CO  Open in Colab

# Part 1: Fundamentals & Classical Statistical Methods

Time Series Analysis in Python

## Setup and Imports

Before we begin, install the required packages:

```
pip install pandas numpy matplotlib statsmodels scipy scikit-learn
```

```
 1 import pandas as pd
 2 import numpy as np
 3 import matplotlib.pyplot as plt
 4 from statsmodels.tsa.seasonal import seasonal_decompose
 5 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
 6 from statsmodels.tsa.stattools import adfuller, kpss
 7 from statsmodels.tsa.arima.model import ARIMA
 8 from statsmodels.tsa.statespace.sarimax import SARIMAX
 9 from statsmodels.stats.diagnostic import acorr_ljungbox
10 from scipy import stats
11 from sklearn.metrics import mean_absolute_error, mean_squared_error
12 import warnings
13 warnings.filterwarnings('ignore')
14
15 plt.style.use('seaborn-v0_8-darkgrid')
```

## 1. What Is a Time Series?

A **time series** is a sequence of observations collected **over time**, usually at regular intervals.

Examples:

- Daily stock prices
- Hourly electricity consumption
- Monthly sales revenue
- Sensor measurements every second

Why Time Series Are Different

- Observations are **dependent**
- The order of data points **cannot be shuffled**
- Past values influence future values

This dependency structure is the core challenge of time series analysis.

## Key Components

Time series typically contain four components:

1. **Trend (T)**: Long-term increase or decrease in the data
2. **Seasonality (S)**: Regular, predictable patterns that repeat over fixed periods (e.g., yearly, monthly)
3. **Cyclicality (C)**: Patterns that repeat but not at fixed intervals (e.g., economic cycles)
4. **Noise/Irregular (I)**: Random variation that cannot be attributed to trend, seasonality, or cyclicality

## ⌄ Mathematical Representation

- **Additive Model**: $Y(t) = T(t) + S(t) + C(t) + I(t)$

  - Use when seasonal variation is roughly constant over time

- **Multiplicative Model**: $Y(t) = T(t) \times S(t) \times C(t) \times I(t)$

  - Use when seasonal variation increases with the level of the series

## Creating a Sample Time Series

Let's create a synthetic time series that contains trend, seasonality, and noise:

```
1  # Set random seed for reproducibility
2  np.random.seed(42)
3
4  # Generate date range
5  date_range = pd.date_range(start='2020-01-01', end='2023-12-31', freq='D'
6  n = len(date_range)
7
8  # Components
9  trend = np.linspace(100, 150, n)  # Linear trend from 100 to 150
10 seasonality = 10 * np.sin(2 * np.pi * np.arange(n) / 365.25)  # Yearly se
```

```
11 noise = np.random.normal(0, 5, n)  # Random noise
12
13 # Combine components (additive model)
14 ts_data = trend + seasonality + noise
15
16 # Create pandas Series with datetime index
17 ts = pd.Series(ts_data, index=date_range, name='Value')
18
19 print("Sample Time Series:")
20 print(ts.head(10))
21 print(f"\nShape: {ts.shape}")
22 print(f"Period: {ts.index.min()} to {ts.index.max()}")
```

```
Sample Time Series:
2020-01-01     102.483571
2020-01-02      99.514941
2020-01-03     103.650916
2020-01-04     108.233733
2020-01-05      99.653774
2020-01-06      99.859609
2020-01-07     109.131857
2020-01-08     105.278161
2020-01-09      99.298455
2020-01-10     104.563060
Freq: D, Name: Value, dtype: float64

Shape: (1461,)
Period: 2020-01-01 00:00:00 to 2023-12-31 00:00:00
```
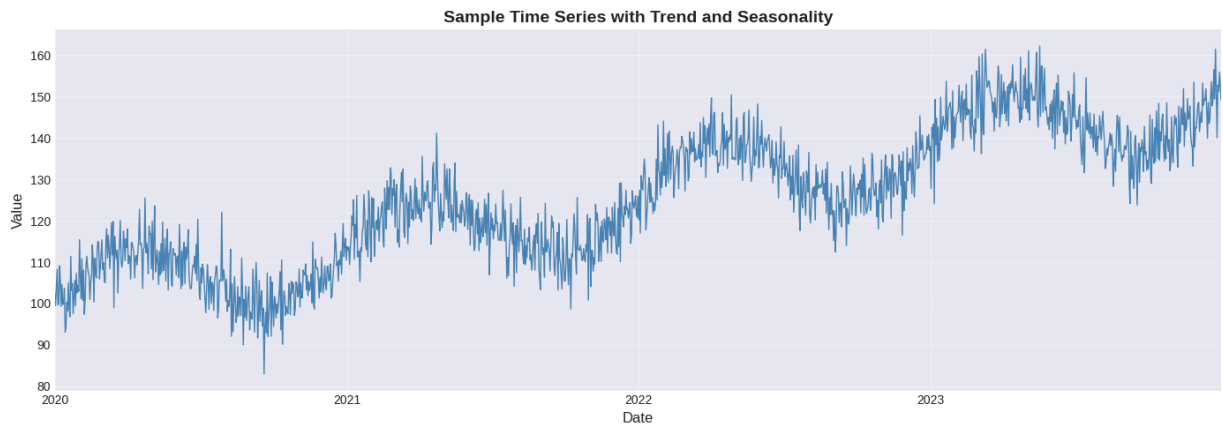
## ⌄ Visualizing the Time Series

```
1 fig, ax = plt.subplots(figsize=(14, 5))
2 ts.plot(ax=ax, linewidth=1, color='steelblue')
3 ax.set_title('Sample Time Series with Trend and Seasonality', fontsize=14,
4 ax.set_xlabel('Date', fontsize=12)
5 ax.set_ylabel('Value', fontsize=12)
6 ax.grid(True, alpha=0.3)
7 plt.tight_layout()
8 plt.show()
```

