

Chapter 1

ARIMA

Amit Dua

Put abstract here with \abstract.

Change with \papernote

Chapter 2

Time Series Analytics

Abstract

This chapter presents time series analysis theory and practice using YouTube, SEO, and NSE data. It provides implementations in both Excel and Python.

1 Foundations of Time Series Analysis

1.1 Theory: What is Time Series Data?

Definition: A time series is a sequence of data points indexed in time order. Unlike cross-sectional data, the temporal ordering matters because values are often dependent on previous observations.

Mathematical Representation:

$$Y_t = f(Y_{t-1}, Y_{t-2}, \dots, Y_{t-k}, \epsilon_t) \quad (2.1)$$

Where:

- Y_t = Value at time t
- ϵ_t = Random error term
- k = Number of lags

Key Properties:

- **Temporal Dependence:** Today's stock price depends on yesterday's price
- **Autocorrelation:** Correlation between observations at different time points
- **Trend:** Long-term movement in data
- **Seasonality:** Regular patterns that repeat over fixed periods

Change with |papernote

Types of Time Series:

- **Stationary:** Constant mean, variance, and covariance over time
- **Non-stationary:** Mean/variance changes over time
- **Seasonal:** Regular patterns (daily, monthly, yearly)
- **Irregular:** Random fluctuations

1.1.1 Python Implementation (Gold Price Data)

```
1 # Using yfinance (Gold Futures GC=F)
2 import yfinance as yf
3
4 gold = yf.download('GC=F', start='2015-06-01',
5 end='2025-06-01', interval='1d')
6 gold = gold.reset_index()[['Date', 'Close']]
7 gold.to_csv('gold_10y.csv', index=False)
```

Notes:

- Uses COMEX Gold Futures as proxy for spot prices
- Automatically handles weekend gaps in futures data
- For physical gold prices, use GLD ETF instead

1.1.2 Python Implementation (NSE Stock Data)

```
1 # Download MRF stock from NSE India
2 mrf = yf.download('MRF.NS', start='2015-06-01',
3 end='2025-06-01', interval='1d')
4 mrf = mrf.reset_index()[['Date', 'Close']]
5 mrf.to_csv('mrf_10y.csv', index=False)
```

Notes:

- Add .NS suffix for any NSE-listed stock
- Contains adjusted closing prices automatically
- 15-minute delay for live Indian market data

1.1.3 Python Implementation (Cryptocurrency Data)

```

1 # Bitcoin historical data (daily)
2 btc = yf.download('BTC-USD', start='2015-06-01',
3 end='2025-06-01', interval='1d')
4 btc = btc.reset_index()[['Date', 'Close']]
5 btc.to_csv('btc_10y.csv', index=False)

```

Notes:

- 24/7 trading data - no date gaps
- For Ethereum: use ETH-USD ticker
- Includes all trading pairs (USD, INR, etc.)

1.1.4 Python Implementation (Economic Data)

```

1 # COVID-19 data from Our World in Data
2 import pandas as pd
3
4 covid_url = "https://covid.ourworldindata.org/data/owid-covid-data.csv"
5 pd.read_csv(covid_url).to_csv('covid_global.csv', index=False)
6
7 # Temperature data (New Delhi example)
8 !pip install meteostat
9 from meteostat import Point, Daily
10
11 data = Daily(Point(28.6139, 77.2090),
12 start='2015-06-01', end='2025-06-01').fetch()
13 data.to_csv('delhi_temp.csv')

```

Notes:

- COVID data updated daily with 100+ metrics
- Temperature data requires latitude/longitude
- Meteorological data has 1-day latency

1.2 Practical: Data Preparation & Visualization

1.2.1 Excel Implementation (NSE MRF Stock Data)

Step 1: Data Import

- 1 1. Download MRF.csv from NSE India.
- 2 2. Data > From Text/CSV > Transform > Load.
- 3 3. Ensure the Date column is in Date format: Format > Cells > Date.

Step 2: Prepare Data for Trend Analysis

- 1 1. Identify columns:
 - 2 - Column A: Date
 - 3 - Column G: close (closing price)
- 4 2. Create a helper column for trend analysis:
 - 5 - In column N (Row Number), enter 1 for the first data row, 2 for the second, and so on (or use =ROW()-1 and fill down).

Step 3: Linear Trend Analysis

- 1 1. Calculate the slope of the trend:
2 =SLOPE(G2:G38, N2:N38)
3
- 4 2. Calculate the intercept:
5 =INTERCEPT(G2:G38, N2:N38)
6
- 7 3. Generate fitted (trend) values for each row:
8 =TREND(G2:G38, N2:N38, N2:N38)
9 (Enter in O2 and fill down to O38)
10
- 11 4. Obtain full regression statistics (slope, intercept, R-squared, etc.):
12 - Select a 5x2 block of empty cells (e.g., P2:Q6)
13 - Enter:
14 =LINEST(G2:G38, N2:N38, TRUE, TRUE)
15 - Press Ctrl+Shift+Enter (for array formula)

Step 4: Summary Statistics

- 1 =AVERAGE(G2:G38) // Mean closing price
- 2 =STDEV.S(G2:G38) // Standard deviation of closing price
- 3 =MAX(G2:G38)-MIN(G2:G38) // Range of closing price

Interpreting the Results:

- **Slope:** The average change in closing price per time period (Row Number). A positive value indicates an upward trend; a negative value indicates a downward trend.
- **Intercept:** The estimated closing price when Row Number is zero (theoretical starting value).
- **Fitted Values:** The predicted closing price for each period, based on the linear trend.

- **LINEST Output:**

- Top left: Slope
 - Top right: Intercept
 - 2nd row: Standard errors
 - 3rd row: R-squared (goodness of fit), standard error of estimate
 - 4th row: F-statistic, degrees of freedom
 - 5th row: Regression sum of squares, residual sum of squares
- **Mean, Standard Deviation, Range:** Basic descriptive statistics for the closing price.

Tip: To visualize the trend, insert a scatter plot of Date vs. close, then add a linear trendline (right-click on data points > Add Trendline > Linear) and display the equation on the chart.

1.2.2 Python Implementation (YouTube Views Data)

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import numpy as np
4  from datetime import datetime
5
6  # Load and prepare data
7  df = pd.read_csv('youtube_views.csv', parse_dates=['Date'],
8  index_col='Date')
9
10 # Basic statistics
11 print("Time Series Summary:")
12 print(f"Mean: {df['Views'].mean():.2f}")
13 print(f"Std Dev: {df['Views'].std():.2f}")
14 print(f"Min: {df['Views'].min()}")
15 print(f"Max: {df['Views'].max()}")
16
17 # Visualization
18 fig, axes = plt.subplots(2, 1, figsize=(12, 8))
19
20 # Time series plot
21 axes[0].plot(df.index, df['Views'], label='Daily Views')
22 axes[0].set_title('YouTube Views Over Time')
23 axes[0].set_ylabel('Views')
24 axes[0].legend()
25

```

