 \text{ (autocorrelation exists)}$$

Test Statistic:

$$Q_{LB} = n(n+2) \sum_{k=1}^h \frac{r_k^2}{n-k}$$

Where:

- n = sample size
- h = number of lags tested
- r_k = sample autocorrelation at lag k

Distribution: Under H_0 : $Q_{LB} \sim \chi_h^2$

Decision Rule:

- If $p\text{-value} < 0.05$: Reject H_0 (autocorrelation exists)
- If $p\text{-value} \geq 0.05$: Fail to reject H_0 (no evidence of autocorrelation)

Brain Teaser

ACF/PACF Pattern Recognition Game!

Match each ACF/PACF pattern combination to the correct time series model:

ACF Pattern	PACF Pattern	Model
Exponential decay	Cut off after lag 1	?
Cut off after lag 1	Exponential decay	?
Cut off after lag 2	Exponential decay	?
Exponential decay	Cut off after lag 2	?
All zeros	All zeros	?

Models to choose from: AR(1), AR(2), MA(1), MA(2), White Noise

Answers: AR(1), MA(1), MA(2), AR(2), White Noise

Memory Trick:

- AR \rightarrow PACF cuts off (AutoRegressive \rightarrow Partial cuts)
- MA \rightarrow ACF cuts off (Moving Average \rightarrow Autocorrelation cuts)

3 Module 2: The Mathematics of Time Series Models

3.1 Stationarity - The Foundation of Time Series Modeling

Professor

We now come to the most crucial concept in time series analysis: stationarity. Think of stationarity as the "steady state" of a time series - a condition where the statistical properties remain constant over time.

Why is this so important? Because most time series models assume that the relationships between variables don't change over time. If the data's behavior keeps changing, our models become unreliable - like trying to hit a moving target!

Let's explore this fundamental concept step by step.

3.1.1 Understanding Stationarity Mathematically

Theory Deep Dive

Strict Stationarity

A time series $\{Y_t\}$ is strictly stationary if the joint probability distribution of any collection of observations is invariant under time shifts.

Mathematical Definition:

$$F(Y_{t_1}, Y_{t_2}, \dots, Y_{t_k}) = F(Y_{t_1+h}, Y_{t_2+h}, \dots, Y_{t_k+h})$$

for any t_1, t_2, \dots, t_k and any lag h .

Weak (Second-Order) Stationarity

A time series is weakly stationary if:

1. **Constant Mean:** $E[Y_t] = \mu$ for all t
2. **Constant Variance:** $Var[Y_t] = \sigma^2$ for all t
3. **Autocovariance depends only on lag:** $Cov[Y_t, Y_{t+k}] = \gamma_k$ for all t

Mathematical Conditions:

$$E[Y_t] = \mu \quad (\text{constant mean}) \tag{4}$$

$$Var[Y_t] = E[(Y_t - \mu)^2] = \sigma^2 \quad (\text{constant variance}) \tag{5}$$

$$\gamma_k = E[(Y_t - \mu)(Y_{t+k} - \mu)] \quad (\text{lag-dependent covariance}) \tag{6}$$

Autocorrelation Function for Stationary Series:

$$\rho_k = \frac{\gamma_k}{\gamma_0} = \frac{\gamma_k}{\sigma^2}$$

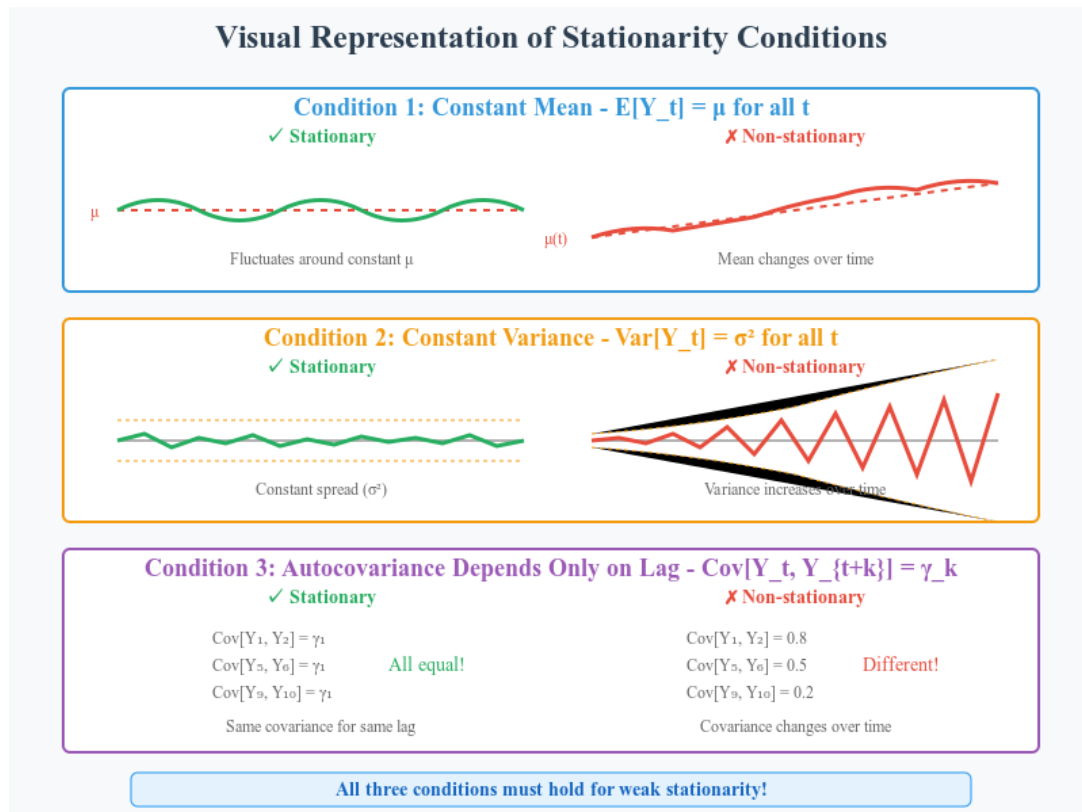


Figure 17: Visual Representation of Stationarity Conditions

Rohan (The Visual Learner)

Professor, I can see the three conditions in the diagram, but I'm having trouble visualizing what non-stationary data looks like. Could you show me examples of data that violate each condition?

Excellent visual question, Rohan! Let me show you clear examples of each type of non-stationarity:

Violation 1 - Non-constant Mean (Trend):

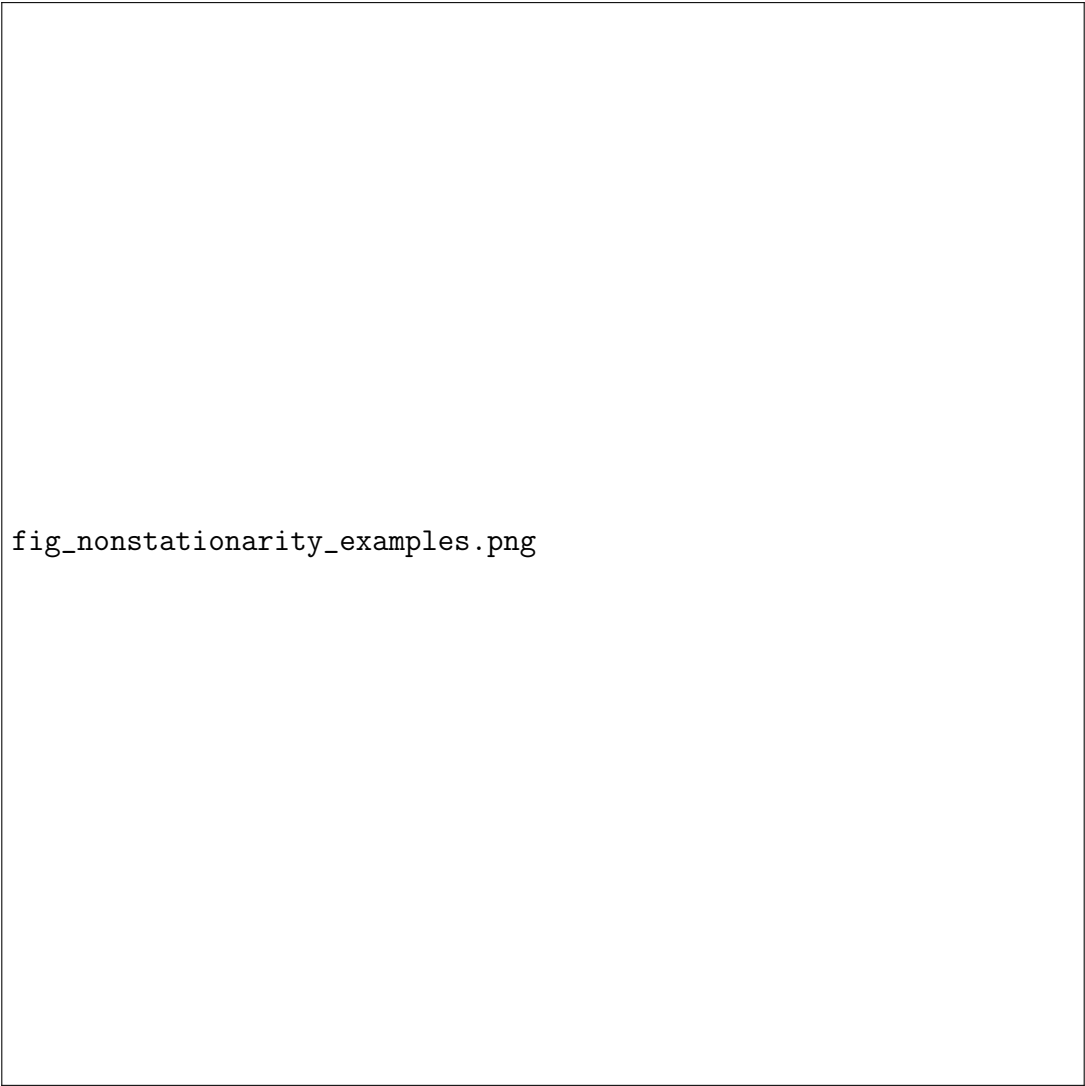
- Example: Stock prices, GDP, population
- Pattern: Clear upward or downward movement over time
- Mathematical: $E[Y_t] = \alpha + \beta t$ (depends on time)

Violation 2 - Non-constant Variance (Heteroscedasticity):

- Example: Stock volatility during crisis periods
- Pattern: Variability changes over time
- Mathematical: $Var[Y_t] = \sigma_t^2$ (time-dependent variance)

Violation 3 - Time-dependent Covariance:

- Example: Structural breaks in relationships
- Pattern: Correlation patterns change over time
- Mathematical: $Cov[Y_t, Y_{t+k}] = \gamma_{k,t}$ (depends on both lag and time)



fig_nonstationarity_examples.png

Figure 18: Examples of Non-Stationary Behavior - Violating Each Condition

3.1.2 Why Stationarity Matters

Professor

You might wonder: "Why should I care about stationarity?" The answer is fundamental to successful time series modeling and forecasting.

Consequences of Non-Stationarity

1. Model Instability: Non-stationary series can lead to:

- Coefficient estimates that change over time
- Models that work well in one period but fail in another
- Unreliable standard errors and confidence intervals

2. Spurious Regression: Two independent non-stationary series often show high correlation:

- High R^2 values even when variables are unrelated
- Misleading statistical significance
- Wrong conclusions about relationships

3. Forecasting Problems:

- Forecast intervals become unreliably wide
- Point forecasts may diverge to infinity
- Model performance degrades over time

4. Invalid Statistical Inference: Standard statistical tests assume stationarity:

- t-tests and F-tests become invalid
- Confidence intervals have wrong coverage
- Hypothesis tests give misleading results

Real-World Example

Real Example: The Spurious Regression Problem

Consider these two completely unrelated time series:

- US GDP (trending upward)
- Number of smartphones sold globally (trending upward)

If we regress one on the other without checking stationarity:

- $R^2 = 0.98$ (seems like strong relationship!)
- p-value ≤ 0.001 (highly significant!)
- Conclusion: GDP causes smartphone sales (completely wrong!)

The high correlation is spurious - both series trend upward over time. After differencing (making stationary), the correlation disappears, revealing no true relationship.

Lesson: Always test for stationarity before modeling relationships!

3.1.3 Testing for Stationarity

Theory Deep Dive

Augmented Dickey-Fuller (ADF) Test

The most common test for stationarity is the ADF test, which tests for the presence of a unit root.

Test Regression:

$$\Delta Y_t = \alpha + \beta Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$$

Where $\Delta Y_t = Y_t - Y_{t-1}$ is the first difference.

Hypotheses:

$H_0 : \beta = 0$ (unit root exists, non-stationary)

$H_1 : \beta < 0$ (no unit root, stationary)

Interpretation of :

- If $\beta = 0$: Series is a random walk (non-stationary)
- If $\beta < 0$: Series has mean reversion (stationary)
- More negative indicates stronger mean reversion

Test Statistic:

$$t_{ADF} = \frac{\hat{\beta}}{SE(\hat{\beta})}$$

Decision Rule:

- If $t_{ADF} < \text{critical value}$: Reject H_0 (stationary)
- If $t_{ADF} \geq \text{critical value}$: Fail to reject H_0 (non-stationary)

Variants:

1. No constant, no trend: $\Delta Y_t = \beta Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$
2. With constant: $\Delta Y_t = \alpha + \beta Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$
3. With constant and trend: $\Delta Y_t = \alpha + \delta t + \beta Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$

Sumit (The Inquisitive)

Professor, I understand the ADF test mechanically, but I'm curious about the intuition. Why does testing whether $\beta = 0$ tell us about stationarity? And how do we choose which variant to use?

Professor

Brilliant question, Sumit! Let me explain the deep intuition behind the ADF test.

The Unit Root Intuition:

The ADF regression can be rewritten as:

$$Y_t = \alpha + (\beta + 1)Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$$

Setting $\phi = \beta + 1$:

$$Y_t = \alpha + \phi Y_{t-1} + \sum_{i=1}^p \gamma_i \Delta Y_{t-i} + \epsilon_t$$

Now:

- If $\beta = 0 \Rightarrow \phi = 1$: Unit root (random walk)
- If $\beta < 0 \Rightarrow \phi < 1$: Mean reversion (stationary)

Economic Interpretation:

- $\phi = 1$: Shocks have permanent effects (non-stationary)
- $\phi < 1$: Shocks die out over time (stationary)
- $|\phi| > 1$: Explosive process (also non-stationary)

Choosing ADF Variant:

1. Plot your data first
2. If no obvious trend \rightarrow use "constant only"
3. If clear trend \rightarrow use "constant and trend"
4. If series fluctuates around zero \rightarrow use "no constant"

Theory Deep Dive

KPSS Test - Alternative Approach

The KPSS test takes the opposite approach to ADF:

Hypotheses:

H_0 : Series is stationary

H_1 : Series has a unit root

Test Statistic:

$$KPSS = \frac{1}{n^2} \frac{\sum_{t=1}^n S_t^2}{\hat{\sigma}^2}$$

Where $S_t = \sum_{i=1}^t \hat{\epsilon}_i$ and $\hat{\epsilon}_t$ are residuals from regression on constant/trend.

Combined Testing Strategy:

ADF Result	KPSS Result	Conclusion
Stationary	Stationary	Stationary
Non-stationary	Non-stationary	Non-stationary
Stationary	Non-stationary	Trend-stationary
Non-stationary	Stationary	Near unit root

Python Code Reference

Python Code Reference: 02_stationarity_tests.py

Stationarity testing:

```
from statsmodels.tsa.stattools import adfuller, kpss

# ADF test
adf_result = adfuller(data, autolag='AIC')
print(f"ADF p-value: {adf_result[1]}")

# KPSS test
kpss_result = kpss(data, regression='ct')
print(f"KPSS p-value: {kpss_result[1]}")
```

This shows how to test both ways for robust conclusions!

3.1.4 Making Data Stationary

Professor

Once we've identified non-stationarity, we need to transform the data to achieve stationarity. There are several transformation techniques, each appropriate for different types of non-stationarity.

Transformation Techniques

1. First Differencing:

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

Effective for: Removing linear trends

2. Seasonal Differencing:

$$\Delta_s Y_t = Y_t - Y_{t-s}$$

Where s is the seasonal period. Effective for: Removing seasonal patterns

3. Combined Differencing:

$$\Delta \Delta_s Y_t = (Y_t - Y_{t-1}) - (Y_{t-s} - Y_{t-s-1})$$

Effective for: Series with both trend and seasonality

4. Log Transformation:

$$Z_t = \log(Y_t)$$

Effective for: Stabilizing variance, converting multiplicative to additive

5. Box-Cox Transformation:

$$Z_t = \begin{cases} \frac{Y_t^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(Y_t) & \text{if } \lambda = 0 \end{cases}$$

Where λ is chosen to maximize normality.

6. Detrending:

$$Z_t = Y_t - \hat{T}_t$$

Where \hat{T}_t is the estimated trend component.

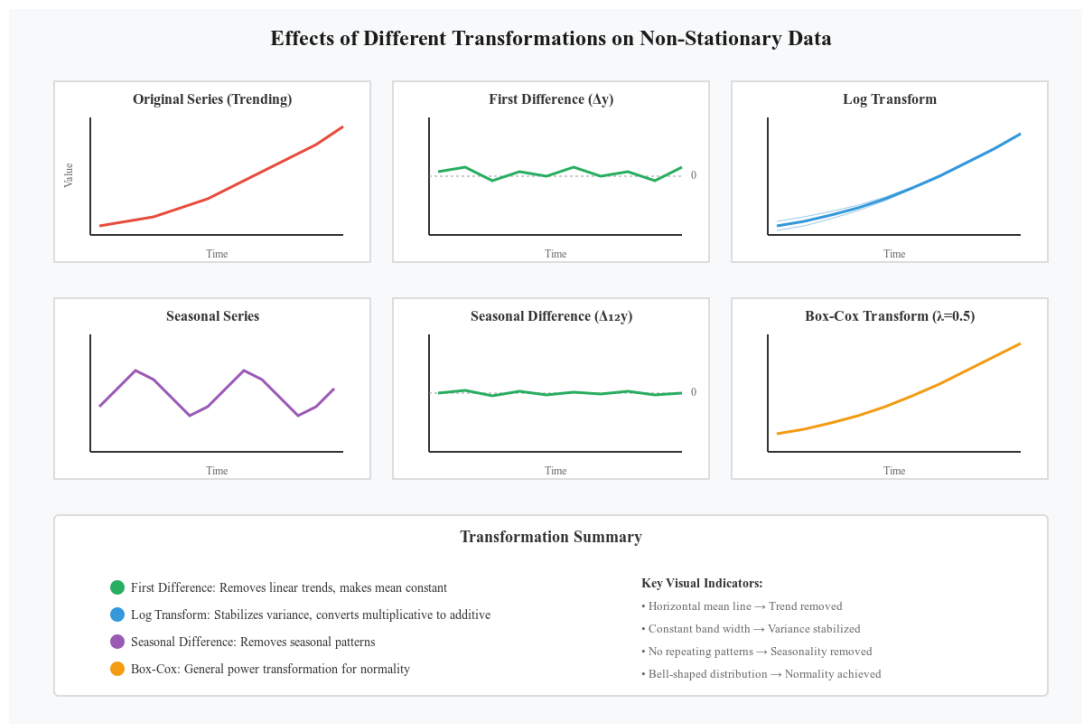


Figure 19: Effects of Different Transformations on Non-Stationary Data

Neha (The Skeptic)

Professor, I'm skeptical about blindly applying these transformations. How do we know which transformation to use? And what if we over-difference and lose important information?