Sample Time Series with Trend and Seasonality

## 2. Time Series Decomposition

### ∨ Theory

Decomposition is the process of separating a time series into its constituent components. This helps us:

- Understand the underlying patterns
- Remove seasonality for better modeling
- Identify anomalies
- Choose appropriate forecasting methods

## Classical Decomposition Methods

The most common method is **moving average decomposition**:

1. Estimate trend using moving average
2. Detrend the series
3. Estimate seasonal component by averaging detrended values for each season

4. Calculate residuals: Residual = Original - Trend - Seasonal

## Performing Decomposition

```
 1  # Perform additive decomposition
 2  # Period = 365 because we have daily data with yearly seasonality
 3  decomposition = seasonal_decompose(ts, model='additive', period=365)
 4
 5  # Extract components
 6  trend_component = decomposition.trend
 7  seasonal_component = decomposition.seasonal
 8  residual_component = decomposition.resid
 9
10  print("Decomposition Components:")
11  print(f"Trend:    {trend_component.shape}")
12  print(f"Seasonal: {seasonal_component.shape}")
13  print(f"Residual: {residual_component.shape}")
```

```
Decomposition Components:
Trend:    (1461,)
Seasonal: (1461,)
Residual: (1461,)
```

## ⌄ Visualizing Decomposition

```
 1 fig, axes = plt.subplots(4, 1, figsize=(14, 10))
 2
 3 # Original
 4 ts.plot(ax=axes[0], linewidth=1, color='black')
 5 axes[0].set_ylabel('Original', fontsize=11, fontweight='bold')
 6 axes[0].set_title('Time Series Decomposition (Additive Model)', fontsize=14
 7
 8 # Trend
 9 trend_component.plot(ax=axes[1], linewidth=1.5, color='orange')
10 axes[1].set_ylabel('Trend', fontsize=11, fontweight='bold')
11
12 # Seasonal
13 seasonal_component.plot(ax=axes[2], linewidth=1, color='green')
14 axes[2].set_ylabel('Seasonal', fontsize=11, fontweight='bold')
15
16 # Residual
17 residual_component.plot(ax=axes[3], linewidth=0.8, color='red', alpha=0.7)
18 axes[3].set_ylabel('Residual', fontsize=11, fontweight='bold')
19 axes[3].set_xlabel('Date', fontsize=12)
20 axes[3].axhline(y=0, color='black', linestyle='--', linewidth=0.8, alpha=0.
21
22 plt.tight_layout()
23 plt.show()
```

## 3. Stationarity

### What is Stationarity?

A time series is **stationary** if its statistical properties (mean, variance, autocorrelation) do not change over time.

### Why Does Stationarity Matter?

Most classical time series models (ARIMA, SARIMA) assume stationarity because:

- Statistical properties are easier to model when constant
- Predictions are more reliable
- Mathematical theory is simpler

### Types of Stationarity

1. **Strict Stationarity**: Joint distribution is time-invariant (very restrictive)
2. **Weak/Covariance Stationarity**: Only mean, variance, and autocorrelation are constant (commonly used)

### Common Non-Stationary Patterns

- **Trend**: Mean changes over time
- **Seasonality**: Pattern repeats at regular intervals
- **Heteroscedasticity**: Variance changes over time

## Augmented Dickey-Fuller (ADF) Test

The ADF test checks the null hypothesis that the series has a unit root (non-stationary).