```
26 # Distribution plot
27 axes[1].hist(df['Views'], bins=30, alpha=0.7)
28 axes[1].set_title('Distribution of Views')
29 axes[1].set_xlabel('Views')
30 axes[1].set_ylabel('Frequency')
31
32 plt.tight_layout()
33 plt.show()
```

Interpretation Guidelines:

- Look for trends (upward/downward movement)
- Identify seasonal patterns (regular cycles)
- Spot outliers (unusual spikes/drops)
- Check for volatility clustering

2 Moving Averages & Smoothing Techniques

2.1 Theory: Moving Averages

Definition: A moving average smooths time series data by creating a series of averages from different subsets of the full dataset.

Mathematical Formulation:

Simple Moving Average (SMA):

$$SMA_t = \frac{1}{k} \sum_{i=0}^{k-1} Y_{t-i} \quad (2.2)$$

Exponential Moving Average (EMA):

$$EMA_t = \alpha Y_t + (1 - \alpha) EMA_{t-1} \quad (2.3)$$

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

Weighted Moving Average (WMA):

$$WMA_t = \frac{\sum_{i=0}^{k-1} w_i Y_{t-i}}{\sum_{i=0}^{k-1} w_i} \quad (2.4)$$

Where w_i are weights assigned to each observation.

Trade-offs:

- **Lag:** All moving averages lag behind the actual data
- **Smoothness vs. Responsiveness:** Longer windows = smoother but less responsive
- **End-point problem:** Cannot compute MA for the most recent $k - 1$ periods

2.2 Practical: Smoothing Techniques

2.2.1 Excel Implementation (Stock Data)

```

1 Simple Moving Average:
2 1. Data > Data Analysis > Moving Average
3 2. Input Range: B2:B100
4 3. Interval: 5 (for 5-day MA)
5 4. Output Range: C2
6 5. Chart Output: Yes

```

2.2.2 Python Implementation (Gold Prices)

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Load gold price data
5 df = pd.read_csv('gold_prices.csv', parse_dates=['Date'],
6 index_col='Date')
7
8 # Simple Moving Averages
9 df['SMA_5'] = df['Price'].rolling(window=5).mean()
10 df['SMA_20'] = df['Price'].rolling(window=20).mean()
11 df['SMA_50'] = df['Price'].rolling(window=50).mean()
12
13 # Exponential Moving Average
14 df['EMA_12'] = df['Price'].ewm(span=12).mean()
15
16 # Adaptive Smoothing using Holt-Winters
17 from statsmodels.tsa.holtwinters import ExponentialSmoothing
18
19 # For trend and seasonality
20 model = ExponentialSmoothing(
21 df['Price'],
22 trend='add',          # Additive trend
23 seasonal='add',       # Additive seasonality
24 seasonal_periods=12    # Monthly seasonality
25 )

```

```

26 fit = model.fit()
27 df['Holt_Winters'] = fit.fittedvalues
28
29 # Visualization
30 plt.figure(figsize=(14, 8))
31 plt.plot(df.index, df['Price'], label='Original', alpha=0.7)
32 plt.plot(df.index, df['SMA_5'], label='SMA(5)')
33 plt.plot(df.index, df['SMA_20'], label='SMA(20)')
34 plt.plot(df.index, df['EMA_12'], label='EMA(12)')
35 plt.plot(df.index, df['Holt_Winters'], label='Holt-Winters')
36 plt.title('Gold Prices: Various Smoothing Techniques')
37 plt.legend()
38 plt.show()
39
40 # Calculate smoothing errors
41 mse_sma5 = ((df['Price'] - df['SMA_5'])**2).mean()
42 mse_ema12 = ((df['Price'] - df['EMA_12'])**2).mean()
43 print(f"MSE SMA(5): {mse_sma5:.2f}")
44 print(f"MSE EMA(12): {mse_ema12:.2f}")

```

Interpretation:

- Shorter MA periods = More responsive to changes, more noise
- Longer MA periods = Smoother, less responsive
- EMA gives more weight to recent observations
- Choose smoothing method based on forecasting horizon

Scenario	Recommended MA	Why?	Example
Day Trading	5-10 period SMA	Need quick signals	Gold intraday
Swing Trading	20-50 period SMA	Balance noise/lag	MRF weekly
Long-term Investing	200 period SMA	Major trends only	Index funds
High Volatility	Longer periods or EMA	Reduce false signals	Crypto
Trend Following	Multiple MAs (50/200)	Confirmation	Stock portfolios
Seasonal Data	Period = Season length	Match cycle	Monthly sales

3 Stationarity & Data Transformations

3.1 Theory: Stationarity

Definition: A time series is stationary if its statistical properties (mean, variance, covariance) remain constant over time.

Mathematical Definition: A time series $\{Y_t\}$ is **strictly stationary** if:

$$P(Y_{t_1}, Y_{t_2}, \dots, Y_{t_k}) = P(Y_{t_1+h}, Y_{t_2+h}, \dots, Y_{t_k+h}) \quad (2.5)$$

For any lag h and any time points t_1, t_2, \dots, t_k .

Weak Stationarity (more practical):

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

Why Stationarity Matters:

- Most time series models assume stationarity
- Non-stationary series can lead to spurious regression
- Forecasting requires stable statistical relationships

Types of Non-Stationarity:

1. **Trend Stationarity:** Deterministic trend, remove by detrending
2. **Difference Stationarity:** Stochastic trend, remove by differencing
3. **Structural Breaks:** Parameters change at specific points

Stationarity Tests:

1. **Augmented Dickey-Fuller (ADF):** Tests for unit root
 - H_0 : Series has unit root (non-stationary)
 - H_1 : Series is stationary
2. **KPSS Test:** Tests for stationarity around trend
 - H_0 : Series is stationary
 - H_1 : Series has unit root

3.2 Why Does Time Series Data Need to Be Stationary?

1. Most Models Assume Stationarity

Statistical models like ARIMA, SARIMA, and many forecasting techniques assume that the data's statistical properties (mean, variance, autocorrelation) do not change over time. If this assumption is violated, model estimates become unreliable, forecasts become inaccurate, and hypothesis tests can be misleading.

2. Stationarity Ensures Reliable Forecasting

Forecasting depends on the idea that patterns in the past will persist into the future. If the mean or variance keeps changing, the model cannot “learn” a stable relationship, making future predictions untrustworthy.

3. Prevents Spurious Regression

Spurious regression is a statistical illusion where two unrelated non-stationary series appear to be strongly related simply because they both wander over time.

Example: Gold prices and global temperature may both trend upward over decades, but this does not mean one causes the other.

Real-World Examples:

- **Gold Prices:** Raw gold price series often show a long-term upward trend due to inflation and macroeconomic factors. Making the series stationary (by differencing or detrending) removes this artificial relationship and reveals the true underlying dynamics.
- **Stock Prices (e.g., MRF):** Stock prices are typically non-stationary. Stock returns (percentage change in price), however, are often stationary. Models like ARIMA or GARCH require stationary input, so we analyze returns, not prices.
- **Cryptocurrency (e.g., Bitcoin):** Crypto prices are highly volatile and often exhibit structural breaks. Stationarity is required to model volatility or forecast future values.
- **Meteorological Data (e.g., Temperature):** Daily temperature shows strong seasonality. Removing seasonality (e.g., subtracting monthly average) makes the series stationary, enabling accurate modeling and anomaly detection.

Data Type	Raw Data Stationary?	How to Achieve Stationarity
Gold Price	No (trending)	Differencing, detrending
Stock Price	No (trending)	Use returns (log or percent diff)
Crypto Price	No (volatile, breaks)	Returns, structural break tests
Temperature	No (seasonal)	Remove seasonality, detrending

Summary Table:

Key Takeaways:

- Stationarity is the foundation for most time series modeling and forecasting.
- Non-stationary data can mislead models, cause spurious results, and produce unreliable forecasts.
- Always check for stationarity using visual inspection and statistical tests (ADF, KPSS).
- Transform your data (differencing, detrending, removing seasonality) before modeling.

Stationarity ensures that the patterns we find in our data are real and persistent, not just artifacts of time. Without it, our models are like compasses spinning without direction.

Practice: Try plotting your gold price, stock price, crypto price, and temperature data. Observe the trends and seasonality. Then, difference or detrend the series and see how the statistical properties stabilize—this is the first step to building robust, reliable time series models!