Professor

Excellent skepticism, Neha! You're absolutely right - inappropriate transformations can destroy valuable information.

Choosing the Right Transformation:

1. Diagnose the Problem First:

- Trending mean → First differencing
- Changing variance → Log or Box-Cox
- Seasonal patterns → Seasonal differencing
- Multiple issues → Combined approach

2. Signs of Over-Differencing:

- ACF becomes strongly negative at lag 1
- Variance increases after differencing
- Model becomes more complex (needs more MA terms)

3. Preservation of Information:

- Always keep track of transformations applied
- Test stationarity after each transformation
- Stop when stationarity is achieved
- Remember to reverse transformations for forecasting!

Golden Rule: Use the minimum transformation necessary to achieve stationarity.

Important Warning

Common Transformation Mistakes

1. **Over-differencing:** Taking differences when data is already stationary
2. **Wrong order:** Applying log after differencing (should be reverse)
3. **Forgetting to reverse:** Not inverting transformations for final forecasts
4. **Mixing transformations:** Using different transformations for different parts of series

Always document your transformation pipeline!

Brain Teaser

Transformation Decision Tree Exercise

For each scenario, choose the appropriate transformation:

1. Stock prices showing exponential growth with increasing volatility
 - Ⓐ First difference only
 - Ⓑ Log then first difference
 - Ⓒ Seasonal difference
2. Monthly sales with strong December peaks that grow proportionally
 - Ⓐ First difference
 - Ⓑ Seasonal difference ($s=12$)
 - Ⓒ Log then seasonal difference
3. Temperature data with no trend but daily patterns
 - Ⓐ No transformation needed
 - Ⓑ First difference
 - Ⓒ Seasonal difference ($s=24$)

Answers: B (stabilize variance first), C (handle multiplicative seasonality), A (already stationary)

3.2 Linear Processes and Moving Average Models

Professor

Now that we understand stationarity, let's dive into the building blocks of time series models. We'll start with linear processes - the foundation upon which all ARMA models are built.

Think of linear processes as recipes where we combine random shocks (ingredients) with different weights (proportions) to create our time series dish!

3.2.1 General Linear Process

Theory Deep Dive

General Linear Process

A time series $\{Y_t\}$ is a linear process if it can be written as:

$$Y_t = \mu + \sum_{j=-\infty}^{\infty} \psi_j \epsilon_{t-j}$$

Where:

- μ = mean of the process
- $\epsilon_t \sim WN(0, \sigma^2)$ = white noise
- ψ_j = weights (impulse response coefficients)
- $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$ (absolute summability)

Special Cases:

1. **Causal Process:** $\psi_j = 0$ for $j < 0$ (depends only on past)
2. **Invertible Process:** Can be written as $AR(\infty)$
3. **Finite Order:** Only finite number of $\psi_j \neq 0$

Properties:

- Mean: $E[Y_t] = \mu$
- Variance: $Var[Y_t] = \sigma^2 \sum_{j=-\infty}^{\infty} \psi_j^2$
- Autocovariance: $\gamma_k = \sigma^2 \sum_{j=-\infty}^{\infty} \psi_j \psi_{j+k}$

3.2.2 Moving Average Process MA(q)

Theory Deep Dive

Moving Average Model of Order q - MA(q)

The MA(q) model represents the current value as a weighted combination of the current and past q random shocks.

Mathematical Definition:

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

Or in compact form:

$$Y_t = \mu + \sum_{j=0}^q \theta_j \epsilon_{t-j}, \quad \theta_0 = 1$$

Using Backshift Operator:

$$Y_t = \mu + \theta(B)\epsilon_t$$

Where $\theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \dots + \theta_q B^q$

Properties of MA(q):

1. Always stationary (no restrictions on θ_j)
2. Mean: $E[Y_t] = \mu$
3. Variance: $Var[Y_t] = \sigma^2(1 + \theta_1^2 + \theta_2^2 + \dots + \theta_q^2)$
4. ACF cuts off after lag q
5. PACF decays exponentially

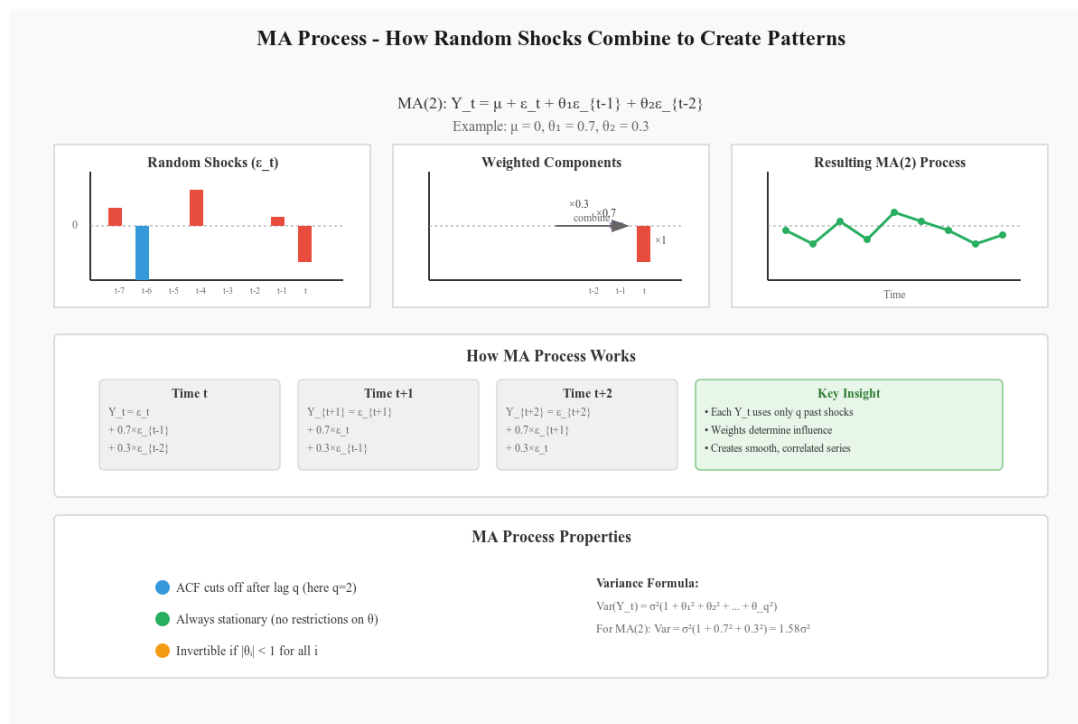


Figure 20: MA Process - How Random Shocks Combine to Create Patterns

Paul (The Innocent)

Professor, I'm confused about the MA model. It says "moving average" but it doesn't look like the moving averages we studied earlier. What's the connection?

Professor

Excellent observation, Paul! The name is indeed confusing. Let me clarify:

Moving Average (Smoothing) vs MA Model:

Smoothing MA: $\bar{Y}_t = \frac{1}{k}(Y_t + Y_{t-1} + \dots + Y_{t-k+1})$

- Averages observed values
- Used for smoothing/trend extraction
- Deterministic calculation

MA Model: $Y_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$

- Weighted sum of random shocks
- Used for modeling/forecasting
- Stochastic process

The name comes from the fact that it's a "weighted moving sum" of error terms. Think of it as nature's way of smoothing random shocks into observable patterns!

3.2.3 MA(1) Process - Deep Dive

Theory Deep Dive

MA(1) Process in Detail

The simplest MA model combines only current and previous shock:

$$Y_t = \mu + \epsilon_t + \theta\epsilon_{t-1}$$

Autocovariance Function:

$$\gamma_0 = \text{Var}[Y_t] = \sigma^2(1 + \theta^2) \quad (7)$$

$$\gamma_1 = \text{Cov}[Y_t, Y_{t-1}] = \sigma^2\theta \quad (8)$$

$$\gamma_k = 0 \text{ for } k > 1 \quad (9)$$

Autocorrelation Function:

$$\rho_1 = \frac{\gamma_1}{\gamma_0} = \frac{\theta}{1 + \theta^2}$$

$$\rho_k = 0 \text{ for } k > 1$$

Key Insight: ACF cuts off after lag 1!

Invertibility Condition: MA(1) is invertible if $|\theta| < 1$, allowing representation:

$$\epsilon_t = Y_t - \theta\epsilon_{t-1} = \sum_{j=0}^{\infty} (-\theta)^j Y_{t-j}$$

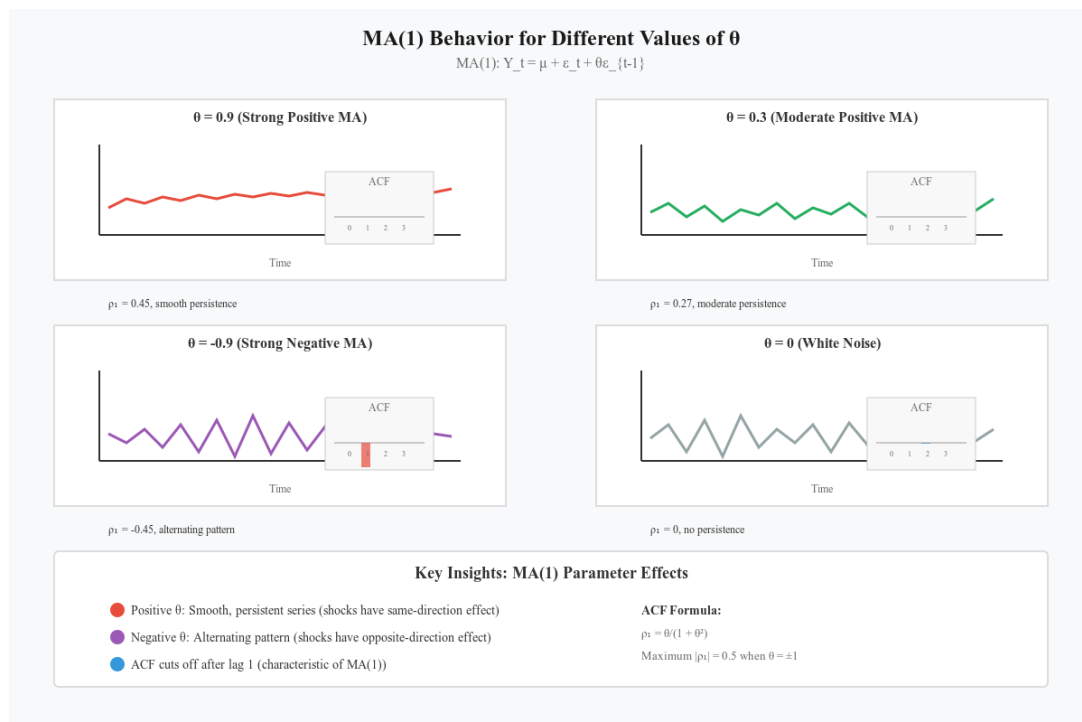


Figure 21: MA(1) Behavior for Different Values of

Sumit (The Inquisitive)

Professor, I notice that both $\theta = 0.5$ and $\theta = 2$ give the same autocorrelation $\rho_1 = 0.4$. How is this possible? And why does invertibility matter?

Professor

Brilliant observation, Sumit! You've discovered the identification problem in MA models.

The Duality Problem: For MA(1), the ACF is: $\rho_1 = \frac{\theta}{1+\theta^2}$
If we solve for θ given ρ_1 :

$$\theta^2 - \frac{\theta}{\rho_1} + 1 = 0$$

This quadratic has two solutions: θ and $1/\theta$

Example: If $\rho_1 = 0.4$:

- $\theta_1 = 0.5$ (invertible)
- $\theta_2 = 2.0 = 1/0.5$ (non-invertible)

Why Invertibility Matters:

1. **Uniqueness:** Only invertible MA has unique representation
2. **Forecasting:** Non-invertible MA gives explosive forecasts
3. **Estimation:** Maximum likelihood naturally selects invertible solution
4. **Interpretation:** Invertible = recent shocks matter more than distant ones

We always choose the invertible solution!

3.3 Autoregressive Models AR(p)

Professor

While MA models express the present as a combination of past shocks, AR models take a more intuitive approach: they express the present as a function of the past values themselves. This is often more natural - today's temperature depends on yesterday's temperature, today's stock price depends on recent prices, and so on.

3.3.1 AR(p) Model Definition

Theory Deep Dive

Autoregressive Model of Order p - AR(p)

An AR(p) model expresses the current value as a linear combination of the past p values plus a random shock.

Mathematical Definition:

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

Where:

- c = constant term
- ϕ_i = autoregressive coefficients
- $\epsilon_t \sim WN(0, \sigma^2)$ = white noise

Mean-Centered Form: If $\mu = E[Y_t]$, then:

$$(Y_t - \mu) = \phi_1 (Y_{t-1} - \mu) + \dots + \phi_p (Y_{t-p} - \mu) + \epsilon_t$$

Backshift Operator Form:

$$\phi(B)Y_t = c + \epsilon_t$$

Where $\phi(B) = 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p$

Properties:

- Stationary if roots of $\phi(z) = 0$ lie outside unit circle
- Mean: $\mu = \frac{c}{1 - \phi_1 - \phi_2 - \dots - \phi_p}$
- ACF decays exponentially (or as damped sinusoid)
- PACF cuts off after lag p

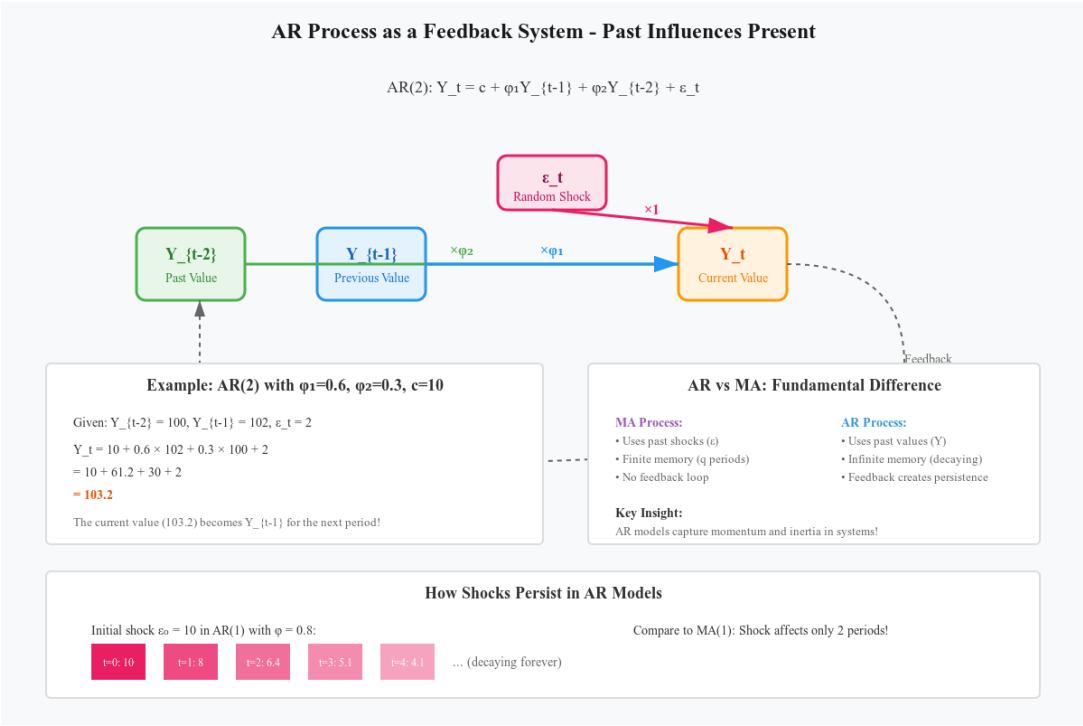


Figure 22: AR Process as a Feedback System - Past Influences Present

3.3.2 AR(1) Process - The Building Block

Theory Deep Dive

AR(1) Process in Detail

The simplest AR model:

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

Stationarity Condition: $|\phi| < 1$

Mean (if stationary):

$$\mu = E[Y_t] = \frac{c}{1 - \phi}$$

Variance (if stationary):

$$\gamma_0 = Var[Y_t] = \frac{\sigma^2}{1 - \phi^2}$$

Autocovariance:

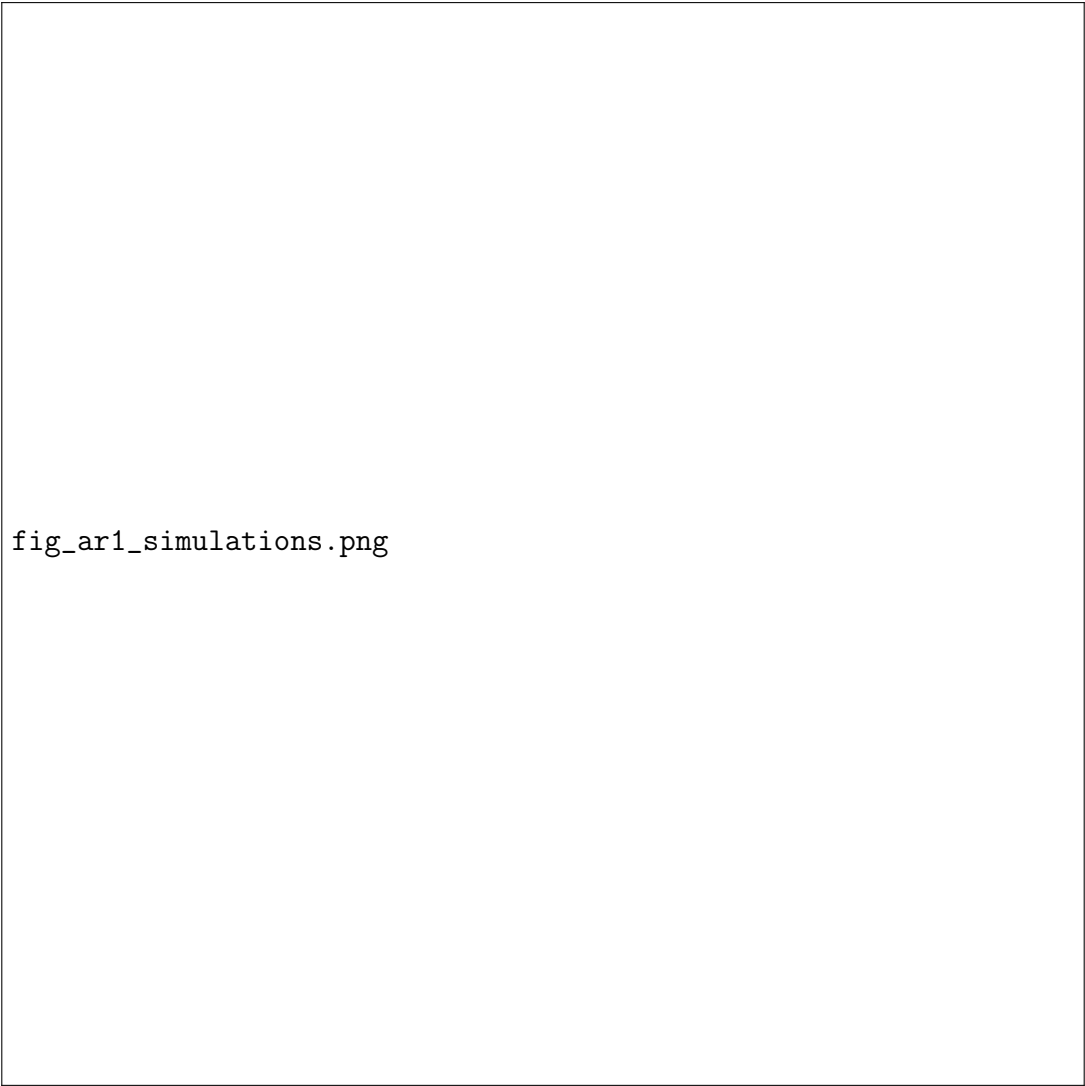
$$\gamma_k = \phi^k \gamma_0 = \phi^k \frac{\sigma^2}{1 - \phi^2}$$

Autocorrelation:

$$\rho_k = \phi^k$$

Key Behaviors:

- $0 < \phi < 1$: Positive autocorrelation, slow decay
- $-1 < \phi < 0$: Alternating autocorrelation
- $\phi = 1$: Random walk (non-stationary)
- $|\phi| > 1$: Explosive (non-stationary)



fig_ar1_simulations.png

Figure 23: AR(1) Process Behavior for Different ρ Values

Rohan (The Visual Learner)

Professor, I can see from the plots how different ρ values create different patterns. But how does this relate to real-world interpretation? What does $\rho = 0.9$ vs $\rho = 0.3$ mean in practice?