- **$H_0$**: Series has a unit root (non-stationary)
- **$H_1$**: Series is stationary

**Decision Rule**: If p-value < 0.05, reject $H_0$ (series is stationary)

```
1 def check_stationarity(timeseries, name='Series'):
2     """
3     Perform Augmented Dickey-Fuller test
4     """
5     print(f"\n{'='*60}")
6     print(f"Stationarity Test: {name}")
```

```
 7      print('='*60)
 8
 9      # Remove NaN values
10      ts_clean = timeseries.dropna()
11
12      # Perform ADF test
13      result = adfuller(ts_clean, autolag='AIC')
14
15      print(f'ADF Statistic:      {result[0]:.6f}')
16      print(f'p-value:            {result[1]:.6f}')
17      print(f'# Lags Used:        {result[2]}')
18      print(f'# Observations:     {result[3]}')
19      print('\nCritical Values:')
20      for key, value in result[4].items():
21          print(f'  {key}: {value:.3f}')
22
23      # Interpretation
24      print('\n' + '-'*60)
25      if result[1] <= 0.05:
26          print(f"✓ STATIONARY (p={result[1]:.4f} < 0.05)")
27          print("  → Reject null hypothesis")
28      else:
29          print(f"✗ NON-STATIONARY (p={result[1]:.4f} > 0.05)")
30          print("  → Fail to reject null hypothesis")
31      print('-'*60)
32
```

## ⌄ Testing Our Series

```
1 # Test original series
2 is_stationary = check_stationarity(ts, 'Original Time Series')
```

```
============================================================
Stationarity Test: Original Time Series
============================================================
ADF Statistic:      -0.835398
p-value:            0.808502
# Lags Used:        12
# Observations:     1448

Critical Values:
  1%: -3.435
  5%: -2.864
  10%: -2.568


------------------------------------------------------------
✗ NON-STATIONARY (p=0.8085 > 0.05)
  → Fail to reject null hypothesis
------------------------------------------------------------
```

## Making a Series Stationary

## Method 1: Differencing

Differencing removes trends and can stabilize the mean:

**First Difference**: Y'(t) = Y(t) – Y(t-1)

```
1 # Apply first differencing
2 ts_diff = ts.diff().dropna()
3
4 # Test differenced series
5 is_stationary_diff = check_stationarity(ts_diff, 'First Differenced Series'
```

```
============================================================
Stationarity Test: First Differenced Series
============================================================
ADF Statistic:     -17.823692
p-value:           0.000000
# Lags Used:       11
# Observations:    1448

Critical Values:
  1%: -3.435
  5%: -2.864
  10%: -2.568


------------------------------------------------------------
✓ STATIONARY (p=0.0000 < 0.05)
  → Reject null hypothesis
------------------------------------------------------------
```

## Visualizing the Transformation

```
1 fig, axes = plt.subplots(2, 1, figsize=(14, 8), sharex=True)
2
3 # Original
4 ts.plot(ax=axes[0], linewidth=1.2, color='steelblue')
5 axes[0].set_title('Original Series (Non-Stationary)', fontsize=12, fontwe:
6 axes[0].set_ylabel('Value', fontsize=11)
7 axes[0].grid(True, alpha=0.3)
8
9 # Differenced
10 ts_diff.plot(ax=axes[1], linewidth=1, color='orange')
11 axes[1].set_title('First Differenced Series (Stationary)', fontsize=12, fo
12 axes[1].set_ylabel('Differenced Value', fontsize=11)
13 axes[1].axhline(y=0, color='red', linestyle='--', linewidth=1, alpha=0.7)
14 axes[1].set_xlabel('Date', fontsize=12)
```

```
15 axes[1].grid(True, alpha=0.3)
16
17 plt.tight_layout()
18 plt.show()
```

**Original Series (Non-Stationary)**

**First Differenced Series (Stationary)**

## 4. Autocorrelation Analysis

## ACF: Autocorrelation Function

The ACF measures the correlation between observations at different time lags:

**ACF(k) = Corr($Y_t$, $Y_{t-k}$)**

- Lag 1: correlation between consecutive observations
- Lag 2: correlation between observations 2 time periods apart
- And so on...

## PACF: Partial Autocorrelation Function

The PACF measures the correlation between $Y_t$ and $Y_{t-k}$ after removing the effect of intermediate lags.

## Why Are ACF and PACF Important?
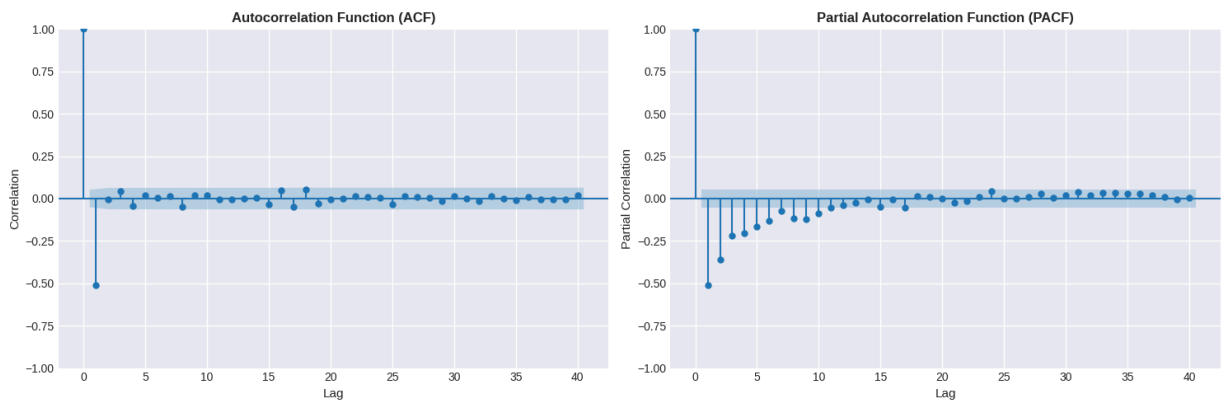
They help identify the order of AR and MA components:

| Pattern | Model Suggestion |
| --- | --- |
| PACF cuts off after lag p, ACF decays | AR(p) |
| ACF cuts off after lag q, PACF decays | MA(q) |
| Both decay gradually | ARMA(p,q) |

## Plotting ACF and PACF

```
1 # Use differenced series (should be stationary)
2 fig, axes = plt.subplots(1, 2, figsize=(15, 5))
3
4 # ACF
5 plot_acf(ts_diff, lags=40, ax=axes[0])
6 axes[0].set_title('Autocorrelation Function (ACF)', fontsize=12, fontweight
7 axes[0].set_xlabel('Lag', fontsize=11)
8 axes[0].set_ylabel('Correlation', fontsize=11)
9
10 # PACF
11 plot_pacf(ts_diff, lags=40, ax=axes[1])
12 axes[1].set_title('Partial Autocorrelation Function (PACF)', fontsize=12, f
13 axes[1].set_xlabel('Lag', fontsize=11)
14 axes[1].set_ylabel('Partial Correlation', fontsize=11)
15
16 plt.tight_layout()
17 plt.show()
```

**Interpretation**:

- Blue shaded area represents the confidence interval
- Bars outside this area are statistically significant
- Look for where bars drop inside the confidence interval (cutoff point)

## 5. Classical Time Series Models

| Model | Full Name | Equation | When to Use |
|-------|-----------|----------|-------------|
| AR(p) | Autoregressive | $Y_t = c + \phi_1 Y_{t-1} + ... + \phi_p Y_{t-p} + \varepsilon_t$ | PACF cuts off |
| MA(q) | Moving Average | $Y_t = \mu + \varepsilon_t + \theta_1 \varepsilon_{t-1} + ... + \theta_\varphi \varepsilon_{t-\varphi}$ | ACF cuts off |
| ARMA(p,q) | AR + MA | Combines both | Both ACF & PACF decay |
| ARIMA(p,d,q) | Integrated ARMA | ARMA on d-times differenced data | Non-stationary series |
| SARIMA(p,d,q)(P,D,Q,s) | Seasonal ARIMA | Adds seasonal components | Seasonal patterns |

## Preparing Data for Modeling

Let's create a simpler dataset to demonstrate the models:

```
1 # Generate AR(1) process
2 np.random.seed(123)
3 n = 300
4 dates = pd.date_range(start='2022-01-01', periods=n, freq='D')
5
6 # AR(1): Y(t) = 0.7 * Y(t-1) + noise
7 ar_coef = 0.7
8 ar_data = [0]
9 for i in range(1, n):
10     ar_data.append(ar_coef * ar_data[i-1] + np.random.normal(0, 1))
11
12 ts_simple = pd.Series(ar_data, index=dates, name='Value')
13
14 # Train-test split (80-20)
15 train_size = int(len(ts_simple) * 0.8)
16 train = ts_simple[:train_size]
17 test = ts_simple[train_size:]
18
19 print(f"Train: {len(train)} observations")
20 print(f"Test:  {len(test)} observations")
```

```
Train: 240 observations
Test:  60 observations
```

## Model 1: AR (Autoregressive)

### Theory

An AR(p) model predicts the current value using p past values:

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

Where:

- c is a constant
- $\phi_1, \phi_2, \dots, \phi_p$ are coefficients
- $\varepsilon_t$ is white noise error

### Implementation

```
1 # Fit AR(1) model
2 # ARIMA(p,d,q) with p=1, d=0, q=0
3 ar_model = ARIMA(train, order=(1, 0, 0))
4 ar_fit = ar_model.fit()
5
6 print("AR(1) Model Summary:")
7 print(ar_fit.summary())
```

```
8 print(f"\nEstimated coefficient: {ar_fit.params[1]:.4f}")
```

```
AR(1) Model Summary:
                           SARIMAX Results
==============================================================================
Dep. Variable:                    Value   No. Observations:                  240
Model:                   ARIMA(1, 0, 0)   Log Likelihood              -345.782
Date:                 Sat, 27 Dec 2025   AIC                          697.563
Time:                         07:05:21   BIC                          708.005
Sample:                       01-01-2022   HQIC                         701.771
                            - 08-28-2022
Covariance Type:                    opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const         -0.0308      0.229     -0.135      0.893     -0.479       0.418
ar.L1          0.7119      0.047     15.074      0.000      0.619       0.805
sigma2         1.0415      0.094     11.072      0.000      0.857       1.226
==============================================================================
Ljung-Box (L1) (Q):                   0.15   Jarque-Bera (JB):                 0
Prob(Q):                              0.70   Prob(JB):                         0
Heteroskedasticity (H):               0.60   Skew:                            -0
Prob(H) (two-sided):                  0.03   Kurtosis:                         3
==============================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-s

Estimated coefficient: 0.7119
True coefficient: 0.7
```