3.3 Python Scripts: Transforming Data to Stationarity

3.3.1 Gold Price (Differencing)

```

1 import pandas as pd
2
3 gold = pd.read_csv('gold_10y.csv', parse_dates=['Date'])
4 gold['Gold_Diff'] = gold['Close'].diff()
5 gold_stationary = gold[['Date', 'Gold_Diff']].dropna()
6 gold_stationary.to_csv('gold_stationary.csv', index=False)

```

3.3.2 Stock Price (Log Returns)

```
1 import pandas as pd
2 import numpy as np
3
4 stock = pd.read_csv('mrf_10y.csv', parse_dates=['Date'])
5 stock['Stock_Returns'] = np.log(stock['Close'] / stock['Close'].shift(1))
6 stock_stationary = stock[['Date', 'Stock_Returns']].dropna()
7 stock_stationary.to_csv('stock_stationary.csv', index=False)
```

3.3.3 Crypto Price (Log Returns)

```
1 import pandas as pd
2 import numpy as np
3
4 btc = pd.read_csv('btc_10y.csv', parse_dates=['Date'])
5 btc['BTC_Returns'] = np.log(btc['Close'] / btc['Close'].shift(1))
6 btc_stationary = btc[['Date', 'BTC_Returns']].dropna()
7 btc_stationary.to_csv('btc_stationary.csv', index=False)
```

3.3.4 Temperature (Remove Seasonality)

```
1 import pandas as pd
2
3 temp = pd.read_csv('delhi_temp_10y.csv', parse_dates=['time'])
4 temp = temp[['time', 'tavg']].dropna()
5 temp['Month'] = temp['time'].dt.month
6 monthly_means = temp.groupby('Month')['tavg'].transform('mean')
7 temp['Deseasonalized'] = temp['tavg'] - monthly_means
8 temp_stationary = temp[['time', 'Deseasonalized']].dropna()
9 temp_stationary.to_csv('temp_stationary.csv', index=False)
```

Colab Download:

```
1 from google.colab import files
2 files.download('gold_stationary.csv')
3 files.download('stock_stationary.csv')
4 files.download('btc_stationary.csv')
5 files.download('temp_stationary.csv')
```

3.4 Practical: Testing and Achieving Stationarity

3.4.1 Excel Implementation (NSE Data Stationarity)

Visual Tests:

```

1 // Plot time series
2 1. Insert > Charts > Line Chart (Date vs Price)
3 2. Add trendline: Right-click series > Add Trendline
4 3. Check if trend is significant
5
6 // Calculate first differences
7 // In column C (assuming prices in B)
8 =B3-B2 // First difference

```

Statistical Tests (Using NumXL Add-in):

```

1 1. Install NumXL add-in
2 2. NumXL > Tests > Unit Root Tests > ADF Test
3 3. Input: B2:B100 (Price series)
4 4. Include: Constant + Trend
5 5. Interpret: p-value < 0.05 indicates stationarity

```

3.4.2 Python Implementation (Comprehensive Stationarity Analysis)

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from statsmodels.tsa.stattools import adfuller, kpss
5 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
6 from scipy import stats
7
8 def test_stationarity(timeseries, title='Time Series'):
9     """
10     Comprehensive stationarity testing and visualization
11     """
12     # Plot the time series
13     fig, axes = plt.subplots(3, 2, figsize=(15, 12))
14
15     # Original series
16     axes[0,0].plot(timeseries)
17     axes[0,0].set_title(f'{title} - Original')
18     axes[0,0].set_ylabel('Value')
19
20     # Rolling statistics
21     rolling_mean = timeseries.rolling(window=12).mean()
22     rolling_std = timeseries.rolling(window=12).std()
23
24     axes[0,1].plot(timeseries, label='Original')
25     axes[0,1].plot(rolling_mean, color='red', label='Rolling Mean')
26     axes[0,1].plot(rolling_std, color='black', label='Rolling Std')
27     axes[0,1].set_title('Rolling Statistics')
28     axes[0,1].legend()

```

```
29
30 # Distribution
31 axes[1,0].hist(timeseries.dropna(), bins=30, alpha=0.7)
32 axes[1,0].set_title('Distribution')
33 axes[1,0].axvline(timeseries.mean(), color='red',
34 linestyle='--', label='Mean')
35
36 # Q-Q plot
37 stats.probplot(timeseries.dropna(), dist="norm",
38 plot=axes[1,1])
39 axes[1,1].set_title('Q-Q Plot')
40
41 # ACF and PACF
42 plot_acf(timeseries.dropna(), ax=axes[2,0], lags=20)
43 plot_pacf(timeseries.dropna(), ax=axes[2,1], lags=20)
44
45 plt.tight_layout()
46 plt.show()
47
48 # Statistical tests
49 print(f"Stationarity Tests for {title}")
50 print("="*50)
51
52 # ADF Test
53 adf_result = adfuller(timeseries.dropna(), autolag='AIC')
54 print(f"ADF Test:")
55 print(f"  Test Statistic: {adf_result[0]:.4f}")
56 print(f"  p-value: {adf_result[1]:.4f}")
57 print(f"  Critical Values:")
58 for key, value in adf_result[4].items():
59     print(f"    {key}: {value:.4f}")
60
61 # KPSS Test
62 kpss_result = kpss(timeseries.dropna(), regression='ct')
63 print(f"\nKPSS Test:")
64 print(f"  Test Statistic: {kpss_result[0]:.4f}")
65 print(f"  p-value: {kpss_result[1]:.4f}")
66 print(f"  Critical Values:")
67 for key, value in kpss_result[3].items():
68     print(f"    {key}: {value:.4f}")
69
70 # Interpretation
71 print("\nInterpretation:")
72 if adf_result[1] < 0.05 and kpss_result[1] > 0.05:
73     print("Series is stationary")
74 elif adf_result[1] > 0.05 and kpss_result[1] < 0.05:
75     print("Series is non-stationary (unit root)")
76 elif adf_result[1] > 0.05 and kpss_result[1] > 0.05:
```



```

77 print("Series is trend stationary")
78 else:
79 print("Results are inconclusive")
80
81 return adf_result, kpss_result
82
83 def make_stationary(timeseries, method='difference'):
84     """
85     Transform series to achieve stationarity
86     """
87     if method == 'difference':
88         # First differencing
89         diff_series = timeseries.diff().dropna()
90         return diff_series, 1
91
92     elif method == 'log_difference':
93         # Log transformation followed by differencing
94         log_series = np.log(timeseries)
95         log_diff = log_series.diff().dropna()
96         return log_diff, 1
97
98     elif method == 'detrend':
99         # Linear detrending
100         from scipy import signal
101         detrended = signal.detrend(timeseries)
102         return pd.Series(detrended, index=timeseries.index), 0
103
104     elif method == 'seasonal_difference':
105         # Seasonal differencing
106         seasonal_diff = timeseries.diff(12).dropna() # Monthly
107         return seasonal_diff, 12
108
109     # Load data
110     df = pd.read_csv('website_traffic.csv', parse_dates=['Date'],
111                     index_col='Date')
112
113     # Test original series
114     print("Testing Original Series")
115     adf_orig, kpss_orig = test_stationarity(df['Traffic'],
116     'Website Traffic')
117
118     # If non-stationary, apply transformations
119     if adf_orig[1] > 0.05:
120         print("\nApplying Transformations...")
121
122     # Try different methods
123     methods = ['difference', 'log_difference', 'detrend']
124     results = {}

```