Professor

Excellent practical question, Rohan! The value of ϕ determines the "persistence" or "memory" of the process.

Real-World Interpretation of :

High (0.8 - 0.95):

- Strong persistence
- Shocks have long-lasting effects
- Examples: GDP, unemployment rate
- Interpretation: "Economic variables adjust slowly"

Moderate (0.3 - 0.7):

- Moderate persistence
- Balance between memory and innovation
- Examples: Inventory levels, sales
- Interpretation: "System responds but with inertia"

Low (0 - 0.3):

- Weak persistence
- Quickly returns to mean
- Examples: Stock returns, temperature deviations
- Interpretation: "System rapidly corrects deviations"

Half-Life of Shocks: Time for shock to decay to half: $h = \frac{\ln(0.5)}{\ln(\phi)}$

- $\phi = 0.9$: Half-life 6.6 periods
- $\phi = 0.5$: Half-life = 1 period
- $\phi = 0.3$: Half-life 0.6 periods

3.3.3 AR(2) Process and Stationarity Conditions

Theory Deep Dive

AR(2) Process

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

Stationarity Conditions: The characteristic equation: $1 - \phi_1 z - \phi_2 z^2 = 0$
Stationarity requires all roots outside unit circle, equivalent to:

1. $\phi_1 + \phi_2 < 1$
2. $\phi_2 - \phi_1 < 1$
3. $|\phi_2| < 1$

Stationarity Triangle: Valid (ϕ_1, ϕ_2) pairs form a triangle in parameter space.

Autocorrelation Function: Satisfies the difference equation:

$$\rho_k = \phi_1 \rho_{k-1} + \phi_2 \rho_{k-2}, \quad k \geq 2$$

With initial conditions: $\rho_0 = 1, \rho_1 = \frac{\phi_1}{1-\phi_2}$

Behavior Types:

- Real distinct roots: Exponential decay
- Complex roots: Damped oscillations
- Frequency of oscillation: $f = \frac{\cos^{-1}(\phi_1/(2\sqrt{-\phi_2}))}{2\pi}$

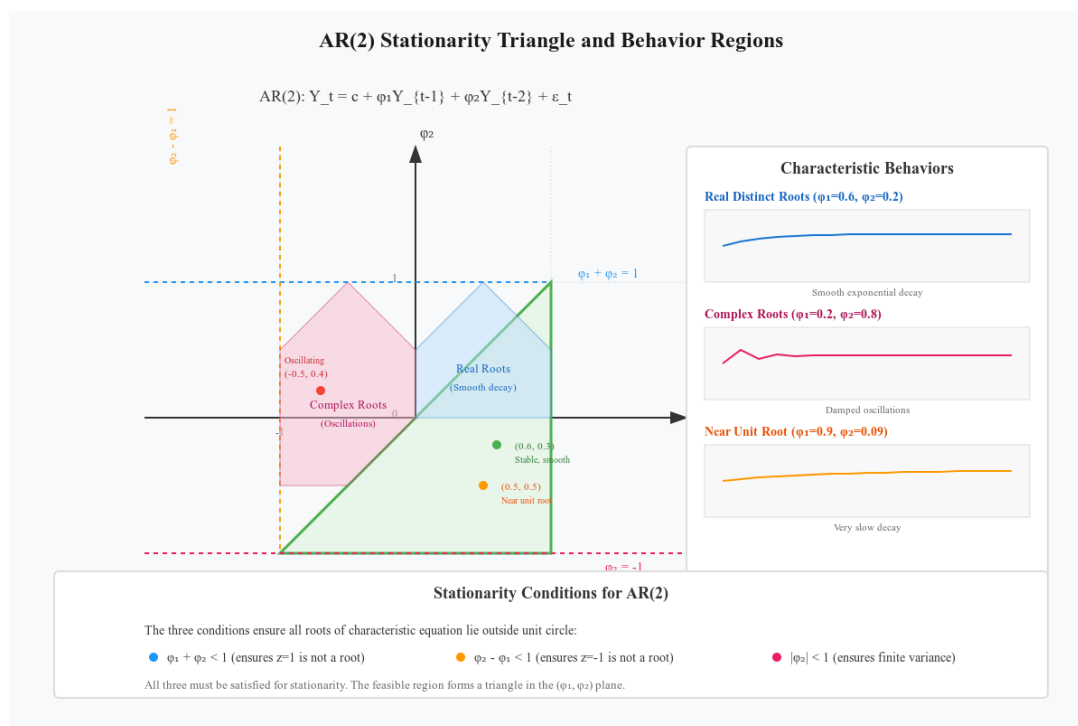


Figure 24: AR(2) Stationarity Triangle and Behavior Regions

Neha (The Skeptic)

Professor, I'm skeptical about these stationarity conditions. They seem arbitrary. Why exactly does $\phi_1 + \phi_2 < 1$ ensure stationarity? Can you prove this rigorously?

Professor

Excellent skepticism, Neha! Let me prove this rigorously.

Rigorous Proof of AR(2) Stationarity:

Starting with: $Y_t = c + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \epsilon_t$

The characteristic equation comes from:

$$(1 - \phi_1 B - \phi_2 B^2)Y_t = c + \epsilon_t$$

Setting $z = B^{-1}$ and solving:

$$z^2 - \phi_1 z - \phi_2 = 0$$

Roots: $z = \frac{\phi_1 \pm \sqrt{\phi_1^2 + 4\phi_2}}{2}$

For stationarity, need $|z_i| > 1$ for both roots.

Deriving the Conditions:

Condition 1: $\phi_1 + \phi_2 < 1$

- This ensures $z = 1$ is not a root
- Substituting: $1 - \phi_1 - \phi_2 \neq 0$

Condition 2: $\phi_2 - \phi_1 < 1$

- This ensures $z = -1$ is not a root
- Substituting: $1 + \phi_1 - \phi_2 \neq 0$

Condition 3: $|\phi_2| < 1$

- Product of roots = $-\phi_2$
- Need $|z_1 z_2| > 1$, so $|\phi_2| < 1$

These aren't arbitrary - they're the exact mathematical requirements for stability!

3.4 ARMA Models - Combining AR and MA

Professor

We've seen that AR models capture persistence and MA models capture the impact of shocks. Real-world data often exhibits both behaviors. ARMA models elegantly combine these two approaches, giving us a flexible framework for modeling a wide variety of time series patterns.

3.4.1 ARMA(p,q) Model Definition

Theory Deep Dive

Autoregressive Moving Average Model - ARMA(p,q)

The ARMA(p,q) model combines p autoregressive terms and q moving average terms:

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

Backshift Notation:

$$\phi(B)Y_t = c + \theta(B)\epsilon_t$$

Where:

- $\phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$ (AR polynomial)
- $\theta(B) = 1 + \theta_1 B + \dots + \theta_q B^q$ (MA polynomial)

Stationarity: Requires roots of $\phi(z) = 0$ outside unit circle **Invertibility:** Requires roots of $\theta(z) = 0$ outside unit circle

Mean (if stationary):

$$\mu = \frac{c}{1 - \phi_1 - \dots - \phi_p}$$

Properties:

- Combines benefits of both AR and MA
- Often more parsimonious than pure AR or MA
- ACF: Decays after lag q (influenced by AR part)
- PACF: Decays after lag p (influenced by MA part)

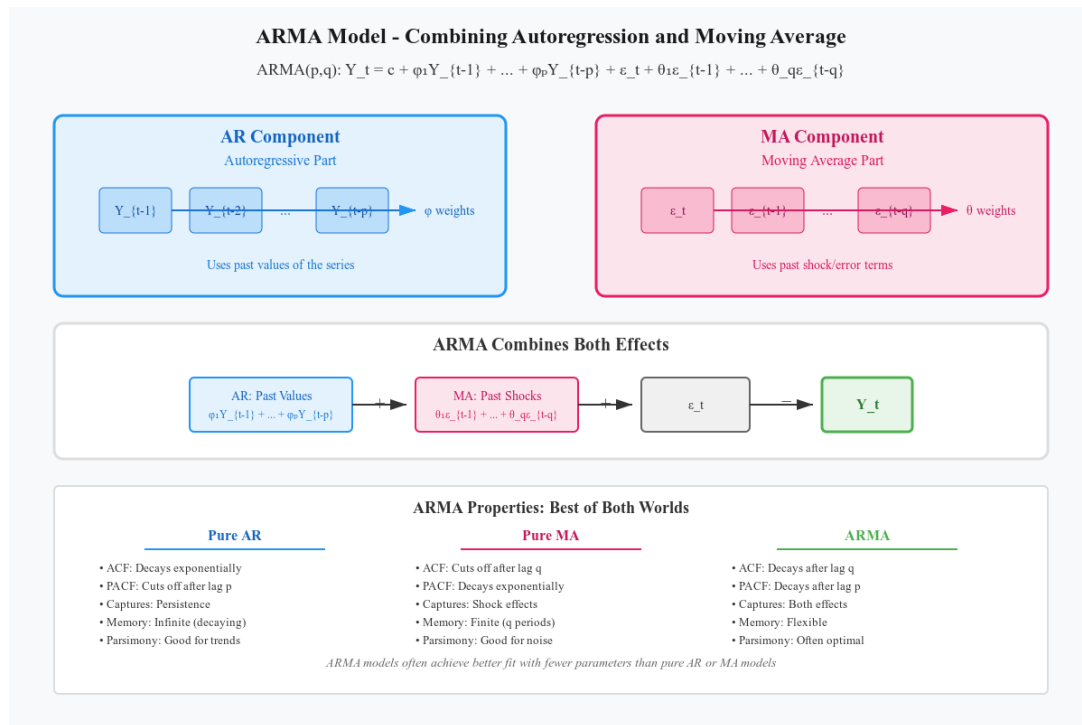


Figure 25: ARMA Model - Combining Autoregression and Moving Average

3.4.2 ARMA(1,1) - The Workhorse Model

Theory Deep Dive

ARMA(1,1) Process

$$Y_t = c + \phi Y_{t-1} + \epsilon_t + \theta \epsilon_{t-1}$$

Or in mean-centered form:

$$(Y_t - \mu) = \phi(Y_{t-1} - \mu) + \epsilon_t + \theta \epsilon_{t-1}$$

Conditions:

- Stationary if $|\phi| < 1$
- Invertible if $|\theta| < 1$

Variance:

$$\gamma_0 = \frac{\sigma^2(1 + 2\phi\theta + \theta^2)}{1 - \phi^2}$$

Autocorrelation Function:

$$\rho_1 = \frac{\phi + \theta + \phi\theta^2}{1 + 2\phi\theta + \theta^2}$$

$$\rho_k = \phi \rho_{k-1} \text{ for } k \geq 2$$

Key Insight: After lag 1, ACF follows AR(1) pattern!

Paul (The Innocent)

Professor, ARMA(1,1) seems more complex than AR(1) or MA(1). When would we use it instead of the simpler models?

Professor

Great practical question, Paul! ARMA(1,1) is indeed more complex, but it's often the "Goldilocks" model - not too simple, not too complex.

When to Use ARMA(1,1):

1. ACF/PACF Patterns:

- Both ACF and PACF decay exponentially
- Neither cuts off cleanly
- Common in real data!

2. Model Parsimony:

- AR(4) with 4 parameters \rightarrow ARMA(1,1) with 2 parameters
- Same fit, fewer parameters = better model

3. Physical Interpretation:

- AR part: System inertia/momentum
- MA part: External shock absorption
- Example: Economic variables with both persistence and policy shocks

Real Example: Inflation often follows ARMA(1,1):

- AR component: Inflation expectations (persistence)
- MA component: Supply shocks (temporary effects)

Brain Teaser

Model Identification Challenge!

Match the ACF/PACF patterns to the correct model:

ACF Pattern	PACF Pattern	Model	Answer
Cuts off after lag 2	Exponential decay	?	
Exponential decay	Cuts off after lag 3	?	
Slow exponential decay	Slow exponential decay	?	
Oscillating decay	Cuts off after lag 2	?	
Cuts off after lag 1	Oscillating decay	?	

Models: MA(2), AR(3), ARMA(2,2) or higher, AR(2) with complex roots, MA(1)

Answers: MA(2), AR(3), ARMA(p,q), AR(2), MA(1)

3.5 Non-Stationary and Seasonal Models

Professor

Most real-world time series are non-stationary and exhibit seasonal patterns. We've learned how to make data stationary through differencing. Now, let's formalize this into the ARIMA and SARIMA frameworks - the workhorses of practical time series modeling.

3.5.1 ARIMA Models - Integration and Differencing

Theory Deep Dive

Autoregressive Integrated Moving Average - ARIMA(p,d,q)

ARIMA models extend ARMA to handle non-stationary series through differencing.

Definition: A series $\{Y_t\}$ follows ARIMA(p,d,q) if:

$$(1 - B)^d Y_t = W_t$$

where $\{W_t\}$ follows ARMA(p,q):

$$\phi(B)W_t = c + \theta(B)\epsilon_t$$

Complete ARIMA Equation:

$$\phi(B)(1 - B)^d Y_t = c + \theta(B)\epsilon_t$$

Parameters:

- p = order of autoregression
- d = degree of differencing
- q = order of moving average

Common ARIMA Models:

- ARIMA(0,1,0): Random walk with drift
- ARIMA(0,1,1): Exponential smoothing
- ARIMA(1,1,1): Damped exponential smoothing
- ARIMA(0,2,2): Linear trend model

Forecasting with ARIMA: After differencing d times, forecast using ARMA, then integrate d times to get original scale.

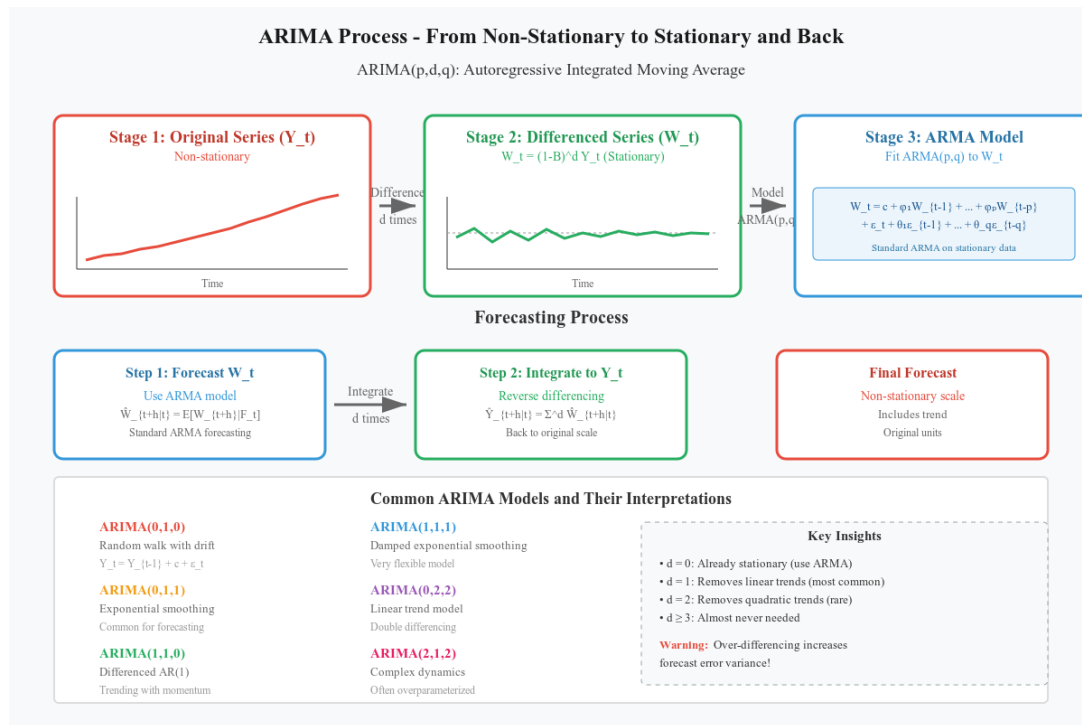


Figure 26: ARIMA Process - From Non-Stationary to Stationary and Back

Sumit (The Inquisitive)

Professor, I understand that d represents the number of differences, but how do we decide the right value of d ? Is there a test, or do we just keep differencing until the data looks stationary?

Excellent analytical question, Sumit! Choosing d is crucial - too little and data remains non-stationary, too much and we over-difference.

Systematic Approach to Choose d :

1. Sequential Testing:

- Test original series with ADF/KPSS
- If non-stationary, take first difference
- Test again; if still non-stationary, take second difference
- Rarely need $d \geq 2$ in practice

2. Visual Inspection:

- $d = 0$: Series fluctuates around fixed mean
- $d = 1$: Series has no trend after differencing
- $d = 2$: Series has no changing trend after differencing

3. Unit Root Tests at Each Stage:

- Original series: Test for one unit root
- If found, test first difference for unit root
- Continue until no unit roots remain

4. Minimize AIC/BIC:

- Fit models with $d = 0, 1, 2$
- Choose d that minimizes information criterion

Warning Signs of Over-Differencing:

- Variance increases after differencing
- ACF becomes strongly negative at lag 1 (≈ -0.5)
- Need many MA terms to model differenced series

3.5.2 Seasonal ARIMA - SARIMA Models

Theory Deep Dive

Seasonal ARIMA - SARIMA(p,d,q)(P,D,Q)

SARIMA models extend ARIMA to handle seasonal patterns.