## Interpreting AR(1) Model Output

**What this really is (in one sentence)**

**AR(1) is just linear regression where the predictor is yesterday's value.**

## Connect it to regular regression (intuition first)

In ordinary regression you write:

> $y = \beta_0 + \beta_1 x_1 + \text{error}$

In AR(1), you write:

> $y_t = \text{constant} + \beta_1 \times y_{t-1} + \text{error}$

Same idea. The "feature" is just the lagged value of y itself.

## Translate the output into regression language

**1. The coefficients**

```
const     = -0.0308
ar.L1     =  0.7119
sigma2    =  1.0415
```

Think of this as a regression table:

| Feature | Coefficient | What it means |
|---|---|---|
| Intercept | -0.031 | Baseline (when $y_{t-1}$ = 0) |
| $y_{t-1}$ | 0.7119 | Today = 71% of yesterday + noise |
| Error variance | 1.042 | Residual spread |

### The prediction equation is:

$y_t$ = -0.031 + 0.7119 × $y_{t-1}$ + $\varepsilon_t$

This is identical to:

sales_today = -0.031 + 0.7119 × sales_yesterday + noise

## 2. Standard errors and p-values

Exactly the same as regression:

```
          coef    std err   p-value
const    -0.031   0.229     0.893
ar.L1     0.712   0.047     0.000
```

**const**: p = 0.893 (high)

Not significant. Could drop it.

**ar.L1**: p = 0.000 (very low)

Highly significant. Yesterday strongly predicts today.

Same interpretation as:

feature_1: p = 0.893 → not useful feature_2: p = 0.000 → very useful

The model found **strong autocorrelation** in the data.

## 3. Model quality metrics

```
AIC = 697.6
BIC = 708.0
Log Likelihood = -345.8
```

These are **goodness-of-fit measures**, not accuracy metrics.

**AIC and BIC**: Lower is better

Use these to compare models:

- AR(1): AIC = 697.6
- AR(2): AIC = 695.3 → better
- MA(1): AIC = 701.2 → worse

Think of AIC like:

> penalized training error

It balances fit with model complexity.

**Log Likelihood**: Higher is better

Similar to minimizing loss in ML. Maximum likelihood estimation finds coefficients that make the observed data most probable.

### 4. Diagnostic tests (are residuals well-behaved?)

```
Ljung-Box (Q):           0.15     p = 0.70
Jarque-Bera (JB):        0.10     p = 0.95
Heteroskedasticity (H):  0.60     p = 0.03
```

Think of these as **residual plots in test form**.

**Ljung-Box test**

Question: "Do residuals have patterns left?"

- $p > 0.05$ → No patterns (good)
- $p < 0.05$ → Patterns remain (bad)

Here: $p = 0.70$ → residuals look random

This is like checking:

> plot(residuals) shows no trend

**Jarque-Bera test**

Question: "Are residuals normally distributed?"

- $p > 0.05$ → Yes (good)
- $p < 0.05$ → No (may need transformation)

Here: $p = 0.95$ → very normal

This is like checking:

> histogram(residuals) looks bell-shaped

**Heteroskedasticity test**

Question: "Is variance constant over time?"

- $p > 0.05 \rightarrow$ Yes (good)
- $p < 0.05 \rightarrow$ No (variance changes)

Here: $p = 0.03 \rightarrow$ some heteroskedasticity detected

This is like seeing:

> residual spread increases over time

**5. Model validation**

```
Estimated coefficient: 0.7119
True coefficient:      0.7000
```

The model recovered the true data-generating process.

In ML terms:

> The model learned the correct function

This would be like:

- You generate data: $y = 2x + noise$
- Your model learns: $y = 1.98x$
- Close match $\rightarrow$ model works

## Why it feels alien

### 1. No feature matrix visible

The "feature" ($y_{t-1}$) is created internally from the time series.

### 2. Maximum likelihood instead of MSE

Same goal (fit the data), different math. Likelihood is more general than squared error.

### 3. Heavy focus on diagnostics

Econometrics cares deeply about **why** the model works, not just **that** it works.

Statistical inference > predictive accuracy

### 4. Different vocabulary

- ML says: "$R^2 = 0.85$"

- Econometrics says: "Log Likelihood = -345, AIC = 697"

Same information, different packaging.