```
125
126 for method in methods:
127     print(f"\nTrying {method}...")
128     transformed, order = make_stationary(df['Traffic'],
129     method)
130     adf_trans, kpss_trans = test_stationarity(
131     transformed, f'Traffic - {method}')
132     results[method] = {
133         'series': transformed,
134         'adf_pvalue': adf_trans[1],
135         'kpss_pvalue': kpss_trans[1],
136         'order': order
137     }
138
139 # Select best transformation
140 best_method = min(results, key=lambda x: results[x]['adf_pvalue'])
141 print(f"\nBest transformation: {best_method}")
142 print(f"ADF p-value: {results[best_method]['adf_pvalue']:.4f}")
143
144 # Box-Cox transformation
145 def box_cox_transform(data):
146     """
147     Apply Box-Cox transformation
148     """
149     # Find optimal lambda
150     fitted_data, fitted_lambda = stats.boxcox(data[data > 0])
151     print(f"Optimal lambda: {fitted_lambda:.4f}")
152     return fitted_data, fitted_lambda
153
154 # Apply Box-Cox
155 transformed_data, lambda_val = box_cox_transform(df['Traffic'])
156
157 # Inverse transformation for forecasting
158 def inverse_box_cox(data, lambda_val):
159     if lambda_val == 0:
160         return np.exp(data)
161     else:
162         return np.power(lambda_val * data + 1, 1/lambda_val)
```

Decision Framework:

- ADF: Stationary, KPSS: Stationary → **Series is stationary**
- ADF: Non-stationary, KPSS: Stationary → **Trend stationary (detrend)**
- ADF: Non-stationary, KPSS: Non-stationary → **Difference stationary (difference)**
- ADF: Stationary, KPSS: Non-stationary → **Need further investigation**

4 ARMA & ARIMA Models

4.1 Theory: ARMA Models

Autoregressive (AR) Model: An AR(p) model expresses current value as a linear combination of past values:

$$Y_t = c + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \epsilon_t \quad (2.6)$$

Where:

- c = constant
- ϕ_i = autoregressive parameters
- ϵ_t = white noise error term

Real-World Example (Gold Prices):

- *AR(2) Model:* Current gold price depends on past 2 months' prices
- $\phi_1 = 0.6$: 60% weight to last month's price
- $\phi_2 = 0.3$: 30% weight to price from two months ago
- Shock decay: \$100 price spike \rightarrow \$80 \rightarrow \$64 \rightarrow ... (stationary)

Analogy:

- Your current mood (Y_t) depends on:
- 70% of yesterday's mood ($\phi_1 = 0.7$)
- 20% of the mood from two days ago ($\phi_2 = 0.2$)
- Plus today's random events (ϵ_t)

Moving Average (MA) Model: An MA(q) model expresses current value as a linear combination of past error terms:

$$Y_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q} \quad (2.7)$$

Where:

- μ = mean of the series

- θ_i = moving average parameters

Real-World Example (Bitcoin Prices):

- *MA(2) Model:* Current BTC price reflects recent market shocks
- $\theta_1 = 0.6$: 60% of yesterday's news impact
- $\theta_2 = 0.3$: 30% of news from two days ago
- \$5,000 hack effect: $-\$3,000 \rightarrow -\$1,500 \rightarrow \$0$ (finite memory)

Analogy:

- Today's stock price (Y_t) reacts to:
- Today's earnings report (ϵ_t)
- 50% residual impact from yesterday's FDA approval ($\theta_1 = 0.5$)
- 20% lingering effect from last week's CEO scandal ($\theta_2 = 0.2$)

ARMA Model: ARMA(p,q) combines both AR and MA components:

$$Y_t = c + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q} \quad (2.8)$$

Stationarity Conditions:

- For AR(p): All roots of $\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p = 0$ must lie outside unit circle
- For MA(q): All roots of $\theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q = 0$ must lie outside unit circle

ARIMA Models: ARIMA(p,d,q) extends ARMA to handle non-stationary data:

- **p:** Order of autoregression
- **d:** Degree of differencing
- **q:** Order of moving average

$$\phi(B)(1-B)^d Y_t = \theta(B)\epsilon_t \quad (2.9)$$

Where B is the backshift operator: $BY_t = Y_{t-1}$.

Model Selection:

- **ACF/PACF Patterns:** Identify model order
- **Information Criteria:** AIC, BIC for model comparison
- **Residual Analysis:** Check for white noise residuals

Autoregressive Integrated Moving Average (ARIMA) Model: An ARIMA(p,d,q) model combines differencing with AR and MA components to handle non-stationary data:

$$(1 - \phi_1 B - \dots - \phi_p B^p)(1 - B)^d Y_t = (1 + \theta_1 B + \dots + \theta_q B^q) \epsilon_t \quad (2.10)$$

Where:

- p = AR order (past values)
- d = Differencing order (make series stationary)
- q = MA order (past errors)
- B = Backshift operator: $BY_t = Y_{t-1}$

Real-World Example (Gold Prices):

- *ARIMA(1,1,1) Model:*
- Raw gold prices need 1 differencing ($d = 1$) to remove trend: $Y'_t = Y_t - Y_{t-1}$
- $\phi_1 = 0.4$: 40% of yesterday's *price change* affects today
- $\theta_1 = 0.3$: 30% of yesterday's shock lingers
- A \$100 price jump becomes $\$40 \rightarrow \$16 \rightarrow \dots$ (controlled decay)

Analogy:

- Your *weight loss journey* (Y_t):
- First difference ($d = 1$): Focus on *weekly changes* rather than absolute weight
- AR(1): 50% of last week's progress ($\phi_1 = 0.5$) influences this week
- MA(1): 30% impact from last week's cheat meal ($\theta_1 = 0.3$)
- Combines persistent habits (AR) with temporary setbacks (MA)

Key Differences from AR/MA:

- Handles *trends* through differencing (*d*)
- Models *changes* rather than raw values
- Example: Bitcoin prices after removing bubble effects

Component	Role	Gold Price Example
AR(p)	Past price changes	60% of last month's change
I(d)	Remove trend	1 differencing
MA(q)	Past market shocks	30% of last month's news impact

4.2 Python Implementation: AR, MA, and ARIMA Models**4.2.1 Autoregressive (AR) Model**

```

1  import pandas as pd
2  from statsmodels.tsa.ar_model import AutoReg
3  import matplotlib.pyplot as plt
4
5  # Load gold price data
6  gold = pd.read_csv('gold_10y.csv', parse_dates=['Date'])
7  gold = gold.set_index('Date').asfreq('B') # business days frequency
8  gold['Close'] = gold['Close'].interpolate() # fill missing values
9
10 # Fit AR(1) model (current value regressed on previous value)
11 model_ar = AutoReg(gold['Close'], lags=1).fit()
12 gold['AR1_pred'] = model_ar.fittedvalues
13
14 # Plot results
15 gold[['Close', 'AR1_pred']].plot(figsize=(12,4),
16 title='AR(1) Model: Gold Price')
17 plt.show()
18
19 # Print model summary
20 print(model_ar.summary())

```

Output Interpretation:

- **Plot Analysis:** The AR(1) prediction line closely follows the actual gold price but with a slight lag

- **Model Behavior:** AR models capture persistence (today's price depends heavily on yesterday's price)
- **Limitation:** Cannot capture sudden jumps or structural breaks well
- **Best Use:** When data shows strong autocorrelation and gradual changes

4.2.2 Moving Average (MA) Model