Complete Model:

$$\phi(B)\Phi(B^s)(1-B)^d(1-B^s)^DY_t = c + \theta(B)\Theta(B^s)\epsilon_t$$

Where:

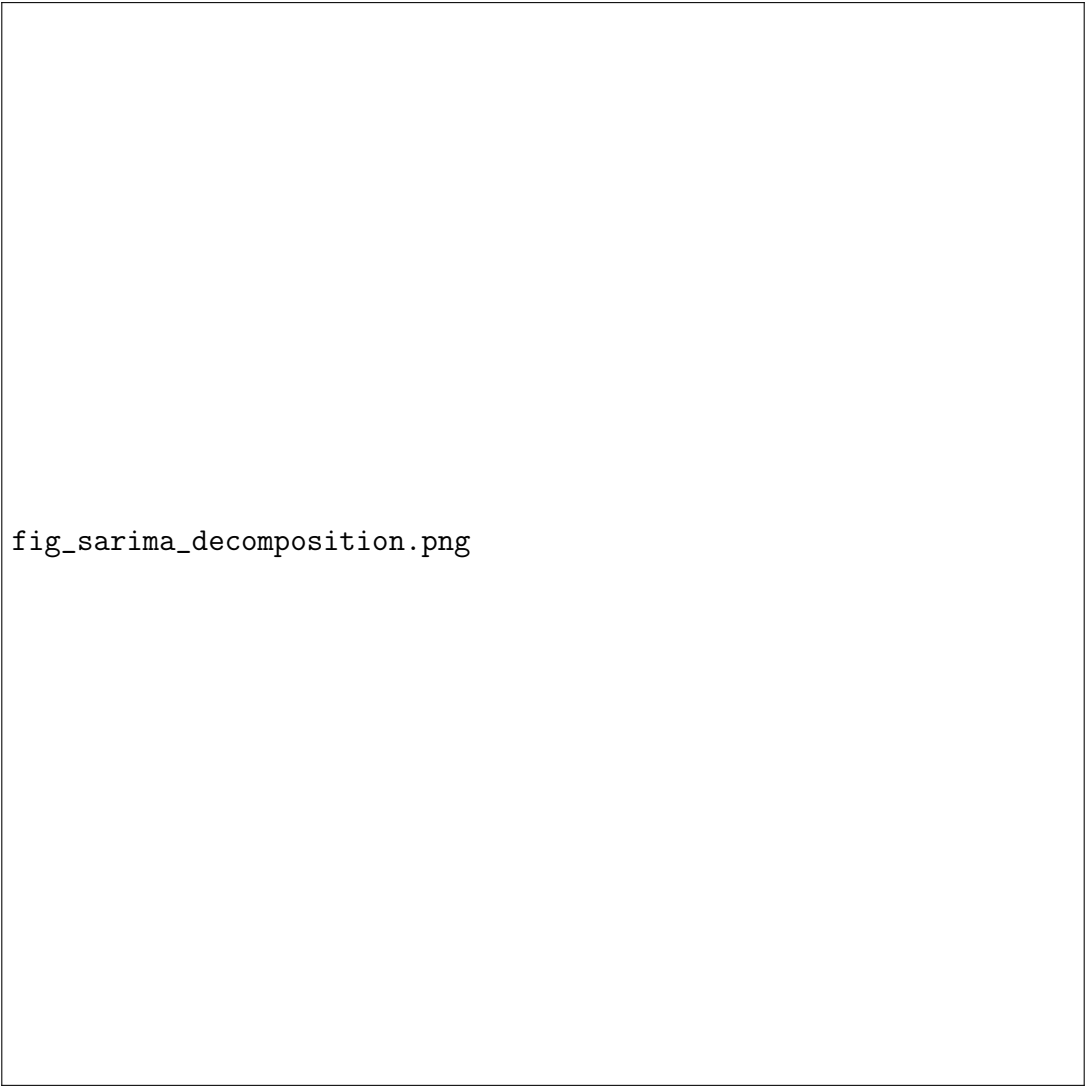
- $\phi(B)$ = non-seasonal AR polynomial of order p
- $\Phi(B^s)$ = seasonal AR polynomial of order P
- $\theta(B)$ = non-seasonal MA polynomial of order q
- $\Theta(B^s)$ = seasonal MA polynomial of order Q
- $(1-B)^d$ = non-seasonal differencing
- $(1-B^s)^D$ = seasonal differencing
- s = seasonal period (12 for monthly, 4 for quarterly)

Example - SARIMA(1,1,1)(1,1,1):

$$(1 - \phi B)(1 - \Phi B^{12})(1 - B)(1 - B^{12})Y_t = (1 + \theta B)(1 + \Theta B^{12})\epsilon_t$$

Common Seasonal Models:

- SARIMA(0,1,1)(0,1,1): Airline model
- SARIMA(1,0,0)(1,0,0): Seasonal AR
- SARIMA(0,0,1)(0,0,1): Quarterly seasonal MA



fig_sarima_decomposition.png

Figure 27: SARIMA Model Components - Regular and Seasonal Parts

Neha (The Skeptic)

Professor, SARIMA notation looks incredibly complex. With all these parameters (p, d, q, P, D, Q, s) , aren't we at risk of overfitting? How can we possibly choose all these values correctly?

Your skepticism is well-founded, Neha! SARIMA models can indeed become overly complex. Here's how we manage this complexity:

Practical SARIMA Modeling Strategy:

1. Start Simple:

- Begin with seasonal differencing if clear seasonality
- Often $(0,1,1)(0,1,1)$ works well (airline model)
- Add complexity only if needed

2. Separate Seasonal and Non-Seasonal:

- First, handle seasonality (choose D, P, Q)
- Then, model remaining non-seasonal part (choose d, p, q)
- This reduces the problem dimension

3. Use Domain Knowledge:

- Monthly data $\rightarrow s = 12$
- Quarterly data $\rightarrow s = 4$
- D rarely exceeds 1
- P and Q rarely exceed 2

4. Automated Selection:

- Use `auto.arima()` in R or `pmdarima` in Python
- These test many combinations efficiently
- Still verify the selected model makes sense!

5. Parsimony Principle:

- Prefer $\text{SARIMA}(1,1,1)(0,1,1)$ over $\text{SARIMA}(3,1,3)(2,1,2)$
- Simpler models generalize better
- Use AIC/BIC to balance fit and complexity

Remember: A simple model that captures main patterns beats a complex model that overfits!

Real-World Example

Real Example: Airline Passenger Data

The famous Box-Jenkins airline data led to the "airline model": SARIMA(0,1,1)(0,1,1)

Applied to log-transformed monthly passengers:

$$(1 - B)(1 - B^{12}) \log(Y_t) = (1 + \theta B)(1 + \Theta B^{12}) \epsilon_t$$

Estimates:

- $\theta = -0.40$ (non-seasonal MA)
- $\Theta = -0.62$ (seasonal MA)

This simple model captures:

- Exponential growth (via log transform)
- Annual seasonality (via seasonal differencing)
- Month-to-month variations (via MA terms)

Despite being developed in 1970, it still works well for many seasonal series!

4 Module 3: Forecasting and Model Selection

4.1 Forecasting in Time Series Models

Professor

Now we arrive at the ultimate goal of time series analysis: forecasting. All our work understanding components, achieving stationarity, and fitting models leads to this moment - predicting the future!

Forecasting is both an art and a science. The science gives us optimal predictions based on mathematical principles. The art lies in choosing the right model, understanding its limitations, and communicating uncertainty.

Let's master both aspects!

4.1.1 Minimum Mean Square Error Forecasting

Theory Deep Dive

Optimal Forecasting Theory

The optimal h-step-ahead forecast minimizes the mean square error (MSE).

Optimization Problem:

$$\hat{Y}_{t+h|t} = \arg \min_g E[(Y_{t+h} - g(\mathcal{F}_t))^2]$$

Where $\mathcal{F}_t = \{Y_t, Y_{t-1}, Y_{t-2}, \dots\}$ is the information set.

Solution:

$$\hat{Y}_{t+h|t} = E[Y_{t+h} | \mathcal{F}_t]$$

The optimal forecast is the conditional expectation!

Properties of Optimal Forecasts:

1. **Unbiased:** $E[\hat{Y}_{t+h|t}] = E[Y_{t+h}]$
2. **Minimum Variance:** Among all unbiased forecasts
3. **Orthogonal Errors:** $E[(Y_{t+h} - \hat{Y}_{t+h|t})Y_{t-j}] = 0$ for all $j \geq 0$

Forecast Error:

$$e_{t+h|t} = Y_{t+h} - \hat{Y}_{t+h|t}$$

Forecast Error Variance:

$$\sigma_h^2 = \text{Var}[e_{t+h|t}] = E[(Y_{t+h} - \hat{Y}_{t+h|t})^2]$$

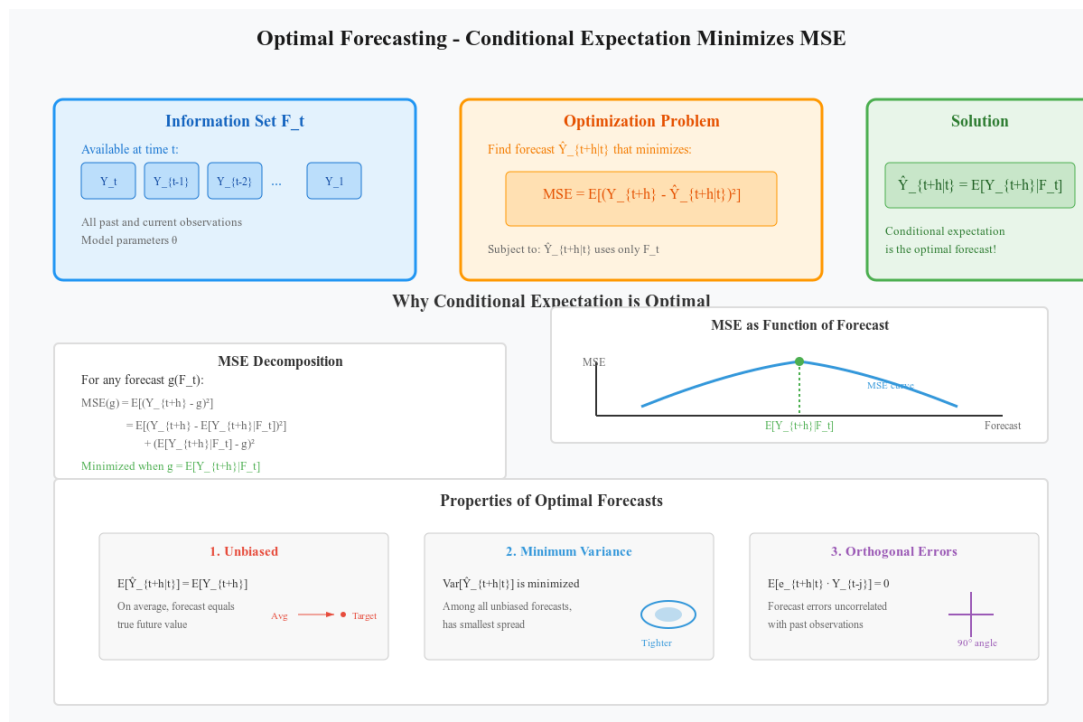


Figure 28: Optimal Forecasting - Conditional Expectation Minimizes MSE

4.1.2 Forecasting with AR Models

Theory Deep Dive

AR(p) Forecasting

For AR(p): $Y_t = c + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + \epsilon_t$

One-Step-Ahead Forecast:

$$\hat{Y}_{t+1|t} = c + \phi_1 Y_t + \phi_2 Y_{t-1} + \dots + \phi_p Y_{t-p+1}$$

Multi-Step-Ahead Forecast: For $h > p$:

$$\hat{Y}_{t+h|t} = c + \phi_1 \hat{Y}_{t+h-1|t} + \dots + \phi_p \hat{Y}_{t+h-p|t}$$

Forecast Error Variance:

$$\sigma_h^2 = \sigma^2 \sum_{j=0}^{h-1} \psi_j^2$$

Where ψ_j are the MA(∞) weights.

Forecast Properties:

- Forecasts converge to unconditional mean as $h \rightarrow \infty$
- Forecast intervals widen with horizon
- For stationary AR, long-run forecast = $\mu = c/(1 - \phi_1 - \dots - \phi_p)$

Paul (The Innocent)

Professor, I understand the formulas, but I'm confused about the recursive forecasting. When we forecast 2 steps ahead, do we use the actual value from 1 step ahead or our forecast?

Professor

Excellent question, Paul! This is a crucial distinction that many students miss.

The Golden Rule of Multi-Step Forecasting: "Use actual values when available, forecasts when not."

Example with AR(2): Model: $Y_t = 10 + 0.6Y_{t-1} + 0.3Y_{t-2} + \epsilon_t$

Current time: $t = 100$ Known values: $Y_{100} = 50$, $Y_{99} = 48$

One-step-ahead (forecasting Y_{101}):

$$\hat{Y}_{101|100} = 10 + 0.6(50) + 0.3(48) = 54.4$$

Uses actual values for both lags.

Two-step-ahead (forecasting Y_{102}):

$$\hat{Y}_{102|100} = 10 + 0.6(54.4) + 0.3(50) = 57.64$$

Uses forecast for Y_{101} and actual for Y_{100} .

Three-step-ahead (forecasting Y_{103}):

$$\hat{Y}_{103|100} = 10 + 0.6(57.64) + 0.3(54.4) = 60.90$$

Uses forecasts for both Y_{102} and Y_{101} .

This recursive approach propagates uncertainty forward!

4.1.3 Forecasting with MA and ARMA Models

Theory Deep Dive

MA(q) Forecasting

For MA(q): $Y_t = \mu + \epsilon_t + \theta_1\epsilon_{t-1} + \dots + \theta_q\epsilon_{t-q}$

Forecast Function:

$$\hat{Y}_{t+h|t} = \begin{cases} \mu + \theta_h\epsilon_t + \dots + \theta_q\epsilon_{t+h-q} & \text{if } h \leq q \\ \mu & \text{if } h > q \end{cases}$$

Key Property: MA(q) forecasts revert to mean after q periods!

ARMA(p,q) Forecasting

Combines AR and MA forecasting:

$$\hat{Y}_{t+h|t} = c + \sum_{i=1}^p \phi_i \hat{Y}_{t+h-i|t}^* + \sum_{j=h}^q \theta_j \epsilon_{t+h-j}$$

Where:

$$\hat{Y}_{t+h-i|t}^* = \begin{cases} Y_{t+h-i} & \text{if } h-i \leq 0 \\ \hat{Y}_{t+h-i|t} & \text{if } h-i > 0 \end{cases}$$

ARIMA Forecasting

For ARIMA(p,d,q):

1. Forecast differenced series $W_t = (1 - B)^d Y_t$ using ARMA
2. Integrate forecasts d times to get original scale
3. Adjust for differencing constants

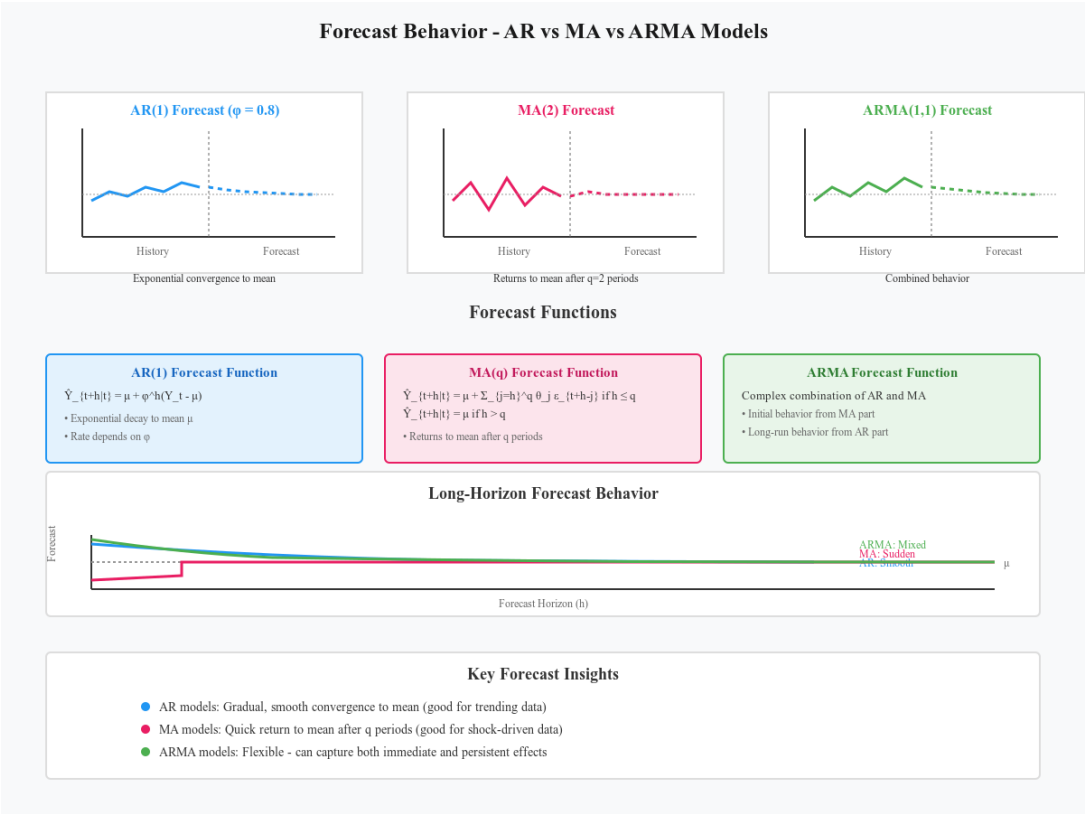


Figure 29: Forecast Behavior - AR vs MA vs ARMA Models

4.1.4 Prediction Intervals

Theory Deep Dive

Forecast Uncertainty Quantification

Point forecasts alone are insufficient - we need prediction intervals!

Theoretical Prediction Interval: Assuming normal errors, the $(1-\alpha)$ prediction interval is:

$$\hat{Y}_{t+h|t} \pm z_{\alpha/2} \sigma_h$$

Where:

- $z_{\alpha/2}$ = critical value from standard normal
- σ_h = forecast standard error at horizon h

Forecast Error Variance Growth: For any ARIMA model:

$$\sigma_h^2 = \sigma^2(1 + \psi_1^2 + \psi_2^2 + \dots + \psi_{h-1}^2)$$

Common Intervals:

- 95% PI: $\hat{Y}_{t+h|t} \pm 1.96\sigma_h$
- 80% PI: $\hat{Y}_{t+h|t} \pm 1.28\sigma_h$
- 50% PI: $\hat{Y}_{t+h|t} \pm 0.67\sigma_h$

Properties:

1. Intervals widen with forecast horizon
2. Width depends on model persistence
3. Assumes constant variance (check this!)

Rohan (The Visual Learner)

Professor, I see in the forecast plots that the prediction intervals look like funnels - narrow at first, then widening. Can you show me visually why this happens and how different models create different funnel shapes?

Professor

Excellent visual observation, Rohan! The "funnel" shape reveals deep truths about forecast uncertainty.

Why Intervals Widen - The Uncertainty Cascade:

Think of forecasting as a game of telephone:

- Step 1: Small uncertainty (just one error term)
- Step 2: Previous uncertainty + new uncertainty
- Step 3: Accumulated uncertainty + new uncertainty
- And so on...

Different Funnel Shapes:

1. MA(1) - Bounded Funnel:

- Widens for 1 period only
- Then stays constant
- Uncertainty has a ceiling

2. AR(1) - Gradual Funnel:

- Widens gradually
- Rate depends on ϕ
- Approaches limit as $h \rightarrow \infty$

3. Random Walk - Linear Funnel:

- Widens indefinitely
- Variance grows linearly with h
- No upper bound on uncertainty

4. ARIMA with $d \geq 0$ - Explosive Funnel:

- Widens rapidly
- Integrated processes have unbounded variance
- Long-term forecasts very uncertain

The funnel shape tells you about the model's "memory" and predictability!