## How to read this output (decision rules)

**Check coefficients:**

- $p < 0.05$ → significant → keep
- $p > 0.05$ → not significant → consider dropping

**Check diagnostics:**

- Ljung-Box $p > 0.05$ → good (no autocorrelation left)
- Ljung-Box $p < 0.05$ → bad (model incomplete)

**Compare models:**

- Lower AIC → better model
- Lower BIC → better model (penalizes complexity more)

**Validate:**

- Residuals should look random
- Residuals should be roughly normal
- Variance should be constant (ideally)

## One grounding sentence you can remember

> AR(1) is linear regression where X is yesterday's y, and the output tells you both fit quality and residual behavior in statistical language instead of ML metrics.

## Quick interpretation of this specific output

**Model found:** $y_t = 0.71 \times y_{t-1} + \text{noise}$

**Constant term:** Not significant ($p = 0.89$) → probably zero in reality

**Autocorrelation:** Strong (coefficient = 0.71, $p < 0.001$) → yesterday matters a lot

**Residuals:** Clean (Ljung-Box $p = 0.70$, JB $p = 0.95$) → no patterns left

**Minor issue:** Slight heteroskedasticity ($p = 0.03$) → variance not perfectly constant

**Conclusion:** This is a good model. The AR(1) structure correctly captures the data-generating process.

## ⌄ Model 2: MA (Moving Average)

## Theory

An MA(q) model predicts the current value using past forecast errors:

$Y_t = \mu + \varepsilon_t + \theta_1\varepsilon_{t-1} + \theta_2\varepsilon_{t-2} + ... + \theta_\varphi\varepsilon_{t-\varphi}$

Where:

- $\mu$ is the mean
- $\theta_1, \theta_2, ..., \theta_\varphi$ are coefficients
- $\varepsilon_t, \varepsilon_{t-1}, ...$ are error terms

## Implementation

```
1 # Fit MA(1) model
2 ma_model = ARIMA(train, order=(0, 0, 1))
3 ma_fit = ma_model.fit()
4
5 print("MA(1) Model Summary:")
6 print(ma_fit.summary())
```

```
MA(1) Model Summary:
                               SARIMAX Results
==============================================================================
Dep. Variable:                    Value   No. Observations:                  240
Model:                   ARIMA(0, 0, 1)   Log Likelihood                -370.431
Date:                 Sat, 27 Dec 2025   AIC                            746.863
Time:                         07:05:22   BIC                            757.305
Sample:                       01-01-2022   HQIC                           751.070
                            - 08-28-2022
Covariance Type:                    opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const         -0.0349      0.120     -0.290      0.771      -0.270       0.200
ma.L1          0.6287      0.051     12.246      0.000       0.528       0.729
sigma2         1.2801      0.119     10.766      0.000       1.047       1.513
==============================================================================
Ljung-Box (L1) (Q):                 10.99   Jarque-Bera (JB):                 0
Prob(Q):                             0.00   Prob(JB):                         0
Heteroskedasticity (H):              0.60   Skew:                            -0
Prob(H) (two-sided):                 0.03   Kurtosis:                         2
==============================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-s
```

## Model 3: ARMA

## Theory

ARMA(p,q) combines AR(p) and MA(q):

$$Y_t = c + \phi_1 Y_{t-1} + ... + \phi_p Y_{t-p} + \theta_1 \varepsilon_{t-1} + ... + \theta_\varphi \varepsilon_{t-\varphi} + \varepsilon_t$$

**Important**: ARMA requires stationary data.

## Implementation

```
1 # Fit ARMA(1,1) model
2 arma_model = ARIMA(train, order=(1, 0, 1))
3 arma_fit = arma_model.fit()
4
5 print("ARMA(1,1) Model Summary:")
6 print(arma_fit.summary())
```

```
ARMA(1,1) Model Summary:
                           SARIMAX Results
================================================================================
Dep. Variable:                     Value   No. Observations:              240
Model:                    ARIMA(1, 0, 1)   Log Likelihood             -345.671
Date:                  Sat, 27 Dec 2025   AIC                         699.341
Time:                          07:05:23   BIC                         713.264
Sample:                        01-01-2022   HQIC                        704.951
                             - 08-28-2022
Covariance Type:                     opg
================================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
--------------------------------------------------------------------------------
const         -0.0308      0.223     -0.138      0.890      -0.468       0.407
ar.L1          0.6912      0.072      9.667      0.000       0.551       0.831
ma.L1          0.0424      0.096      0.441      0.659      -0.146       0.231
sigma2         1.0405      0.094     11.083      0.000       0.857       1.225
================================================================================
Ljung-Box (L1) (Q):                  0.00   Jarque-Bera (JB):               0
Prob(Q):                             0.96   Prob(JB):                       0
Heteroskedasticity (H):              0.60   Skew:                          -0
Prob(H) (two-sided):                 0.02   Kurtosis:                       3
================================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-s
```

## Model 4: ARIMA

## Theory

ARIMA(p,d,q) is ARMA applied to d-times differenced data:

- **p**: Order of autoregressive part
- **d**: Degree of differencing

- **q**: Order of moving average part

**Steps**:

1. Difference the series d times to make it stationary
2. Apply ARMA(p,q) to the differenced series

# Implementation

```python
1 # Fit ARIMA(1,1,1)
2 arima_model = ARIMA(train, order=(1, 1, 1))
3 arima_fit = arima_model.fit()
4
5 print("ARIMA(1,1,1) Model Summary:")
6 print(arima_fit.summary())
7 print(f"\nAIC: {arima_fit.aic:.2f}")
8 print(f"BIC: {arima_fit.bic:.2f}")
9
10 # Forecast
11 forecast_steps = len(test)
12 arima_forecast = arima_fit.forecast(steps=forecast_steps)
```

```
ARIMA(1,1,1) Model Summary:
                               SARIMAX Results
==============================================================================
Dep. Variable:                    Value   No. Observations:              240
Model:                   ARIMA(1, 1, 1)   Log Likelihood             -346.335
Date:                 Sat, 27 Dec 2025   AIC                         698.670
Time:                         07:05:25   BIC                         709.100
Sample:                     01-01-2022   HQIC                        702.873
                          - 08-28-2022
Covariance Type:                    opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1          0.7190      0.052     13.919      0.000       0.618       0.820
ma.L1         -1.0000      6.486     -0.154      0.877     -13.712      11.712
sigma2         1.0459      6.772      0.154      0.877     -12.227      14.319
==============================================================================
Ljung-Box (L1) (Q):                  0.08   Jarque-Bera (JB):                @
Prob(Q):                             0.78   Prob(JB):                        @
Heteroskedasticity (H):              0.59   Skew:                           -@
Prob(H) (two-sided):                 0.02   Kurtosis:                        3
==============================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-s

AIC: 698.67
BIC: 709.10
```
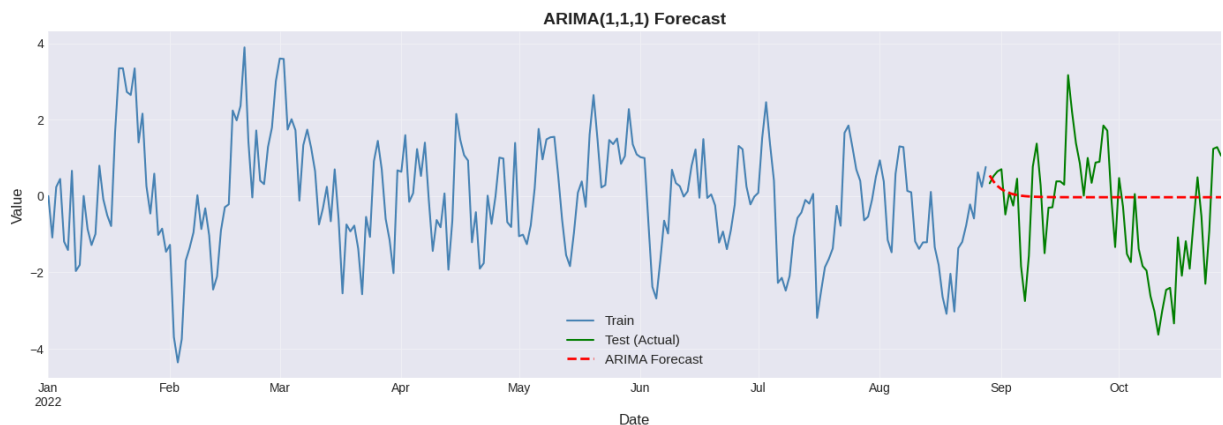
## Visualizing Forecast

```
 1 fig, ax = plt.subplots(figsize=(14, 5))
 2
 3 train.plot(ax=ax, label='Train', linewidth=1.5, color='steelblue')
 4 test.plot(ax=ax, label='Test (Actual)', linewidth=1.5, color='green')
 5 arima_forecast.plot(ax=ax, label='ARIMA Forecast', linewidth=2,
 6                     linestyle='--', color='red')
 7
 8 ax.set_title('ARIMA(1,1,1) Forecast', fontsize=14, fontweight='bold')
 9 ax.set_xlabel('Date', fontsize=12)
10 ax.set_ylabel('Value', fontsize=12)
11 ax.legend(fontsize=11)
12 ax.grid(True, alpha=0.3)
13 plt.tight_layout()
14 plt.show()
```



## Model 5: SARIMA (Seasonal ARIMA)

## Theory

SARIMA(p,d,q)(P,D,Q,s) extends ARIMA to handle seasonality:

**Non-seasonal part**: (p,d,q) **Seasonal part**: (P,D,Q,s)

- P: Seasonal AR order
- D: Seasonal differencing order
- Q: Seasonal MA order
- s: Seasonal period (e.g., 12 for monthly data with yearly seasonality)