```

1  from statsmodels.tsa.arima.model import ARIMA
2
3  # Fit MA(1) model (current value depends on current and previous error)
4  model_ma = ARIMA(gold['Close'], order=(0,0,1)).fit()
5  gold['MA1_pred'] = model_ma.fittedvalues
6
7  # Plot results
8  gold[['Close', 'MA1_pred']].plot(figsize=(12,4),
9  title='MA(1) Model: Gold Price')
10 plt.show()
11
12 # Print model summary
13 print(model_ma.summary())

```

Output Interpretation:

- **Plot Analysis:** MA prediction appears smoother and may lag behind sudden price movements
- **Model Behavior:** Captures short-term shocks and their immediate effects
- **Strength:** Good for data where recent surprises influence current values
- **Limitation:** May not capture long-term trends as effectively as AR models

4.2.3 ARIMA Model

```

1  # First, difference the series to remove trend (d=1)
2  gold['Close_diff'] = gold['Close'].diff()
3
4  # Fit ARIMA(1,1,1): AR(1) + 1 difference + MA(1)
5  model_arima = ARIMA(gold['Close'], order=(1,1,1)).fit()
6  gold['ARIMA_pred'] = model_arima.fittedvalues
7
8  # Plot differenced series and ARIMA predictions
9  gold[['Close_diff', 'ARIMA_pred']].dropna().plot(figsize=(12,4),

```

```
10 title='ARIMA(1,1,1) Model: Gold Price Changes')
11 plt.show()
12
13 # Print model summary
14 print(model_arma.summary())
```

Output Interpretation:

- **Plot Analysis:** Shows the differenced series (price changes) and model predictions
- **Key Insight:** ARIMA models changes, not levels—focuses on whether price goes up/down
- **Advantage:** Handles non-stationary data by differencing first
- **Practical Use:** Best for trending data like financial prices

4.3 Model Comparison Guidelines

What to Look For in Plots:

- **Closeness of Fit:** How well does the predicted line track the actual data?
- **Lag Behavior:** Does the model prediction lag behind sudden changes?
- **Residual Patterns:** Are there systematic patterns in prediction errors?

Model Selection Criteria:

Model	Best For	Key Indicator
AR(p)	Persistent, stationary data	High autocorrelation
MA(q)	Shock-driven, stationary data	Short memory effects
ARIMA(p,d,q)	Trending, non-stationary data	Unit root present

Diagnostic Checks:

```
1 # Check residuals for white noise
2 residuals = model_arma.resid
3 print("Ljung-Box Test p-value:",
4       acorr_ljungbox(residuals, lags=10, return_df=True)['lb_pvalue'].iloc[0])
5
6 # Plot residuals
7 residuals.plot(title='Model Residuals')
8 plt.show()
```

Note: A good model should have residuals that look like white noise (random, no patterns). If patterns remain, consider different model orders or transformations.

4.4 Practical: Building ARMA/ARIMA Models

4.4.1 Excel Implementation (Using NumXL Add-in)

ARMA Model Building:

- 1 1. Ensure data is stationary (use ADF test from Module 3)
- 2 2. NumXL > Model > ARMA
- 3 3. Input Range: Stationary series
- 4 4. Model Order: Start with ARMA(1,1)
- 5 5. Estimation Method: Maximum Likelihood
- 6 6. Generate Forecasts: Yes

Model Diagnostics:

- 1 1. NumXL > Model > Residual Analysis
- 2 2. Check ACF of residuals (should be white noise)
- 3 3. Ljung-Box test: p-value > 0.05 for good model

4.4.2 Python Implementation (Complete ARIMA Workflow)

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from statsmodels.tsa.arima.model import ARIMA
5  from statsmodels.tsa.stattools import adfuller
6  from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
7  from statsmodels.stats.diagnostic import acorr_ljungbox
8  import warnings
9  warnings.filterwarnings('ignore')
10
11 # Load YouTube views data
12 df = pd.read_csv('youtube_views.csv', parse_dates=['Date'],
13 index_col='Date')
14
15 def identify_arima_order(data, max_p=5, max_q=5):
16     """
17     Identify optimal ARIMA order using ACF/PACF and
18     information criteria
19     """
20     # Plot ACF and PACF
21     fig, axes = plt.subplots(1, 2, figsize=(15, 5))
22     plot_acf(data.dropna(), ax=axes[0], lags=20)
23     plot_pacf(data.dropna(), ax=axes[1], lags=20)
24     plt.show()
25
26     # Grid search for optimal parameters
27     best_aic = np.inf

```

```
28 best_order = None
29 aic_results = []
30
31 for p in range(max_p + 1):
32     for q in range(max_q + 1):
33         try:
34             model = ARIMA(data, order=(p, 0, q))
35             results = model.fit()
36             aic = results.aic
37             aic_results.append((p, q, aic))
38
39             if aic < best_aic:
40                 best_aic = aic
41                 best_order = (p, 0, q)
42         except:
43             continue
44
45 # Display AIC table
46 aic_df = pd.DataFrame(aic_results,
47                       columns=['p', 'q', 'AIC'])
48 print("AIC Values for Different Orders:")
49 print(aic_df.pivot(index='p', columns='q',
50                   values='AIC').round(2))
51
52 print(f"\nBest order: ARMA{best_order[0], best_order[2]}")
53 print(f"Best AIC: {best_aic:.2f}")
54
55 return best_order
56
57 def fit_arma_model(data, order):
58     """
59     Fit ARIMA model and perform diagnostics
60     """
61     model = ARIMA(data, order=order)
62     results = model.fit()
63
64     print(results.summary())
65
66     # Diagnostic plots
67     fig, axes = plt.subplots(2, 2, figsize=(15, 10))
68
69     # Residuals plot
70     residuals = results.resid
71     axes[0, 0].plot(residuals)
72     axes[0, 0].set_title('Residuals')
73     axes[0, 0].axhline(y=0, color='red', linestyle='--')
74
75     # Residuals distribution
```

```

76 axes[0, 1].hist(residuals, bins=30, alpha=0.7)
77 axes[0, 1].set_title('Residuals Distribution')
78
79 # ACF of residuals
80 plot_acf(residuals, ax=axes[1, 0], lags=20)
81 axes[1, 0].set_title('ACF of Residuals')
82
83 # Q-Q plot
84 from scipy import stats
85 stats.probplot(residuals, plot=axes[1, 1])
86 axes[1, 1].set_title('Q-Q Plot')
87
88 plt.tight_layout()
89 plt.show()
90
91 # Ljung-Box test
92 lb_test = acorr_ljungbox(residuals, lags=10,
93 return_df=True)
94 print("\nLjung-Box Test:")
95 print(lb_test)
96
97 if lb_test['lb_pvalue'].min() > 0.05:
98     print("Residuals appear to be white noise (Good!)")
99 else:
100     print("Residuals show autocorrelation (Model may need improvement)")
101
102 return results
103
104 def forecast_arima(model_results, steps=30):
105     """
106     Generate forecasts with confidence intervals
107     """
108     forecast = model_results.forecast(steps=steps)
109     conf_int = model_results.get_forecast(steps=steps).conf_int()
110
111     # Plot forecasts
112     plt.figure(figsize=(12, 6))
113
114     # Historical data (last 100 points)
115     plt.plot(model_results.fittedvalues.index[-100:],
116 model_results.fittedvalues[-100:],
117 label='Fitted Values', color='red')
118
119     # Forecasts
120     forecast_index = pd.date_range(
121 start=model_results.fittedvalues.index[-1] + pd.Timedelta(days=1),
122 periods=steps, freq='D')
123     plt.plot(forecast_index, forecast, label='Forecast',