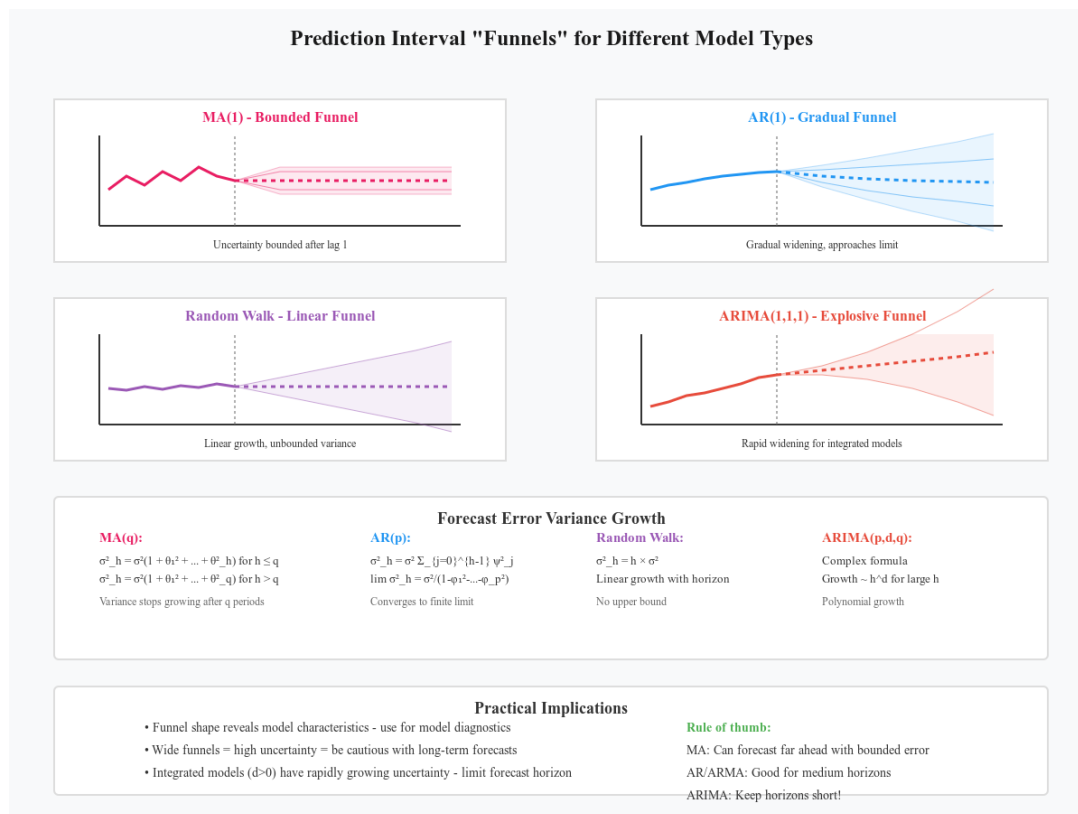


Figure 30: Prediction Interval "Funnels" for Different Model Types

4.2 Durbin-Levinson Algorithm

Professor

The Durbin-Levinson algorithm is an elegant recursive method for solving the Yule-Walker equations. It's not just a computational tool - it provides deep insights into the structure of time series models and connects to the partial autocorrelation function we studied earlier.

The Durbin-Levinson Algorithm

Problem: Given autocorrelations $\rho_1, \rho_2, \dots, \rho_p$, find AR(p) coefficients.

Yule-Walker Equations:

$$\begin{bmatrix} 1 & \rho_1 & \rho_2 & \cdots & \rho_{p-1} \\ \rho_1 & 1 & \rho_1 & \cdots & \rho_{p-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_{p-1} & \rho_{p-2} & \rho_{p-3} & \cdots & 1 \end{bmatrix} \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_p \end{bmatrix} = \begin{bmatrix} \rho_1 \\ \rho_2 \\ \vdots \\ \rho_p \end{bmatrix}$$

Algorithm:

1. Initialize: $\phi_{11} = \rho_1$, $v_0 = 1$, $v_1 = 1 - \rho_1^2$
2. For $k = 2, 3, \dots, p$:

$$\phi_{kk} = \frac{\rho_k - \sum_{j=1}^{k-1} \phi_{k-1,j} \rho_{k-j}}{v_{k-1}} \quad (10)$$

$$\phi_{kj} = \phi_{k-1,j} - \phi_{kk} \phi_{k-1,k-j} \text{ for } j = 1, \dots, k-1 \quad (11)$$

$$v_k = v_{k-1}(1 - \phi_{kk}^2) \quad (12)$$

Output:

- AR(p) coefficients: $\phi_{p1}, \phi_{p2}, \dots, \phi_{pp}$
- Partial autocorrelations: $\phi_{11}, \phi_{22}, \dots, \phi_{pp}$
- Prediction error variances: v_1, v_2, \dots, v_p

Sumit (The Inquisitive)

Professor, the algorithm looks very mechanical. What's the intuition behind these recursive formulas? And why do the diagonal elements give us the PACF?

Professor

Brilliant question, Sumit! Let me reveal the beautiful intuition behind Durbin-Levinson.

The Deep Intuition:

Think of fitting increasingly complex AR models:

- AR(1): Best linear predictor using just Y_{t-1}
- AR(2): Best linear predictor using Y_{t-1} and Y_{t-2}
- AR(k): Best linear predictor using Y_{t-1}, \dots, Y_{t-k}

The Recursive Magic: When going from AR(k-1) to AR(k), we ask: "How much does adding Y_{t-k} improve our prediction?"

The answer is ϕ_{kk} - the PACF at lag k!

Why Diagonal = PACF:

- ϕ_{11} : Direct effect of Y_{t-1} on Y_t
- ϕ_{22} : Additional effect of Y_{t-2} after accounting for Y_{t-1}
- ϕ_{kk} : Additional effect of Y_{t-k} after accounting for all intermediate lags

Prediction Error Reduction: $v_k = v_{k-1}(1 - \phi_{kk}^2)$

This shows:

- Each lag reduces prediction variance by factor $(1 - \phi_{kk}^2)$
- If $\phi_{kk} = 0$, no improvement (stop here!)
- If $|\phi_{kk}|$ large, significant improvement

The algorithm builds models incrementally, keeping what works!

Python Code Reference: 07_durbin_levinson.py

Implementation of Durbin-Levinson:

```

def durbin_levinson(acf, p):
    """
    Durbin-Levinson algorithm
    Returns: AR coefficients, PACF, prediction variances
    """
    phi = np.zeros((p+1, p+1))
    v = np.zeros(p+1)

    # Initialize
    phi[1,1] = acf[1]
    v[0] = 1
    v[1] = 1 - acf[1]**2

    # Recursion
    for k in range(2, p+1):
        # Calculate PACF
        num = acf[k] - sum(phi[k-1,j] * acf[k-j]
                           for j in range(1,k))
        phi[k,k] = num / v[k-1]

        # Update coefficients
        for j in range(1, k):
            phi[k,j] = phi[k-1,j] - phi[k,k] * phi[k-1,k-j]

        # Update variance
        v[k] = v[k-1] * (1 - phi[k,k]**2)

    return phi[p,1:p+1], [phi[i,i] for i in range(1,p+1)], v

```

4.3 Parameter Estimation Methods

Professor

Having identified the model structure, we now face the crucial task of estimating parameters. The quality of our forecasts depends critically on accurate parameter estimation. Let's explore the main methods, their properties, and when to use each.

4.3.1 Method of Moments Estimation

Theory Deep Dive

Method of Moments (MOM) Estimation

MOM matches theoretical moments with sample moments.

For AR(p) Models: Use Yule-Walker equations with sample autocorrelations:

$$\hat{\phi} = \hat{\mathbf{R}}^{-1} \hat{\rho}$$

Where:

- $\hat{\mathbf{R}}$ = sample autocorrelation matrix
- $\hat{\rho} = [\hat{\rho}_1, \hat{\rho}_2, \dots, \hat{\rho}_p]'$

For MA(q) Models: Solve nonlinear equations:

$$\rho_k(\theta) = \hat{\rho}_k \text{ for } k = 1, 2, \dots, q$$

For ARMA(p,q) Models:

1. Estimate AR part from extended Yule-Walker
2. Estimate MA part from residuals

Properties:

- Always produces stationary/invertible estimates
- Computationally simple for AR
- Inefficient for finite samples
- Complex for MA and ARMA

4.3.2 Least Squares Estimation

Theory Deep Dive

Least Squares (LS) Estimation

Conditional Least Squares (CLS): Minimize conditional sum of squares:

$$S(\boldsymbol{\theta}) = \sum_{t=m+1}^n (Y_t - \hat{Y}_{t|t-1})^2$$

Where $m = \max(p, q)$ and $\hat{Y}_{t|t-1}$ is one-step forecast.

For AR(p): Linear regression problem:

$$\hat{\boldsymbol{\phi}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Where:

$$\mathbf{X} = \begin{bmatrix} Y_p & Y_{p-1} & \cdots & Y_1 \\ Y_{p+1} & Y_p & \cdots & Y_2 \\ \vdots & \vdots & \ddots & \vdots \\ Y_{n-1} & Y_{n-2} & \cdots & Y_{n-p} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} Y_{p+1} \\ Y_{p+2} \\ \vdots \\ Y_n \end{bmatrix}$$

Properties:

- Asymptotically equivalent to MLE
- Simple computation for AR
- May produce non-stationary estimates
- Ignores initial observations

4.3.3 Maximum Likelihood Estimation

Theory Deep Dive

Maximum Likelihood Estimation (MLE)

Assumes Gaussian errors and maximizes likelihood function.

Log-Likelihood:

$$\ell(\boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Gamma}| - \frac{1}{2} (\mathbf{Y} - \boldsymbol{\mu})' \boldsymbol{\Gamma}^{-1} (\mathbf{Y} - \boldsymbol{\mu})$$

Where $\boldsymbol{\Gamma}$ is the covariance matrix of $\mathbf{Y} = [Y_1, \dots, Y_n]'$.

Innovations Algorithm: For computational efficiency, use prediction error decomposition:

$$\ell(\boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{t=1}^n \log v_{t-1} - \frac{1}{2} \sum_{t=1}^n \frac{e_t^2}{v_{t-1}}$$

Where:

- $e_t = Y_t - \hat{Y}_{t|t-1}$ = prediction error
- $v_{t-1} = \text{Var}[Y_t | Y_1, \dots, Y_{t-1}]$ = prediction variance

Properties:

- Asymptotically efficient
- Provides standard errors
- Handles all observations optimally
- Computationally intensive
- May converge to local maxima

Neha (The Skeptic)

Professor, with three different estimation methods, I'm skeptical about which one to trust. They seem to give different results on the same data. How do we know which estimates are "correct"?

Professor

Your skepticism is justified, Neha! Different methods can indeed give different estimates, especially in small samples. Here's how to think about it:

No Single "Correct" Answer: Each method makes different trade-offs:

- MOM: Prioritizes matching correlations
- LS: Prioritizes prediction accuracy
- MLE: Prioritizes probabilistic fit

When Methods Agree:

- Large samples ($n \geq 200$):
- Simple models (AR(1), MA(1))
- Well-behaved data (no outliers)

When to Use Each Method:

Use MOM when:

- Quick initial estimates needed
- Guaranteed stationarity required
- Working with AR models only

Use LS when:

- Focusing on forecast accuracy
- Computational simplicity needed
- Sample size is moderate to large

Use MLE when:

- Need optimal properties
- Constructing confidence intervals
- Sample size is small
- Final model estimation

Practical Strategy:

1. Start with MOM for initial values
2. Refine with MLE
3. Compare results - large differences indicate model uncertainty
4. Use bootstrap to assess estimation uncertainty

Python Code Reference

Python Code Reference: 08_parameter_estimation.py

Comparing estimation methods:

```
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.ar_model import AutoReg

# Method of Moments (via Yule-Walker)
ar_mom = AutoReg(data, lags=2, trend='c')
res_mom = ar_mom.fit(method='yw')

# Least Squares
ar_ls = AutoReg(data, lags=2, trend='c')
res_ls = ar_ls.fit(method='ols')

# Maximum Likelihood
model_mle = ARIMA(data, order=(2,0,0))
res_mle = model_mle.fit(method='mle')

# Compare estimates
print("MOM:", res_mom.params)
print("LS:", res_ls.params)
print("MLE:", res_mle.params)
```

4.4 Model Selection and Comparison

Professor

One of the most challenging aspects of time series modeling is choosing the right model. Should we use AR(2) or AR(3)? Is ARMA(1,1) better than MA(2)? Model selection is about finding the sweet spot between underfitting (missing important patterns) and overfitting (fitting noise). Let's explore the scientific criteria that guide these decisions.

4.4.1 Information Criteria

Theory Deep Dive

Information Criteria for Model Selection

Information criteria balance goodness-of-fit against model complexity.

Akaike Information Criterion (AIC):

$$AIC = -2 \log L(\hat{\theta}) + 2k$$

Where:

- $L(\hat{\theta})$ = maximized likelihood
- k = number of parameters

For normal errors:

$$AIC = n \log(\hat{\sigma}^2) + 2k$$

Bayesian Information Criterion (BIC):

$$BIC = -2 \log L(\hat{\theta}) + k \log(n)$$

Or: $BIC = n \log(\hat{\sigma}^2) + k \log(n)$

Corrected AIC (AICc): For small samples:

$$AICc = AIC + \frac{2k(k+1)}{n-k-1}$$

Final Prediction Error (FPE):

$$FPE = \hat{\sigma}^2 \cdot \frac{n+k}{n-k}$$

Interpretation:

- Lower values indicate better models
- Compare only models fit to same data
- Differences ≥ 2 are significant

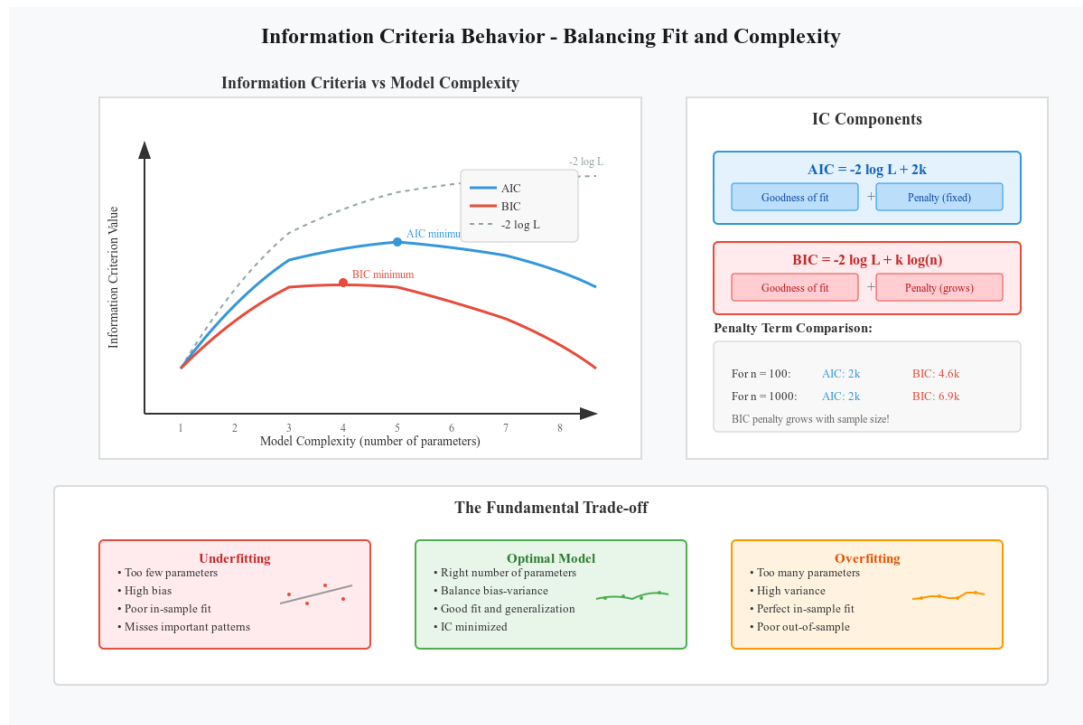


Figure 31: Information Criteria Behavior - Balancing Fit and Complexity

Paul (The Innocent)

Professor, I see that AIC and BIC have different penalty terms ($2k$ vs $k \log n$). Which one should I use? Do they ever disagree on the best model?

Professor

Excellent observation, Paul! AIC and BIC often disagree, and understanding why helps you choose appropriately.

The Fundamental Difference:

- AIC penalty: $2k$ (fixed)
- BIC penalty: $k \log(n)$ (grows with sample size)

For $n \geq 8$, BIC penalty $>$ AIC penalty, so BIC chooses simpler models.

When They Disagree: Example scenario:

- Model A: AR(2), AIC = 100, BIC = 110
- Model B: AR(4), AIC = 98, BIC = 115
- AIC chooses Model B (lower AIC)
- BIC chooses Model A (lower BIC)

Which to Use:

Use AIC when:

- Focusing on prediction accuracy
- Sample size is small/moderate
- Concerned about underfitting
- Doing short-term forecasting

Use BIC when:

- Seeking true model
- Sample size is large
- Concerned about overfitting
- Need parsimonious model

Practical Advice:

- Check both - if they agree, decision is clear
- If they disagree, consider your goals
- Large disagreement suggests model uncertainty
- Consider ensemble of top models

4.4.2 Model Selection Strategy

Theory Deep Dive

Box-Jenkins Model Selection Procedure

Step 1: Identification

- Plot data and check stationarity
- Determine differencing order (d)
- Examine ACF and PACF patterns
- Propose candidate models

Step 2: Estimation

- Estimate parameters for each candidate
- Check convergence and standard errors
- Verify stationarity/invertibility

Step 3: Diagnostic Checking

- Analyze residuals
- Test for white noise
- Check for remaining patterns

Step 4: Selection

- Compare information criteria
- Consider parsimony
- Evaluate forecast performance

Modern Automated Approach:

1. Set maximum orders: p_{max}, q_{max}
2. Fit all combinations where $p \leq p_{max}, q \leq q_{max}$
3. Select model minimizing chosen criterion
4. Verify with diagnostic checks

Important Warning

Common Model Selection Pitfalls

1. **Over-reliance on automated selection:** Always verify the selected model makes sense
2. **Ignoring domain knowledge:** Physical constraints should guide selection
3. **Fitting too many models:** Increases chance of spurious selection
4. **Using same data for selection and evaluation:** Leads to overoptimistic assessment
5. **Ignoring seasonal patterns:** Check for seasonality before ARIMA selection

Brain Teaser

Model Selection Practice

Given these model fits for the same dataset (n=200):

Model	Parameters	Log-Like	AIC	BIC	RMSE
AR(1)	2	-250.3	504.6	511.2	1.45
AR(2)	3	-248.1	502.2	512.1	1.42
MA(1)	2	-252.7	509.4	516.0	1.48
MA(2)	3	-249.2	504.4	514.3	1.43
ARMA(1,1)	3	-247.8	501.6	511.5	1.41

Questions:

1. Which model does AIC select?
2. Which model does BIC select?
3. Which is most parsimonious with good fit?
4. Calculate AICc for AR(1):

Answers: ARMA(1,1), AR(1), AR(2), $504.6 + 6/197 = 504.63$

4.5 Residual Analysis and Diagnostic Checking

Professor

After fitting a model, our job isn't done! We must verify that our model adequately captures the data's patterns. Residual analysis is like a detective investigating whether we've missed any clues. If the model is correct, residuals should be white noise - random, unpredictable, and patternless.

4.5.1 Properties of Good Residuals

Theory Deep Dive

Residual Properties for Adequate Models

If the model is correctly specified, residuals should satisfy:

1. Zero Mean:

$$E[\hat{\epsilon}_t] = 0$$

Test: One-sample t-test

2. Constant Variance (Homoscedasticity):

$$Var[\hat{\epsilon}_t] = \sigma^2 \text{ for all } t$$

Test: Goldfeld-Quandt, Breusch-Pagan

3. No Autocorrelation:

$$Corr[\hat{\epsilon}_t, \hat{\epsilon}_{t-k}] = 0 \text{ for all } k > 0$$

Test: Ljung-Box, ACF plots

4. Normality:

$$\hat{\epsilon}_t \sim N(0, \sigma^2)$$

Test: Shapiro-Wilk, QQ plots

5. No Patterns: Residuals vs fitted values should show no patterns Test: Visual inspection, runs test



Figure 32: Residual Diagnostic Plots - What to Look For

4.5.2 Ljung-Box Test for Residuals

Theory Deep Dive

Ljung-Box Test on Residuals

Tests whether residual autocorrelations are jointly zero.

Test Statistic:

$$Q_{LB} = n(n+2) \sum_{k=1}^h \frac{\hat{\rho}_k^2(\hat{\epsilon})}{n-k}$$

Distribution: Under H_0 : $Q_{LB} \sim \chi_{h-p-q}^2$

Note: Degrees of freedom reduced by number of estimated parameters!

Modified Test for Small Samples: Use $h = \min(10, n/5)$ for better size properties

Interpretation:

- p-value ≥ 0.05 : Residuals are white noise (good!)
- p-value < 0.05 : Residuals have structure (model inadequate)

Testing Strategy:

1. Test at multiple lags ($h = 5, 10, 15, 20$)
2. All should pass for model adequacy
3. If fail at specific lag, check ACF at that lag

Rohan (The Visual Learner)

Professor, I'm looking at the residual plots, but I'm not sure what patterns I should be worried about. Could you show me examples of "bad" residual patterns and what they indicate?

Excellent visual question, Rohan! Let me show you the common problematic patterns and their diagnoses.

Pattern 1: Trending Residuals

- Appearance: Residuals drift up or down
- Diagnosis: Missing trend component
- Fix: Add differencing or trend term

Pattern 2: Cyclical Residuals

- Appearance: Sine-wave pattern
- Diagnosis: Missing seasonal/cyclical component
- Fix: Add seasonal terms or use SARIMA

Pattern 3: Changing Variance

- Appearance: Funnel shape in residual plot
- Diagnosis: Heteroscedasticity
- Fix: Transform data (log, Box-Cox) or use GARCH

Pattern 4: Autocorrelated Residuals

- Appearance: Runs of positive/negative residuals
- Diagnosis: Model order too low
- Fix: Increase p or q

Pattern 5: Outliers

- Appearance: Isolated extreme residuals
- Diagnosis: Unusual events or model misspecification
- Fix: Investigate outliers, consider robust methods

Pattern 6: Nonlinearity

- Appearance: Curved pattern vs fitted values
- Diagnosis: Linear model inappropriate
- Fix: Transform variables or use nonlinear models

Remember: Perfect residuals are rare - look for major violations!

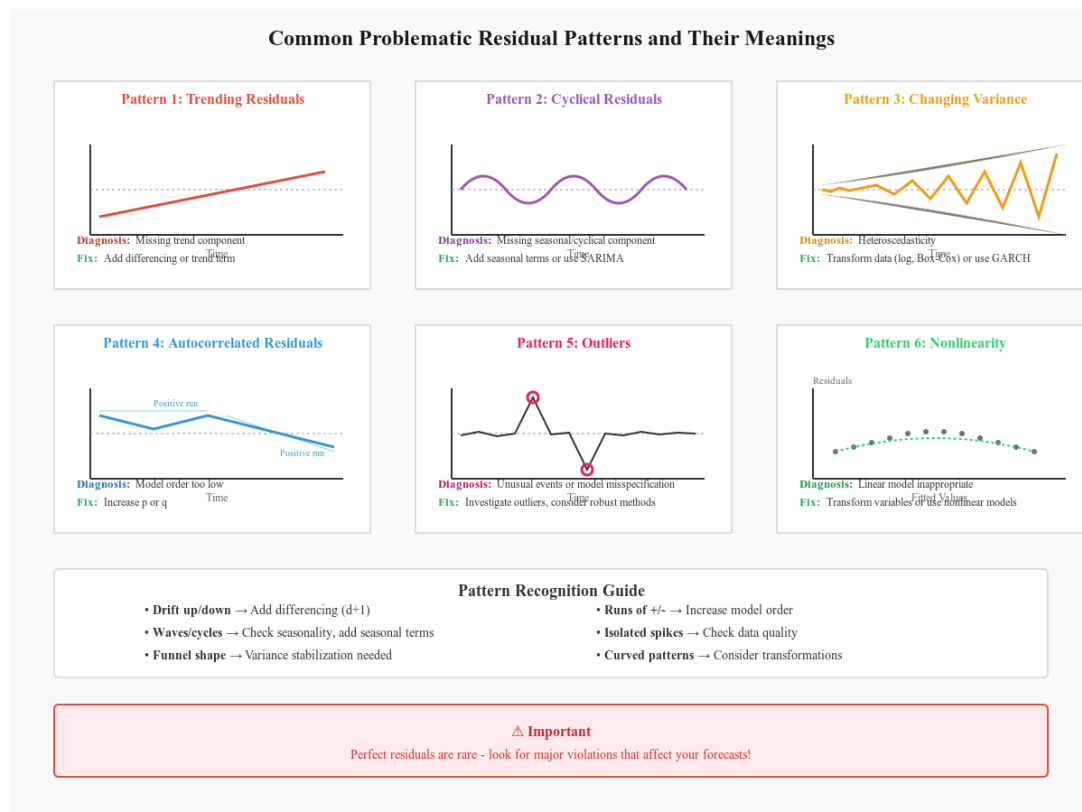


Figure 33: Common Problematic Residual Patterns and Their Meanings

4.5.3 Comprehensive Diagnostic Strategy

Theory Deep Dive

Complete Model Diagnostic Checklist

1. Residual Plots:

- Time series plot of residuals
- ACF and PACF of residuals
- Residuals vs fitted values
- QQ plot for normality
- Histogram of residuals

2. Statistical Tests:

- Ljung-Box at multiple lags
- Shapiro-Wilk for normality
- ARCH-LM for heteroscedasticity
- Runs test for randomness

3. Influence Diagnostics:

- Check for influential observations
- Examine parameter stability
- Look for structural breaks

4. Out-of-Sample Validation:

- Reserve last 10-20
- Compare forecast errors
- Check prediction interval coverage

Python Code Reference: 09_residual_diagnostics.py

Comprehensive residual diagnostics:

```
import matplotlib.pyplot as plt
from statsmodels.stats.diagnostic import acorr_ljungbox
from scipy import stats

def comprehensive_diagnostics(residuals, p, q):
    """Complete residual diagnostic suite"""

    # Create figure with subplots
    fig, axes = plt.subplots(2, 3, figsize=(15, 10))

    # 1. Time series plot
    axes[0,0].plot(residuals)
    axes[0,0].set_title('Residual Series')
    axes[0,0].axhline(y=0, color='r', linestyle='--')

    # 2. ACF plot
    plot_acf(residuals, ax=axes[0,1], lags=20)
    axes[0,1].set_title('Residual ACF')

    # 3. QQ plot
    stats.probplot(residuals, dist="norm", plot=axes[0,2])
    axes[0,2].set_title('Normal QQ Plot')

    # 4. Histogram
    axes[1,0].hist(residuals, bins=20, density=True)
    axes[1,0].set_title('Residual Distribution')

    # 5. Residuals vs fitted
    axes[1,1].scatter(fitted_values, residuals, alpha=0.5)
    axes[1,1].axhline(y=0, color='r', linestyle='--')
    axes[1,1].set_title('Residuals vs Fitted')

    # 6. Ljung-Box test results
    lb_test = acorr_ljungbox(residuals, lags=20,
                             model_df=p+q, return_df=True)
    axes[1,2].plot(lb_test.index, lb_test['lb_pvalue'])
    axes[1,2].axhline(y=0.05, color='r', linestyle='--')
    axes[1,2].set_title('Ljung-Box p-values')

    plt.tight_layout()
    return fig, lb_test
```

4.6 Unit Root Tests

Professor

Unit root tests are the gatekeepers of time series modeling. They determine whether we need differencing (ARIMA) or can use stationary models (ARMA). Beyond the ADF test we've seen, let's explore the full toolkit of unit root tests and their strategic use.

4.6.1 Phillips-Perron Test

Theory Deep Dive

Phillips-Perron (PP) Test

A non-parametric alternative to ADF that handles serial correlation and heteroscedasticity.

Test Regression (same as ADF):

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

Key Difference: Instead of adding lags, PP modifies test statistics:

$$Z_t = t_\beta \left(\frac{\gamma_0}{s^2} \right)^{1/2} - \frac{T(s^2 - \gamma_0)}{2s\sqrt{T^{-2} \sum Y_{t-1}^2}}$$

Where:

- γ_0 = consistent estimate of error variance
- s^2 = Newey-West long-run variance estimate

Advantages over ADF:

- No need to specify lag length
- Robust to heteroscedasticity
- Often more powerful

Disadvantages:

- Size distortions in small samples
- Sensitive to structural breaks

4.6.2 KPSS Test Revisited

Theory Deep Dive

Kwiatkowski-Phillips-Schmidt-Shin (KPSS) Test

Tests null hypothesis of stationarity (opposite of ADF/PP).

Test Statistic:

$$KPSS = \frac{1}{T^2} \frac{\sum_{t=1}^T S_t^2}{\hat{\sigma}^2}$$

Where $S_t = \sum_{i=1}^t \hat{e}_i$ (partial sum of residuals)

Interpretation Strategy:

ADF/PP	KPSS	Conclusion
Reject H_0	Fail to reject H_0	Series is stationary
Fail to reject H_0	Reject H_0	Series has unit root
Reject H_0	Reject H_0	Series is trend stationary
Fail to reject H_0	Fail to reject H_0	Data is inconclusive

Sumit (The Inquisitive)

Professor, I find it confusing that ADF and KPSS have opposite null hypotheses. Why would statisticians design tests this way, and how do I reconcile conflicting results?

Professor

Brilliant question, Sumit! The opposite null hypotheses are actually a feature, not a bug. Let me explain the deep statistical philosophy here.

The Power Problem:

- ADF has low power against near-unit-root alternatives
- Example: $\phi = 0.95$ is stationary but ADF often fails to reject
- KPSS has low power against near-stationary alternatives
- Example: $\phi = 1.02$ has unit root but KPSS might not reject

Complementary Testing Strategy: Think of it like a criminal trial:

- ADF: "Innocent (stationary) until proven guilty (unit root)"
- KPSS: "Guilty (unit root) until proven innocent (stationary)"
- Need strong evidence from both perspectives!

Reconciling Conflicts:

Both reject their nulls:

- Likely trend-stationary
- Solution: Include trend in model or detrend first

Neither rejects:

- Near unit root ($\phi \approx 1$)
- Solution: Try both $I(0)$ and $I(1)$ models, compare forecasts

ADF rejects, KPSS doesn't:

- Clear evidence of stationarity
- Solution: Use ARMA models

ADF doesn't reject, KPSS does:

- Clear evidence of unit root
- Solution: Use ARIMA with $d = 1$

This dual testing approach gives more robust conclusions!

4.6.3 Testing for Multiple Unit Roots

Theory Deep Dive

Testing for I(2) vs I(1)

Some series need differencing twice (I(2)).

Sequential Testing Procedure:

1. Test Y_t for unit root
2. If unit root found, test ΔY_t for unit root
3. If second unit root found, series is I(2)

Common I(2) Series:

- Nominal price levels (not logs)
- Cumulative totals
- Some monetary aggregates

Warning: Over-differencing is more common than under-differencing!

Brain Teaser

Unit Root Testing Practice

Given these test results for a financial time series:

Series	ADF p-value	KPSS p-value	PP p-value
Level Y_t	0.35	0.02	0.41
First diff ΔY_t	0.001	0.45	0.001
Second diff $\Delta^2 Y_t$	0.0001	0.50	0.0001

Questions:

1. Is the original series stationary? _____
2. What is the order of integration I(d)? _____
3. Should we use ARMA or ARIMA? _____
4. What value of d for ARIMA? _____

Answers: No, I(1), ARIMA, d=1

5 Module 4: Advanced Time Series Models

5.1 Multivariate Time Series and VAR Models

Professor

Welcome to the frontier of time series analysis! So far, we've focused on single variables evolving through time. But in the real world, economic and financial variables don't exist in isolation - they influence each other in complex webs of relationships.

Think about it: Interest rates affect exchange rates, which affect inflation, which affects interest rates again. Multivariate time series models, particularly Vector Autoregression (VAR), help us understand these interconnected systems.

5.1.1 Introduction to Multivariate Time Series

Theory Deep Dive

Multivariate Time Series Fundamentals

A multivariate time series is a collection of multiple time series observed simultaneously.

Notation:

$$\mathbf{Y}_t = \begin{bmatrix} Y_{1t} \\ Y_{2t} \\ \vdots \\ Y_{kt} \end{bmatrix}$$

Where each Y_{it} is a univariate time series.

Cross-Correlation Function:

$$\rho_{ij}(h) = \frac{\gamma_{ij}(h)}{\sqrt{\gamma_{ii}(0)\gamma_{jj}(0)}}$$

Where $\gamma_{ij}(h) = Cov(Y_{it}, Y_{jt-h})$

Key Concepts:

1. **Contemporaneous correlation:** Correlation at same time
2. **Lead-lag relationships:** One series predicts another
3. **Feedback effects:** Bidirectional causality
4. **Common shocks:** External factors affecting all series

Stationarity for Multivariate Series:

- $E[\mathbf{Y}_t] = \boldsymbol{\mu}$ (constant mean vector)
- $Cov(\mathbf{Y}_t, \mathbf{Y}_{t-h}) = \boldsymbol{\Gamma}(h)$ (depends only on lag)

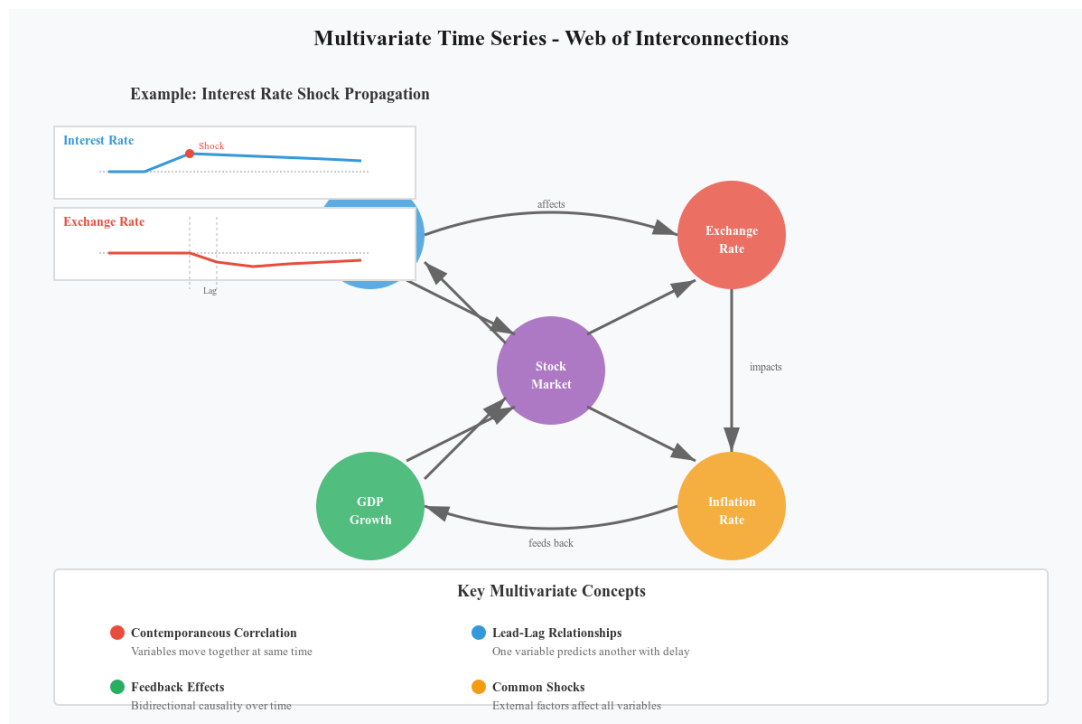


Figure 34: Multivariate Time Series - Web of Interconnections

5.1.2 Vector Autoregression (VAR) Models

Theory Deep Dive

VAR(p) Model Definition

A VAR(p) model expresses each variable as a linear function of past values of all variables:

$$\mathbf{Y}_t = \mathbf{c} + \mathbf{A}_1 \mathbf{Y}_{t-1} + \mathbf{A}_2 \mathbf{Y}_{t-2} + \dots + \mathbf{A}_p \mathbf{Y}_{t-p} + \boldsymbol{\epsilon}_t$$

Where:

- $\mathbf{Y}_t = k \times 1$ vector of variables
- $\mathbf{c} = k \times 1$ vector of constants
- $\mathbf{A}_i = k \times k$ coefficient matrices
- $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}, \boldsymbol{\Sigma}) = k \times 1$ error vector

VAR(1) Example with 2 variables:

$$\begin{bmatrix} Y_{1t} \\ Y_{2t} \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} Y_{1,t-1} \\ Y_{2,t-1} \end{bmatrix} + \begin{bmatrix} \epsilon_{1t} \\ \epsilon_{2t} \end{bmatrix}$$

Expanded:

$$Y_{1t} = c_1 + a_{11}Y_{1,t-1} + a_{12}Y_{2,t-1} + \epsilon_{1t} \quad (13)$$

$$Y_{2t} = c_2 + a_{21}Y_{1,t-1} + a_{22}Y_{2,t-1} + \epsilon_{2t} \quad (14)$$

Key Properties:

- Each equation has same regressors
- Allows for rich dynamics
- Can capture feedback effects
- OLS estimation equation by equation

Neha (The Skeptic)

Professor, VAR seems like we're just running multiple regressions. Why is this considered special? And with so many parameters (k^2p), aren't we guaranteed to overfit?

Your skepticism highlights crucial issues, Neha! Let me address both concerns.

Why VAR is Special:

1. System Perspective:

- Not just multiple regressions - it's a system
- Captures dynamic interdependencies
- Shocks propagate through the system
- Can trace impulse responses

2. Unique Capabilities:

- Granger causality testing
- Impulse response functions
- Forecast error variance decomposition
- No need to specify exogenous/endogenous

The Curse of Dimensionality:

You're right - parameters grow as k^2p :

- 2 variables, lag 2: 8 slope parameters
- 5 variables, lag 4: 100 slope parameters!
- 10 variables, lag 4: 400 slope parameters!!