```

```
124 color='blue')
125 plt.fill_between(forecast_index,
126 conf_int.iloc[:, 0],
127 conf_int.iloc[:, 1],
128 alpha=0.3, color='blue')
129
130 plt.title('ARIMA Forecast')
131 plt.legend()
132 plt.show()
133
134 return forecast, conf_int
135
136 # Complete workflow
137 print("Step 1: Check Stationarity")
138 adf_result = adfuller(df['Views'])
139 print(f"ADF p-value: {adf_result[1]:.4f}")
140
141 if adf_result[1] > 0.05:
142     print("Series is non-stationary, applying first difference")
143     df['Views_diff'] = df['Views'].diff().dropna()
144     data_for_modeling = df['Views_diff'].dropna()
145     d = 1
146 else:
147     print("Series is stationary")
148     data_for_modeling = df['Views']
149     d = 0
150
151 print("\nStep 2: Identify Model Order")
152 if d == 0:
153     best_order = identify_arma_order(data_for_modeling)
154 else:
155     # For differenced data, identify ARMA order then add back integration
156     arma_order = identify_arma_order(data_for_modeling)
157     best_order = (arma_order[0], d, arma_order[2])
158
159 print(f"\nStep 3: Fit ARIMA{best_order} Model")
160 final_model = fit_arma_model(df['Views'], best_order)
161
162 print("\nStep 4: Generate Forecasts")
163 forecasts, conf_intervals = forecast_arma(final_model, steps=30)
164
165 # Model evaluation metrics
166 def evaluate_model(actual, fitted):
167     """Calculate evaluation metrics"""
168     mse = np.mean((actual - fitted)**2)
169     rmse = np.sqrt(mse)
170     mae = np.mean(np.abs(actual - fitted))
171     mape = np.mean(np.abs((actual - fitted) / actual)) * 100
```

```

172
173 print("\nModel Evaluation Metrics:")
174 print(f"MSE: {mse:.2f}")
175 print(f"RMSE: {rmse:.2f}")
176 print(f"MAE: {mae:.2f}")
177 print(f"MAPE: {mape:.2f}%")
178
179 # Evaluate on in-sample data
180 evaluate_model(df['Views'], final_model.fittedvalues)

```

Interpretation Guidelines:

- **Parameters:** Significant coefficients (p-value < 0.05)
- **Residuals:** Should be white noise (no autocorrelation)
- **Model Selection:** Lower AIC/BIC indicates better model
- **Forecast Intervals:** Wider intervals indicate more uncertainty

5 Seasonal ARIMA (SARIMA) Models

5.1 Theory: Seasonal ARIMA

Definition: SARIMA models extend ARIMA to handle seasonal patterns in time series data.

Mathematical Formulation: SARIMA(p, d, q) \times (P, D, Q)_s model:

$$\phi(B)\Phi(B^s)(1-B)^d(1-B^s)^DY_t = \theta(B)\Theta(B^s)\epsilon_t \quad (2.11)$$

Where:

- (p, d, q) = Non-seasonal ARIMA order
- (P, D, Q) = Seasonal ARIMA order
- s = Seasonal period (e.g., 12 for monthly data)
- $\Phi(B^s)$ = Seasonal AR polynomial
- $\Theta(B^s)$ = Seasonal MA polynomial

Identification Process:

1. Check for seasonal patterns in ACF/PACF

2. Apply seasonal differencing if needed
3. Identify seasonal and non-seasonal orders
4. Estimate parameters jointly

5.2 Practical: Building SARIMA Models

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from statsmodels.tsa.statespace.sarimax import SARIMAX
5 from statsmodels.tsa.seasonal import seasonal_decompose
6 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
7 import warnings
8 warnings.filterwarnings('ignore')
9
10 # Load seasonal data (e.g., monthly sales)
11 df = pd.read_csv('monthly_sales.csv', parse_dates=['Date'],
12 index_col='Date')
13 df = df.asfreq('MS') # Monthly start frequency
14
15 def seasonal_analysis(data, period=12):
16     """
17     Perform seasonal decomposition and analysis
18     """
19     # Seasonal decomposition
20     decomposition = seasonal_decompose(data, model='multiplicative',
21 period=period)
22
23     # Plot components
24     fig, axes = plt.subplots(4, 1, figsize=(12, 10))
25
26     data.plot(ax=axes[0], title='Original Series')
27     decomposition.trend.plot(ax=axes[1], title='Trend')
28     decomposition.seasonal.plot(ax=axes[2], title='Seasonal')
29     decomposition.resid.plot(ax=axes[3], title='Residual')
30
31     plt.tight_layout()
32     plt.show()
33
34     # Check seasonal strength
35     seasonal_strength = 1 - (decomposition.resid.var() /
36 (decomposition.resid + decomposition.seasonal).var())
37     print(f"Seasonal Strength: {seasonal_strength:.3f}")
38
39     return decomposition
```

```

40
41 def identify_sarima_order(data, seasonal_period=12, max_order=2):
42     """
43     Identify SARIMA order using grid search
44     """
45     # Define parameter ranges
46     p = d = q = range(0, max_order + 1)
47     P = D = Q = range(0, 2)
48
49     # Grid search
50     best_aic = np.inf
51     best_order = None
52     best_seasonal_order = None
53
54     for param in itertools.product(p, d, q):
55         for param_seasonal in itertools.product(P, D, Q):
56             try:
57                 model = SARIMAX(data,
58                                 order=param,
59                                 seasonal_order=param_seasonal + (seasonal_period,),
60                                 enforce_stationarity=False,
61                                 enforce_invertibility=False)
62                 results = model.fit(dispatch=False)
63
64                 if results.aic < best_aic:
65                     best_aic = results.aic
66                     best_order = param
67                     best_seasonal_order = param_seasonal + (seasonal_period,)
68             except:
69                 continue
70
71     print(f"Best SARIMA order: {best_order} x {best_seasonal_order}")
72     print(f"Best AIC: {best_aic:.2f}")
73
74     return best_order, best_seasonal_order
75
76 def fit_sarima_model(data, order, seasonal_order):
77     """
78     Fit SARIMA model with diagnostics
79     """
80     model = SARIMAX(data,
81                     order=order,
82                     seasonal_order=seasonal_order,
83                     enforce_stationarity=False,
84                     enforce_invertibility=False)
85
86     results = model.fit()
87

```

```
88 print(results.summary())
89
90 # Diagnostic plots
91 results.plot_diagnostics(figsize=(15, 12))
92 plt.show()
93
94 return results
95
96 # Analysis workflow
97 print("Seasonal Time Series Analysis")
98 print("="*50)
99
100 # Step 1: Seasonal decomposition
101 decomposition = seasonal_analysis(df['Sales'])
102
103 # Step 2: Identify SARIMA order
104 order, seasonal_order = identify_sarima_order(df['Sales'])
105
106 # Step 3: Fit model
107 sarima_model = fit_sarima_model(df['Sales'], order, seasonal_order)
108
109 # Step 4: Forecast
110 forecast_steps = 24 # 2 years ahead
111 forecast = sarima_model.get_forecast(steps=forecast_steps)
112 forecast_mean = forecast.predicted_mean
113 forecast_conf_int = forecast.conf_int()
114
115 # Plot forecast
116 plt.figure(figsize=(14, 8))
117 plt.plot(df.index, df['Sales'], label='Historical')
118 plt.plot(sarima_model.fittedvalues.index, sarima_model.fittedvalues,
119 label='Fitted', alpha=0.7)
120 plt.plot(forecast_mean.index, forecast_mean, label='Forecast',
121 color='red')
122 plt.fill_between(forecast_conf_int.index,
123 forecast_conf_int.iloc[:, 0],
124 forecast_conf_int.iloc[:, 1],
125 alpha=0.3, color='red')
126 plt.title('SARIMA Forecast')
127 plt.legend()
128 plt.show()
129
130 # Seasonal evaluation
131 def evaluate_seasonal_model(actual, fitted, period=12):
132     """
133     Evaluate seasonal forecasting performance
134     """
135     # Overall metrics
```