Solutions to Overfitting:

1. **Keep models small:** 2-5 variables typical
2. **Use information criteria:** BIC penalizes complexity
3. **Bayesian VAR:** Shrinkage priors
4. **Factor models:** Reduce dimensionality first
5. **Sparse VAR:** LASSO penalties

The art is balancing richness with parsimony!

5.1.3 Stability and Stationarity of VAR

Theory Deep Dive

VAR Stability Conditions

Companion Form: Rewrite VAR(p) as VAR(1) in companion form:

$$\mathbf{Z}_t = \mathbf{F}\mathbf{Z}_{t-1} + \mathbf{V}_t$$

Where:

$$\mathbf{Z}_t = \begin{bmatrix} \mathbf{Y}_t \\ \mathbf{Y}_{t-1} \\ \vdots \\ \mathbf{Y}_{t-p+1} \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_{p-1} & \mathbf{A}_p \\ \mathbf{I}_k & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{I}_k & \mathbf{0} \end{bmatrix}$$

Stability Condition: All eigenvalues of \mathbf{F} must be inside unit circle:

$$|\lambda_i| < 1 \text{ for all } i$$

Alternative Check: Roots of $\det(\mathbf{I}_k - \mathbf{A}_1 z - \cdots - \mathbf{A}_p z^p) = 0$ outside unit circle.

Consequences of Instability:

- Explosive behavior
- Non-stationary system
- Invalid inference
- Meaningless impulse responses

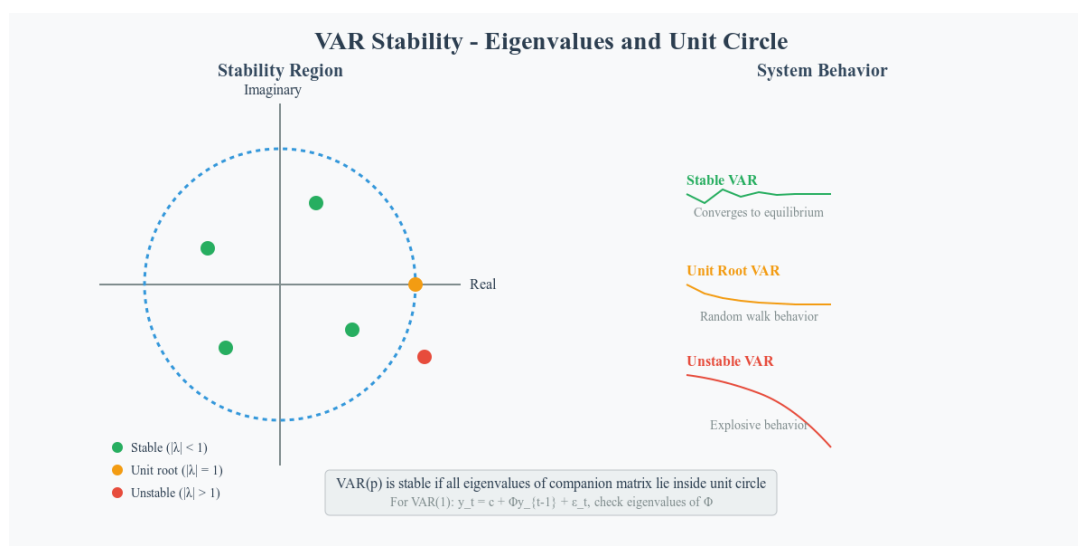


Figure 35: VAR Stability - Eigenvalues and Unit Circle

5.1.4 Impulse Response Functions

Theory Deep Dive

Impulse Response Function (IRF)

IRF traces the effect of a one-time shock to one variable on all variables over time.

Mathematical Definition: The IRF at horizon h is:

$$\frac{\partial \mathbf{Y}_{t+h}}{\partial \boldsymbol{\epsilon}'_t} = \boldsymbol{\Psi}_h$$

Moving Average Representation: VAR can be written as VMA(∞):

$$\mathbf{Y}_t = \boldsymbol{\mu} + \sum_{i=0}^{\infty} \boldsymbol{\Psi}_i \boldsymbol{\epsilon}_{t-i}$$

Where $\boldsymbol{\Psi}_0 = \mathbf{I}_k$ and $\boldsymbol{\Psi}_i$ are obtained recursively:

$$\boldsymbol{\Psi}_s = \sum_{j=1}^s \boldsymbol{\Psi}_{s-j} \mathbf{A}_j$$

Orthogonalized IRF: Since $\boldsymbol{\epsilon}_t$ are correlated, use Cholesky decomposition:

$$\boldsymbol{\Sigma} = \mathbf{P}\mathbf{P}'$$

Then orthogonalized shocks: $\mathbf{u}_t = \mathbf{P}^{-1}\boldsymbol{\epsilon}_t$

Interpretation:

- Shows dynamic effects of shocks
- Reveals transmission mechanisms
- Quantifies persistence
- Tests economic theories

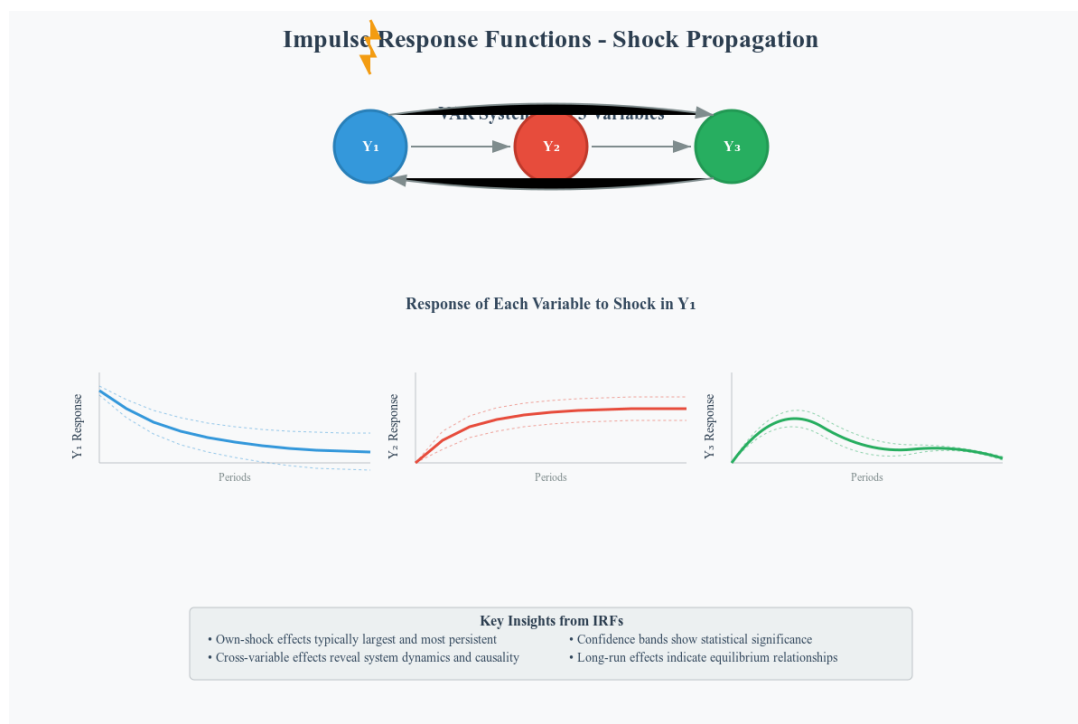


Figure 36: Impulse Response Functions - Shock Propagation Through System

Paul (The Innocent)

Professor, I'm confused about orthogonalized shocks. Why can't we just shock one variable at a time in the original system?

Excellent question, Paul! This gets at a fundamental issue in VAR analysis.

The Correlation Problem:

In real data, errors are contemporaneously correlated:

$$Corr(\epsilon_{1t}, \epsilon_{2t}) \neq 0$$

This means:

- A "shock" to variable 1 naturally comes with a shock to variable 2
- Can't have one without the other in the data
- Like trying to increase temperature without affecting pressure

Example - Interest Rates and Exchange Rates:

- Errors are correlated (= -0.6)
- When interest rates get positive shock, exchange rates typically get negative shock
- Orthogonalization asks: "What if we could shock interest rates alone?"

Cholesky Solution:

- Transforms correlated shocks into uncorrelated ones
- Assumes causal ordering (variable 1 affects 2 contemporaneously, not vice versa)
- Order matters! Different orderings give different IRFs

Alternative Approaches:

- Structural VAR: Use economic theory to identify shocks
- Sign restrictions: Identify shocks by their effects
- External instruments: Use outside information

Think of orthogonalization as creating "laboratory conditions" where we can study pure effects!

5.1.5 Granger Causality

Theory Deep Dive

Granger Causality Testing

Tests whether past values of one variable help predict another.

Definition: Y_2 does not Granger-cause Y_1 if:

$$E[Y_{1t}|\mathcal{I}_{t-1}] = E[Y_{1t}|\mathcal{I}_{t-1} \setminus \{Y_{2,t-j}, j \geq 1\}]$$

Where \mathcal{I}_{t-1} is all information up to time $t - 1$.

Testing in VAR: Test H_0 : Y_2 does not Granger-cause Y_1

In VAR equation for Y_1 :

$$Y_{1t} = c_1 + \sum_{j=1}^p a_{11,j} Y_{1,t-j} + \sum_{j=1}^p a_{12,j} Y_{2,t-j} + \epsilon_{1t}$$

Test: $H_0 : a_{12,1} = a_{12,2} = \dots = a_{12,p} = 0$

F-Test Statistic:

$$F = \frac{(RSS_r - RSS_u)/p}{RSS_u/(T - 2p - 1)}$$

Where:

- RSS_r = restricted model (without Y_2 lags)
- RSS_u = unrestricted model (full VAR)

Important Notes:

- Granger causality \neq true causality
- Tests predictive power, not structural relationships
- Sensitive to omitted variables
- Requires correct lag specification

Sumit (The Inquisitive)

Professor, I've always been confused by the term "Granger causality." If it's not real causality, why call it causality at all? Can you give an example where Granger causality exists but true causality doesn't?

Your confusion is justified, Sumit! The terminology is indeed misleading. Let me clarify with examples.

What Granger Causality Really Means: "X Granger-causes Y" = "Past X contains information useful for predicting Y beyond what's in past Y"

Classic Example - Rooster and Sunrise:

- Rooster crows Granger-cause sunrise? YES
- Rooster crows truly cause sunrise? NO
- Why? Rooster anticipates sunrise, creating predictive power

Financial Example - Stock Prices and News:

- Stock prices often Granger-cause news announcements
- Do price changes cause news? NO
- Why? Insider information leaks into prices before announcements

Economic Example - Leading Indicators:

- Building permits Granger-cause GDP growth
- Do permits directly cause all GDP? NO
- Why? Permits are early signal of future construction activity

When Granger Causality True Causality:

- Controlled experiments
- Clear temporal precedence
- No confounding variables
- Strong theoretical backing

Better Terminology Would Be: "Granger predictability" or "Granger precedence"

But we're stuck with the historical term!

Brain Teaser

VAR Model Interpretation Challenge

Given this estimated VAR(1) model for GDP growth (g) and inflation (π):

$$\begin{bmatrix} g_t \\ \pi_t \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.0 \end{bmatrix} + \begin{bmatrix} 0.6 & -0.2 \\ 0.1 & 0.8 \end{bmatrix} \begin{bmatrix} g_{t-1} \\ \pi_{t-1} \end{bmatrix} + \begin{bmatrix} \epsilon_{1t} \\ \epsilon_{2t} \end{bmatrix}$$

Questions:

1. Does inflation Granger-cause GDP growth? _____
2. Does GDP growth Granger-cause inflation? _____
3. What's the long-run mean of GDP growth? _____
4. Is the system stable? (Hint: eigenvalues are $0.7 \pm 0.1i$) _____

Answers: Yes (-0.20), Yes (0.10), 2.0, Yes (—j1)

5.2 Vector ARMA Models

Professor

While VAR models are workhorses of multivariate analysis, they can require many parameters. Vector ARMA (VARMA) models extend the parsimony of univariate ARMA to multiple series, often achieving similar fit with fewer parameters.

Theory Deep Dive

VARMA(p,q) Model

The VARMA(p,q) model combines vector autoregression and vector moving average:

$$\mathbf{Y}_t = \mathbf{c} + \sum_{i=1}^p \mathbf{A}_i \mathbf{Y}_{t-i} + \boldsymbol{\epsilon}_t + \sum_{j=1}^q \mathbf{M}_j \boldsymbol{\epsilon}_{t-j}$$

Or using lag polynomials:

$$\mathbf{A}(B)\mathbf{Y}_t = \mathbf{c} + \mathbf{M}(B)\boldsymbol{\epsilon}_t$$

Where:

- $\mathbf{A}(B) = \mathbf{I} - \mathbf{A}_1 B - \dots - \mathbf{A}_p B^p$
- $\mathbf{M}(B) = \mathbf{I} + \mathbf{M}_1 B + \dots + \mathbf{M}_q B^q$

Identification Issues:

- More complex than univariate case
- Multiple representations for same model
- Need restrictions for uniqueness
- Often use Kronecker indices

When to Use VARMA:

- VAR requires too many lags
- Strong MA behavior in residuals
- Parsimony is crucial
- Theoretical reasons for MA structure

5.3 Conditional Heteroscedastic Models

Professor

Traditional time series models assume constant variance. But financial data shouts otherwise! Stock returns show volatility clustering - periods of high volatility followed by high volatility, calm followed by calm.

ARCH and GARCH models revolutionized finance by modeling this changing variance. They're not about predicting returns - they're about predicting risk!

5.3.1 ARCH Models

Theory Deep Dive

Autoregressive Conditional Heteroscedasticity - ARCH(q)

Models time-varying variance as function of past squared shocks.

Model Structure:

$$\begin{aligned} Y_t &= \mu_t + \epsilon_t \\ \epsilon_t &= \sigma_t z_t, \quad z_t \sim N(0, 1) \\ \sigma_t^2 &= \omega + \alpha_1 \epsilon_{t-1}^2 + \alpha_2 \epsilon_{t-2}^2 + \dots + \alpha_q \epsilon_{t-q}^2 \end{aligned}$$

Where:

- $\mu_t = E[Y_t | \mathcal{F}_{t-1}]$ (conditional mean)
- $\sigma_t^2 = Var[Y_t | \mathcal{F}_{t-1}]$ (conditional variance)
- z_t = standardized innovation

ARCH(1) Example:

$$\sigma_t^2 = \omega + \alpha_1 \epsilon_{t-1}^2$$

Parameter Restrictions:

- $\omega > 0$ (positive variance)
- $\alpha_i \geq 0$ for all i (non-negative weights)
- $\sum_{i=1}^q \alpha_i < 1$ (stationarity)

Properties:

- $E[\epsilon_t] = 0$
- $Var[\epsilon_t] = \frac{\omega}{1 - \sum \alpha_i}$ (unconditional)
- Captures volatility clustering
- Fat-tailed distributions

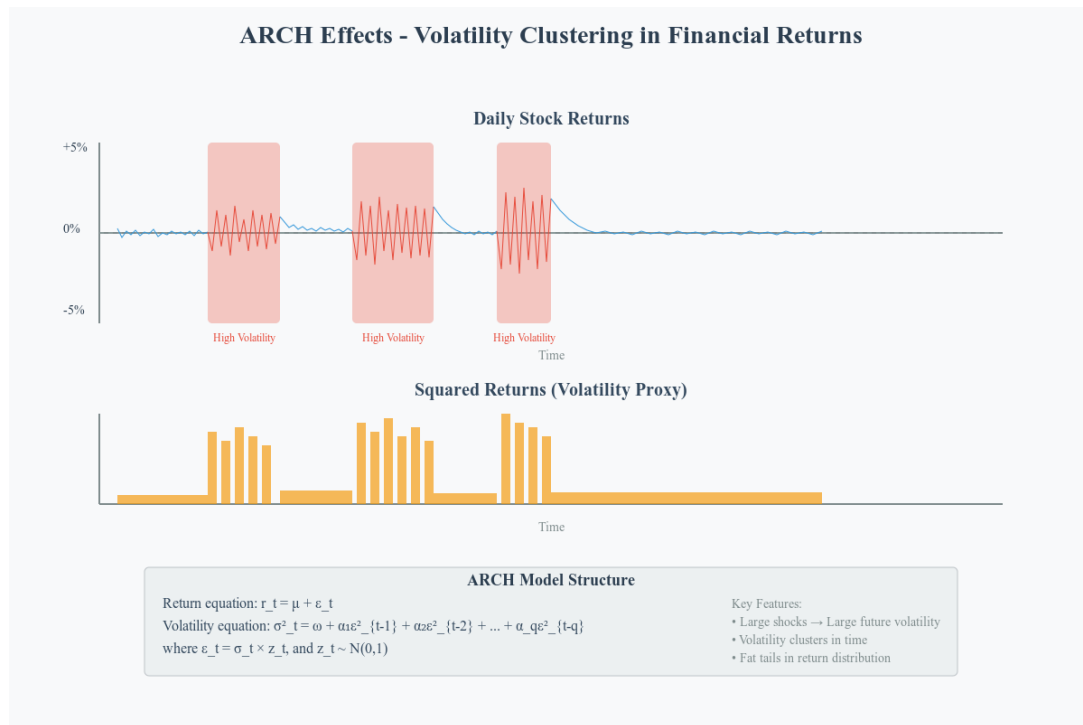


Figure 37: ARCH Effects - Volatility Clustering in Financial Returns

Rohan (The Visual Learner)

Professor, I can see the volatility clustering in the plot, but I don't understand why we square the past errors. Why not use absolute values or some other function?

Great visual insight, Rohan! The choice of squaring is both practical and theoretical.

Why Squared Errors:

1. Statistical Properties:

- Squared errors = variance (natural measure of spread)
- Differentiable (needed for optimization)
- Connects to normal distribution theory

2. Symmetry:

- Positive and negative shocks have same effect
- $(-2)^2 = (+2)^2 = 4$
- Makes sense for many financial applications

3. Mathematical Convenience:

- Closed-form solutions possible
- Links to chi-squared distributions
- Easier statistical inference

Alternative Specifications:

- **Absolute value:** Used in some models (harder math)
- **Asymmetric:** EGARCH, GJR-GARCH (different effects for + vs -)
- **Power ARCH:** $\sigma_t^\delta = \omega + \alpha |\epsilon_{t-1}|^\delta$

Intuition: Think of variance as "energy" in the system:

- Big shock (± 5) \rightarrow High energy (25) \rightarrow Expects more volatility
- Small shock (± 0.5) \rightarrow Low energy (0.25) \rightarrow Expects calm
- Energy dissipates over time unless refreshed by new shocks

Squaring amplifies large shocks, which matches how markets react!

5.3.2 GARCH Models

Theory Deep Dive

Generalized ARCH - GARCH(p,q)

GARCH adds lagged conditional variances to ARCH, achieving parsimony.

Model Structure:

$$\sigma_t^2 = \omega + \sum_{i=1}^q \alpha_i \epsilon_{t-i}^2 + \sum_{j=1}^p \beta_j \sigma_{t-j}^2$$

GARCH(1,1) - The Workhorse:

$$\sigma_t^2 = \omega + \alpha_1 \epsilon_{t-1}^2 + \beta_1 \sigma_{t-1}^2$$

Parameter Restrictions:

- $\omega > 0$
- $\alpha_i \geq 0, \beta_j \geq 0$
- $\sum_{i=1}^q \alpha_i + \sum_{j=1}^p \beta_j < 1$ (stationarity)

Persistence: Total persistence = $\sum \alpha_i + \sum \beta_j$

- Close to 1: High persistence (volatility shocks last long)
- Close to 0: Low persistence (quick mean reversion)
- = 1: IGARCH (integrated GARCH)

Unconditional Variance:

$$E[\epsilon_t^2] = \frac{\omega}{1 - \sum \alpha_i - \sum \beta_j}$$

Why GARCH(1,1) Often Suffices:

- Captures persistence parsimoniously
- One parameter for immediate impact (α_1)
- One parameter for persistence (β_1)
- Fits most financial data well

Neha (The Skeptic)

Professor, I'm skeptical about GARCH models. They seem to assume that only the magnitude of past shocks matters, not their direction. But in real markets, don't negative shocks (bad news) create more volatility than positive shocks?

Brilliant observation, Neha! You've identified the key limitation of standard GARCH - the symmetry assumption. Real markets absolutely show asymmetric responses!

The Leverage Effect:

- Bad news (negative returns) → Higher volatility
- Good news (positive returns) → Lower volatility increase
- Named because falling prices increase financial leverage

Asymmetric GARCH Models:

1. EGARCH (Exponential GARCH):

$$\log(\sigma_t^2) = \omega + \alpha \left[\frac{|\epsilon_{t-1}|}{\sigma_{t-1}} - \sqrt{\frac{2}{\pi}} \right] + \gamma \frac{\epsilon_{t-1}}{\sigma_{t-1}} + \beta \log(\sigma_{t-1}^2)$$

- $\gamma < 0$ captures leverage effect
- No restrictions needed (log ensures positivity)

2. GJR-GARCH (Threshold GARCH):

$$\sigma_t^2 = \omega + (\alpha + \gamma I_{t-1}) \epsilon_{t-1}^2 + \beta \sigma_{t-1}^2$$

Where $I_{t-1} = 1$ if $\epsilon_{t-1} < 0$

- $\gamma < 0$ means negative shocks have larger impact
- More intuitive than EGARCH

Empirical Evidence:

- Stock indices: Strong leverage effect (significant)
- Exchange rates: Weak/no leverage effect
- Commodities: Mixed results

Your skepticism is well-founded - symmetric GARCH is often too simple for equity markets!

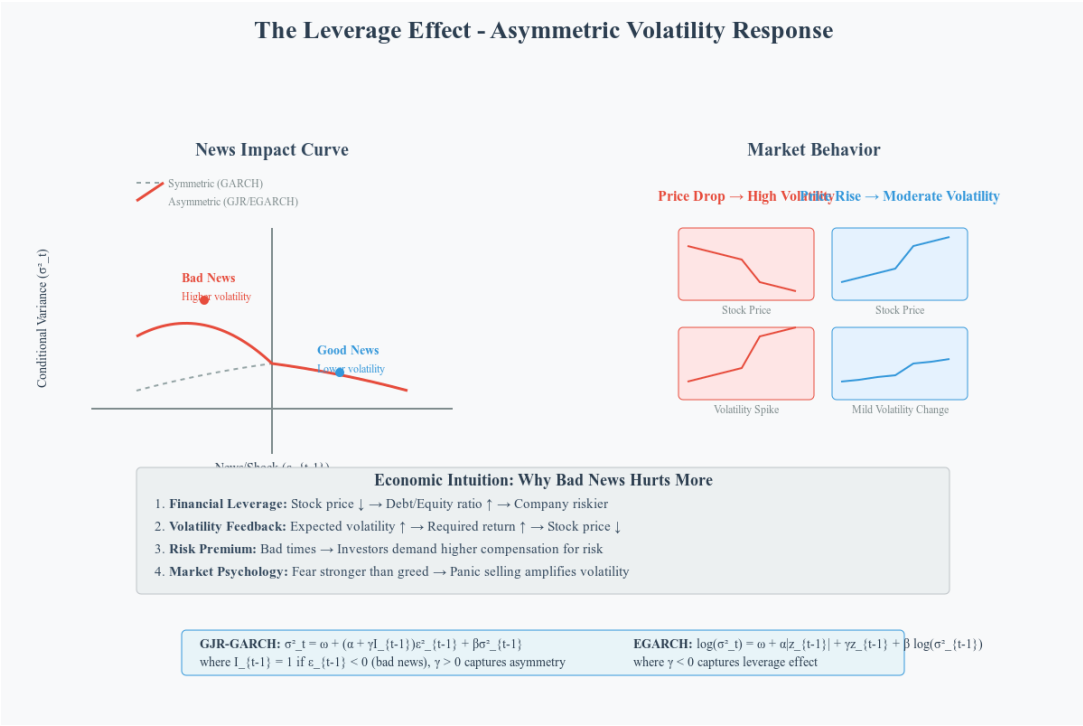


Figure 38: The Leverage Effect - Asymmetric Volatility Response

5.3.3 GARCH Model Estimation

Theory Deep Dive

Maximum Likelihood Estimation for GARCH

Log-Likelihood Function:

$$\ell(\boldsymbol{\theta}) = -\frac{1}{2} \sum_{t=1}^T \left[\log(2\pi) + \log(\sigma_t^2) + \frac{\epsilon_t^2}{\sigma_t^2} \right]$$

Where $\boldsymbol{\theta} = (\omega, \alpha_1, \dots, \alpha_q, \beta_1, \dots, \beta_p)'$

Estimation Procedure:

1. Estimate mean equation (AR, MA, ARMA, etc.)
2. Calculate residuals $\hat{\epsilon}_t$
3. Initialize σ_0^2 (often sample variance)
4. Maximize likelihood numerically
5. Check parameter constraints

Quasi-Maximum Likelihood: Even if errors aren't normal, QML estimates are consistent if:

- Model correctly specifies conditional mean and variance
- Use robust standard errors (Bollerslev-Wooldridge)

Model Diagnostics:

- Standardized residuals: $\tilde{z}_t = \epsilon_t / \sigma_t$
- Should be iid $N(0,1)$
- Test with Ljung-Box on \tilde{z}_t and \tilde{z}_t^2
- ARCH-LM test on standardized residuals

5.3.4 GARCH Forecasting

Theory Deep Dive

Volatility Forecasting with GARCH

One-Step-Ahead Forecast: For GARCH(1,1):

$$\hat{\sigma}_{t+1|t}^2 = \omega + \alpha \epsilon_t^2 + \beta \sigma_t^2$$

Multi-Step-Ahead Forecast:

$$\hat{\sigma}_{t+h|t}^2 = \omega \sum_{i=0}^{h-2} (\alpha + \beta)^i + (\alpha + \beta)^{h-1} \hat{\sigma}_{t+1|t}^2$$

As $h \rightarrow \infty$:

$$\hat{\sigma}_{t+h|t}^2 \rightarrow \frac{\omega}{1 - \alpha - \beta} = \text{unconditional variance}$$

Value-at-Risk (VaR) Application:

$$VaR_{t+1|t}^{(p)} = \mu_{t+1|t} + z_p \hat{\sigma}_{t+1|t}$$

Where z_p is the p-th quantile of return distribution.

Forecast Evaluation:

- Mincer-Zarnowitz regression
- Realized volatility comparison
- VaR backtesting
- Diebold-Mariano tests

Paul (The Innocent)

Professor, I understand GARCH forecasts volatility, but how do traders actually use these forecasts? Can you give a concrete example?

Professor

Excellent practical question, Paul! Let me show you exactly how GARCH transforms from academic model to trading tool.

Example: Options Trading

Suppose you're trading Nifty options on Friday:

- Current Nifty: 18,000
- GARCH forecast: $\sigma = 1.5$
- Normal times: $\sigma = 0.8$

Trading Decisions:

1. Option Pricing:

- Black-Scholes needs volatility input
- Market implying 1.2
- Your GARCH says 1.5
- Conclusion: Options underpriced, buy volatility!

2. Position Sizing:

- Normal volatility: Risk 100,000 per position
- High volatility period: Reduce to 65,000
- Maintains constant risk despite changing volatility

3. Stop-Loss Placement:

- Normal: Stop at $2\sigma = 1.6$
- High volatility: Stop at $2\sigma = 3.0$
- Avoids premature exits during volatile periods

4. Portfolio Risk Management:

- Calculate portfolio VaR using GARCH
- If VaR exceeds limit, reduce positions
- Dynamically adjust to market conditions

Real Trader's Rule: "When GARCH volatility doubles, halve your position size!"

Brain Teaser

GARCH Model Challenge

Given GARCH(1,1) estimates for a stock:

- $\omega = 0.00002$
- $\alpha_1 = 0.05$
- $\beta_1 = 0.93$
- Current variance: $\sigma_t^2 = 0.0004$
- Current shock: $\epsilon_t = -0.03$ (3

Calculate:

1. Persistence measure: _____
2. Unconditional variance: _____
3. Next period variance forecast: _____
4. Is this high or low persistence? _____

Answers: 0.98, 0.001, 0.000442, High persistence

Python Code Reference

Python Code Reference: 10_garch_modeling.py

GARCH implementation:

```
from arch import arch_model
import numpy as np

# Simulate returns data
returns = np.random.normal(0, 1, 1000)

# Fit GARCH(1,1)
model = arch_model(returns, vol='Garch', p=1, q=1)
results = model.fit(disp='off')

# Print results
print(results.summary())

# Forecast volatility
forecasts = results.forecast(horizon=5)
print(f"Volatility forecasts: {np.sqrt(forecasts.variance.values[-1])}")

# Plot conditional volatility
results.plot()
plt.show()
```

6 Course Summary and Integration

Professor

Congratulations! You've journeyed through the complete landscape of time series analysis. From basic decomposition to advanced GARCH models, you now possess a powerful toolkit for understanding and predicting temporal data. Let's integrate everything we've learned and see how these pieces fit together in real-world applications.



Figure 39: Complete Time Series Analysis Framework - From Data to Decisions

Integrated Time Series Analysis Framework

1. Exploratory Analysis:

- Decomposition (Module 1)
- ACF/PACF analysis
- Stationarity testing

2. Model Selection:

- Univariate: ARIMA/SARIMA (Module 2)
- Multivariate: VAR/VARMA (Module 4)
- Volatility: ARCH/GARCH (Module 4)

3. Estimation Diagnostics:

- Parameter estimation (Module 3)
- Residual analysis
- Model validation

4. Forecasting Decision Making:

- Point and interval forecasts
- Risk assessment
- Strategic applications

Real-World Example

Complete Real-World Example: Forex Trading Strategy Let's apply our entire toolkit to USD/INR exchange rate: **Step 1: Data Exploration**

- Daily data shows trend (depreciation)
- No seasonality (daily financial data)
- Volatility clustering evident

Step 2: Stationarity Transformation

- Level: Non-stationary (ADF p-value = 0.82)
- Returns: Stationary (ADF p-value = 0.001)
- Decision: Model returns, not levels

Step 3: Model Selection

- Mean equation: ARMA(1,1) based on AIC
- Variance equation: GARCH(1,1) for volatility
- Combined model: ARMA(1,1)-GARCH(1,1)

Step 4: Results Application

- Mean forecast: 0.02
- Volatility forecast: 0.6
- Trading rule: Go long USD when forecast ≥ 0.5
- Risk management: Position size inversely proportional to GARCH volatility

Paul (The Innocent)

Professor, this course covered so much! How do I remember when to use which technique? Is there a decision tree or checklist I can follow?

Professor

Excellent question, Paul! Let me give you a practical decision framework that will guide you through any time series problem. **The Master Decision Tree:**

- 1. How many series?**

- One → Univariate methods (ARIMA)
- Multiple → Multivariate methods (VAR)

- 2. Is variance constant?**

- Yes → Standard models
- No → Add GARCH component

- 3. Any seasonality?**

- Yes → Seasonal models (SARIMA)
- No → Non-seasonal models

- 4. Is it stationary?**

- Yes → ARMA/VAR
- No → Difference first (ARIMA)

Quick Reference Checklist:

1. Plot the data
2. Test stationarity
3. Check seasonality
4. Examine ACF/PACF
5. Fit candidate models
6. Diagnose residuals
7. Compare criteria
8. Validate forecasts

Remember: Start simple, add complexity only when needed!

Rohan (The Visual Learner)

Professor, as a visual learner, I'd love to see how all these models relate to each other. Is there a visual map of the time series universe?

Professor

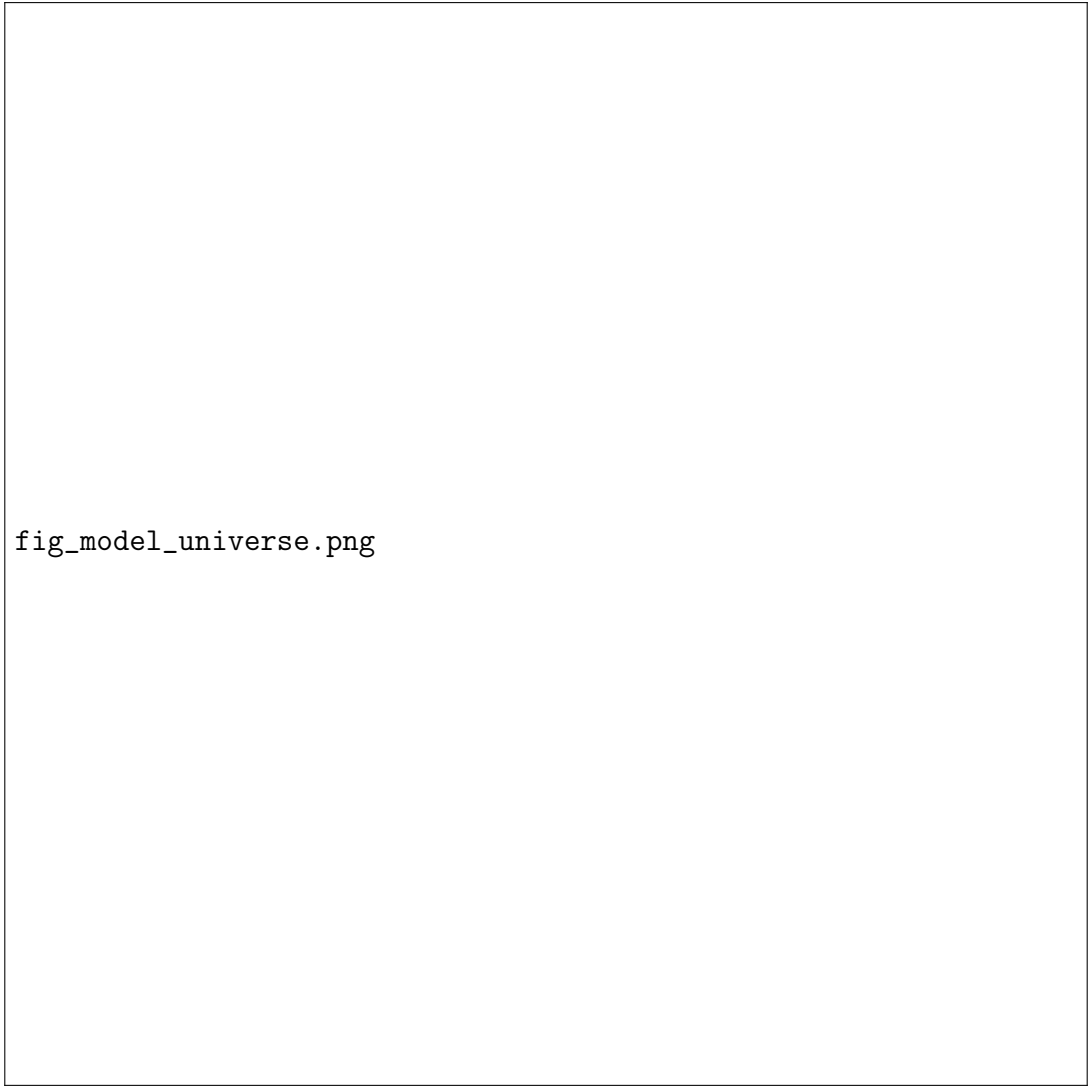
Perfect final question, Rohan! Let me show you the complete visual taxonomy of everything we've learned. **The Time Series Model Universe:** Think of it as a tree:

- **Root:** General linear process
- **Branches:** Different restrictions give different models
- **Leaves:** Specific applications

Key Relationships:

- AR() MA() (under invertibility)
- ARMA ARIMA (d=0 case)
- VAR VARMA (q=0 case)
- All models \rightarrow ARCH/GARCH can be added

This visual map is your compass in the time series world!



fig_model_universe.png

Figure 40: The Complete Time Series Model Universe - A Visual Taxonomy

Neha (The Skeptic)

Professor, one final skeptical question: With all these sophisticated models, are we really any better at predicting the future than simple methods? What's the evidence that this complexity is worth it?

Professor

A fitting final question, Neha! Your skepticism touches on a fundamental issue in forecasting. **The Honest Truth: When Simple Models Win:**

- Short horizons (1-2 periods)
- Noisy data
- Structural breaks
- Limited data

When Complex Models Shine:

- Rich patterns (seasonality, cycles)
- Long, stable history
- Multiple related series
- Risk measurement (GARCH)

Empirical Evidence:

- M-Competitions: Simple methods often competitive
- But: Complex models better at capturing uncertainty
- Domain-specific: Finance needs GARCH, economics needs VAR

The Wisdom: "All models are wrong, but some are useful" - George Box **My Advice:**

1. Start simple (naive, MA, exponential smoothing)
2. Add complexity incrementally
3. Always benchmark against simple methods
4. Focus on confidence intervals, not just point forecasts
5. Remember: Understanding the process \hat{y} mechanical forecasting

The journey through these models hasn't just taught you techniques - it's taught you to think probabilistically about time and uncertainty. That's the real value!

Sumit (The Inquisitive)

Professor, thank you for this incredible journey! One last question: What's next? Where should we go from here to deepen our time series expertise?

What a wonderful question to end on, Sumit! Your curiosity will take you far. Here's your roadmap forward: **Advanced Topics to Explore:** **1. Modern Methods:**

- State space models and Kalman filtering
- Machine learning for time series (LSTM, Prophet)
- Functional time series
- High-frequency data analysis

2. Specialized Domains:

- Financial econometrics (stochastic volatility, copulas)
- Macroeconometrics (DSGE models)
- Spatial-temporal models
- Bayesian time series

3. Practical Skills:

- Real-time forecasting systems
- Big data time series
- Ensemble methods
- Forecast combination

Resources:

- Books: Hamilton, Tsay, Lütkepohl
- Journals: Journal of Time Series Analysis, IJF
- Communities: ISF, time series conferences
- Practice: Kaggle competitions, real data projects

Final Wisdom: "The best way to learn time series is to do time series. Find data you're passionate about and start forecasting!" Thank you all for being such engaged students. May your forecasts be accurate and your residuals white noise!

"Time is the wisest counselor of all." - Pericles "In time series, the past whispers secrets about the future. Our job is to listen carefully."