```

136 rmse = np.sqrt(np.mean((actual - fitted)**2))
137 mae = np.mean(np.abs(actual - fitted))
138
139 # Seasonal metrics
140 seasonal_errors = []
141 for i in range(period):
142     seasonal_actual = actual[i::period]
143     seasonal_fitted = fitted[i::period]
144     if len(seasonal_actual) > 0 and len(seasonal_fitted) > 0:
145         seasonal_rmse = np.sqrt(np.mean((seasonal_actual - seasonal_fitted)**2))
146         seasonal_errors.append(seasonal_rmse)
147
148 print("Model Evaluation:")
149 print(f"Overall RMSE: {rmse:.2f}")
150 print(f"Overall MAE: {mae:.2f}")
151 print(f"Average Seasonal RMSE: {np.mean(seasonal_errors):.2f}")
152
153 evaluate_seasonal_model(df['Sales'], sarima_model.fittedvalues)

```

Excel Implementation:

```

1 // Seasonal modeling in Excel:
2 1. Create seasonal dummy variables
3 2. Use Data > Data Analysis > Regression with seasonal dummies
4 3. Create seasonal dummy variables for each month/quarter
5 4. NumXL > SARIMA Model (if available)

```

6 Transfer Function Models

6.1 Theory: Transfer Function Models

Definition: Transfer function models relate an output time series to one or more input time series, capturing the dynamic relationship between variables.

Mathematical Formulation:

$$Y_t = \frac{\omega(B)}{\delta(B)} B^b X_t + \frac{\theta(B)}{\phi(B)} \epsilon_t \quad (2.12)$$

Where:

- Y_t = Output series
- X_t = Input series
- $\omega(B), \delta(B)$ = Transfer function polynomials

- $\theta(B), \phi(B)$ = Noise model polynomials
- b = Pure delay parameter
- ϵ_t = White noise

Components:

1. **Transfer Function:** $\frac{\omega(B)}{\delta(B)}B^b$ describes how input affects output
2. **Noise Model:** $\frac{\theta(B)}{\phi(B)}\epsilon_t$ captures unexplained variation

Applications:

- Marketing: How advertising spend affects sales
- Economics: How interest rates affect inflation
- Digital: How content creation affects YouTube views

Model Building Steps:

1. **Prewhitening:** Transform input to white noise
2. **Cross-correlation:** Identify delay and transfer function structure
3. **Model Estimation:** Fit transfer function parameters
4. **Residual Analysis:** Check noise model adequacy

6.2 Practical: Building Transfer Function Models

6.2.1 Python Implementation (YouTube Views vs. Ad Spend)

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from statsmodels.tsa.statespace.sarimax import SARIMAX
5 from statsmodels.tsa.arima.model import ARIMA
6 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
7 from scipy.stats import pearsonr
8 import warnings
9 warnings.filterwarnings('ignore')
10
11 # Load data with input (Ad_Spend) and output (Views)
12 df = pd.read_csv('youtube_ad_data.csv', parse_dates=['Date'],
```

```

13 index_col='Date')
14
15 def prewhiten_series(input_series, max_order=3):
16     """
17     Prewhiten input series by fitting ARIMA model
18     """
19     print("Prewhitening Input Series")
20     print("="*30)
21
22     # Find best ARIMA model for input
23     best_aic = np.inf
24     best_order = None
25
26     for p in range(max_order + 1):
27         for d in range(2):
28             for q in range(max_order + 1):
29                 try:
30                     model = ARIMA(input_series, order=(p, d, q))
31                     results = model.fit()
32                     if results.aic < best_aic:
33                         best_aic = results.aic
34                         best_order = (p, d, q)
35                 except:
36                     continue
37
38     print(f"Best order for input: ARIMA{best_order}")
39
40     # Fit model and get residuals (prewhitened series)
41     model = ARIMA(input_series, order=best_order)
42     results = model.fit()
43     prewhitened = results.resid
44
45     return prewhitened, results
46
47 def cross_correlation_analysis(input_clean, output_clean, max_lag=20):
48     """
49     Calculate and plot cross-correlation function
50     """
51     # Calculate cross-correlations
52     cross_corrs = []
53     lags = range(-max_lag, max_lag + 1)
54
55     for lag in lags:
56         if lag > 0:
57             corr, _ = pearsonr(input_clean[:-lag], output_clean[lag:])
58         else:
59             corr, _ = pearsonr(input_clean[-lag:], output_clean[:lag])
60     cross_corrs.append(corr)

```

```
61
62 # Plot cross-correlations
63 plt.figure(figsize=(12, 6))
64 plt.stem(lags, cross_corrs)
65 plt.axhline(y=0, color='black', linestyle='--', alpha=0.3)
66 plt.axhline(y=1.96/np.sqrt(len(input_clean)), color='red',
67             linestyle='--', alpha=0.5)
68 plt.axhline(y=-1.96/np.sqrt(len(input_clean)), color='red',
69             linestyle='--', alpha=0.5)
70 plt.title('Cross-Correlation Function')
71 plt.xlabel('Lag')
72 plt.ylabel('Cross-Correlation')
73 plt.grid(True, alpha=0.3)
74 plt.show()
75
76 # Identify significant lags
77 threshold = 1.96/np.sqrt(len(input_clean))
78 significant_lags = [lag for lag, corr in zip(lags, cross_corrs)
79                  if abs(corr) > threshold]
80
81 print(f"Significant lags: {significant_lags}")
82
83 return cross_corrs, lags
84
85 def fit_transfer_function(output_series, input_series, delay=0,
86                           transfer_order=(1,1), noise_order=(1,1)):
87     """
88     Fit transfer function model using SARIMAX
89     """
90     print(f"Fitting Transfer Function Model")
91     print(f"Delay: {delay}, Transfer Order: {transfer_order}, Noise Order: {
92           noise_order}")
93     print("="*50)
94
95     # Prepare input with delay
96     if delay > 0:
97         input_delayed = input_series.shift(delay)
98     else:
99         input_delayed = input_series
100
101     # Align data
102     aligned_data = pd.concat([output_series, input_delayed], axis=1).dropna()
103     y = aligned_data.iloc[:, 0]
104     x = aligned_data.iloc[:, 1]
105
106     # Fit SARIMAX model (which can handle exogenous variables)
107     model = SARIMAX(y,
108                     exog=x,
```

```

108 order=noise_order + (0,), # (p,d,q) for noise model
109 enforce_stationarity=False,
110 enforce_invertibility=False)
111
112 results = model.fit()
113
114 print(results.summary())
115
116 # Diagnostics
117 results.plot_diagnostics(figsize=(15, 12))
118 plt.show()
119
120 return results
121
122 def transfer_function_forecast(model_results, future_input, steps=30):
123     """
124     Generate forecasts using transfer function model
125     """
126     print("Transfer Function Forecasting")
127     print("="*30)
128
129     # Generate forecasts
130     forecast = model_results.get_forecast(steps=steps, exog=future_input)
131     forecast_mean = forecast.predicted_mean
132     forecast_ci = forecast.conf_int()
133
134     # Plot results
135     fig, axes = plt.subplots(2, 1, figsize=(15, 12))
136
137     # Historical and forecast for output
138     axes[0].plot(model_results.data.endog[-60:], label='Historical Output',
139                 color='blue')
140     axes[0].plot(model_results.fittedvalues[-60:], label='Fitted',
141                 color='red', alpha=0.7)
142
143     forecast_index = pd.date_range(
144         start=model_results.data.endog.index[-1] + pd.Timedelta(days=1),
145         periods=steps, freq='D')
146     axes[0].plot(forecast_index, forecast_mean, label='Forecast',
147                 color='green', linewidth=2)
148     axes[0].fill_between(forecast_index,
149                         forecast_ci.iloc[:, 0],
150                         forecast_ci.iloc[:, 1],
151                         alpha=0.3, color='green')
152     axes[0].set_title('Output Forecast')
153     axes[0].legend()
154
155     # Input series

```

```
156 axes[1].plot(model_results.data.exog[-60:], label='Historical Input',
157 color='orange')
158 axes[1].plot(forecast_index, future_input, label='Future Input',
159 color='purple', linewidth=2)
160 axes[1].set_title('Input Series')
161 axes[1].legend()
162
163 plt.tight_layout()
164 plt.show()
165
166 return forecast_mean, forecast_ci
167
168 # Complete Transfer Function Workflow
169 print("Transfer Function Modeling: YouTube Views vs Ad Spend")
170 print("="*60)
171
172 print("\nStep 1: Data Exploration")
173 # Plot input and output series
174 fig, axes = plt.subplots(2, 1, figsize=(15, 8))
175 axes[0].plot(df.index, df['Ad_Spend'], label='Ad Spend')
176 axes[0].set_title('Input Series: Ad Spend')
177 axes[0].legend()
178
179 axes[1].plot(df.index, df['Views'], label='Views')
180 axes[1].set_title('Output Series: Views')
181 axes[1].legend()
182 plt.tight_layout()
183 plt.show()
184
185 # Check correlation
186 correlation, p_value = pearsonr(df['Ad_Spend'], df['Views'])
187 print(f"Overall correlation: {correlation:.4f} (p-value: {p_value:.4f})")
188
189 print("\nStep 2: Prewhitening Input Series")
190 prewhitened_input, input_model = prewhiten_series(df['Ad_Spend'])
191
192 print("\nStep 3: Cross-Correlation Analysis")
193 cross_corrs, lags = cross_correlation_analysis(prewhitened_input, df['Views'])
194
195 print("\nStep 4: Fit Transfer Function Model")
196 # Based on cross-correlation, determine delay and orders
197 # For demonstration, using delay=1 and simple orders
198 tf_model = fit_transfer_function(df['Views'], df['Ad_Spend'],
199 delay=1,
200 transfer_order=(1,1),
201 noise_order=(1,0))
202
203 print("\nStep 5: Forecasting")
```

```

204 # Create future input scenario (e.g., constant ad spend)
205 future_ad_spend = pd.Series([df['Ad_Spend'].mean()] * 30)
206 forecast_mean, forecast_ci = transfer_function_forecast(tf_model,
207 future_ad_spend)
208
209 # Model evaluation
210 def evaluate_transfer_function(actual, fitted):
211     """Evaluate transfer function model"""
212     mse = np.mean((actual - fitted)**2)
213     rmse = np.sqrt(mse)
214     mae = np.mean(np.abs(actual - fitted))
215     mape = np.mean(np.abs((actual - fitted) / actual)) * 100
216
217     print("\nTransfer Function Model Evaluation:")
218     print(f"MSE: {mse:.2f}")
219     print(f"RMSE: {rmse:.2f}")
220     print(f"MAE: {mae:.2f}")
221     print(f"MAPE: {mape:.2f}%")
222
223     evaluate_transfer_function(df['Views'], tf_model.fittedvalues)
224
225 # Scenario analysis
226 def scenario_analysis(model, base_input, scenarios):
227     """
228     Analyze different input scenarios
229     """
230     print("\nScenario Analysis")
231     print("="*20)
232
233     results = {}
234     for scenario_name, multiplier in scenarios.items():
235         scenario_input = base_input * multiplier
236         forecast = model.get_forecast(steps=len(scenario_input),
237 exog=scenario_input)
238         results[scenario_name] = {
239             'forecast': forecast.predicted_mean,
240             'total_impact': forecast.predicted_mean.sum()
241         }
242         print(f"{scenario_name}: Total predicted views = {results[scenario_name]['
total_impact']:.0f}")
243
244     return results
245
246 # Run scenario analysis
247 base_spend = pd.Series([df['Ad_Spend'].mean()] * 30)
248 scenarios = {
249     'Low Spend (50%)': 0.5,
250     'Current Spend (100%)': 1.0,

```

```
251     'High Spend (150%)': 1.5,
252     'Maximum Spend (200%)': 2.0
253 }
254
255 scenario_results = scenario_analysis(tf_model, base_spend, scenarios)
```

Excel Implementation:

```
1 // Cross-correlation in Excel:
2 1. Use CORREL function with shifted data
3 2. =CORREL(A2:A100, B3:B101) for lag 1
4 3. Plot correlations vs lags
5 4. Use Regression with lagged inputs for transfer function
```

7 Comprehensive Project: Indian Stock Market Analysis

7.1 Final Project: NSE MRF Stock Prediction

```
1 # Complete end-to-end time series analysis
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6
7 # Load NSE MRF data
8 df = pd.read_csv('MRF_stock_data.csv', parse_dates=['Date'],
9 index_col='Date')
10
11 print("="*60)
12 print("COMPREHENSIVE TIME SERIES ANALYSIS: MRF STOCK")
13 print("="*60)
14
15 # 1. Exploratory Data Analysis
16 print("\n1. EXPLORATORY DATA ANALYSIS")
17 print("-" * 40)
18 print(f>Data Range: {df.index.min()} to {df.index.max()}")
19 print(f>Total Observations: {len(df)}")
20 print(f>Missing Values: {df.isnull().sum().sum()}")
21
22 # 2. Decomposition
23 print("\n2. TIME SERIES DECOMPOSITION")
24 from statsmodels.tsa.seasonal import seasonal_decompose
25 decomposition = seasonal_decompose(df['Close'],
26 model='multiplicative',
27 period=252) # Annual
28 decomposition.plot()
```



```

29 plt.show()
30
31 # 3. Stationarity Tests
32 print("\n3. STATIONARITY ANALYSIS")
33 from statsmodels.tsa.stattools import adfuller
34 adf_result = adfuller(df['Close'])
35 print(f"ADF Test p-value: {adf_result[1]:.6f}")
36
37 # 4. ARIMA Modeling
38 print("\n4. ARIMA MODELING")
39 # Auto ARIMA approach
40 from statsmodels.tsa.arima.model import ARIMA
41 # ... (use previous code)
42
43 # 5. SARIMA for Seasonality
44 print("\n5. SEASONAL MODELING")
45 # ... (use previous SARIMA code)
46
47 # 6. Transfer Function with Market Index
48 print("\n6. TRANSFER FUNCTION WITH NIFTY INDEX")
49 # Load NIFTY data and model relationship
50 # ... (use previous transfer function code)
51
52 # 7. Final Forecast
53 print("\n7. FINAL FORECAST AND TRADING STRATEGY")
54 # Combine models for robust forecast
55 # ... (implement ensemble approach)
56
57 print("\n" + "="*60)
58 print("ANALYSIS COMPLETE")
59 print("="*60)

```

8 Resources & Next Steps

8.1 Recommended Books

- *Forecasting: Principles and Practice* by Hyndman & Athanasopoulos
- *Time Series Analysis* by Hamilton
- *Applied Time Series Analysis* by Cryer & Chan

8.2 Python Libraries

- statsmodels: Statistical modeling

- pmdarima: Auto ARIMA
- prophet: Facebook's forecasting tool
- arch: GARCH models for volatility

8.3 Practice Datasets

- NSE historical data
- RBI economic indicators
- Weather data from IMD
- COVID-19 statistics

8.4 Advanced Topics (Next Level)

- GARCH models for volatility
- Vector Autoregression (VAR)
- State Space Models
- Machine Learning for Time Series