## Creating Seasonal Data

```python
1 # Generate monthly data with seasonality
2 n_months = 240  # 20 years
3 dates_monthly = pd.date_range(start='2004-01-01', periods=n_months, freq='M
4
5 # Components
6 trend = np.linspace(50, 100, n_months)
7 seasonal = 15 * np.sin(2 * np.pi * np.arange(n_months) / 12)
8 noise = np.random.normal(0, 3, n_months)
9
10 ts_seasonal = pd.Series(trend + seasonal + noise,
11                         index=dates_monthly,
12                         name='Sales')
13
14 # Split
15 train_seas = ts_seasonal[:int(len(ts_seasonal) * 0.8)]
16 test_seas = ts_seasonal[int(len(ts_seasonal) * 0.8):]
17
18 print(f"Seasonal series shape: {ts_seasonal.shape}")
19 print(f"Train: {len(train_seas)}, Test: {len(test_seas)}")
```

```
Seasonal series shape: (240,)
Train: 192, Test: 48
```

## Fitting SARIMA

```python
1 # Fit SARIMA(1,1,1)(1,1,1,12)
2 # Seasonal period = 12 months
3 sarima_model = SARIMAX(train_seas,
4                        order=(1, 1, 1),
5                        seasonal_order=(1, 1, 1, 12))
6 sarima_fit = sarima_model.fit(disp=False)
7
8 print("SARIMA(1,1,1)(1,1,1,12) Model Summary:")
9 print(sarima_fit.summary())
10
11 # Forecast
12 sarima_forecast = sarima_fit.forecast(steps=len(test_seas))
```

```
SARIMA(1,1,1)(1,1,1,12) Model Summary:
                            SARIMAX Results
==============================================================================
Dep. Variable:                         Sales   No. Observations:
Model:             SARIMAX(1, 1, 1)x(1, 1, 1, 12)   Log Likelihood
Date:                      Sat, 27 Dec 2025   AIC
Time:                              07:05:31   BIC
Sample:                          01-01-2004   HQIC
                               - 12-01-2019
Covariance Type:                        opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1         -0.0127      0.084     -0.151      0.880      -0.178       0.153
ma.L1         -0.9801      0.040    -24.629      0.000      -1.058      -0.902
ar.S.L12      -0.0239      0.104     -0.231      0.817      -0.227       0.179
ma.S.L12      -0.9156      0.140     -6.531      0.000      -1.190      -0.641
sigma2         8.9150      1.336      6.671      0.000       6.296      11.534
==============================================================================
Ljung-Box (L1) (Q):                 0.03   Jarque-Bera (JB):             1
Prob(Q):                            0.87   Prob(JB):                     0
Heteroskedasticity (H):             1.30   Skew:                        -0
Prob(H) (two-sided):                0.31   Kurtosis:                     2
==============================================================================

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-s
```
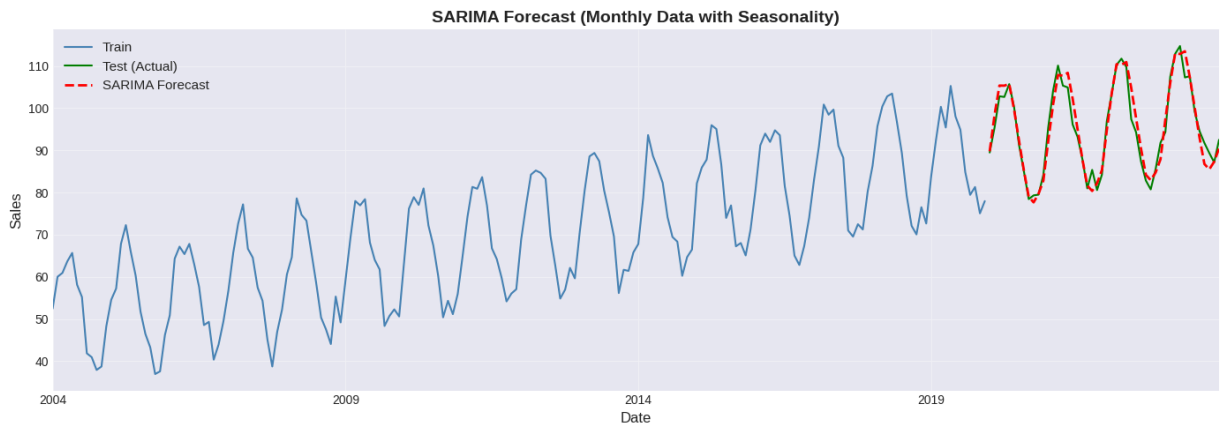
## Visualizing Seasonal Forecast

```python
1 fig, ax = plt.subplots(figsize=(14, 5))
2
3 train_seas.plot(ax=ax, label='Train', linewidth=1.5, color='steelblue')
4 test_seas.plot(ax=ax, label='Test (Actual)', linewidth=1.5, color='green')
5 sarima_forecast.plot(ax=ax, label='SARIMA Forecast', linewidth=2,
6                      linestyle='--', color='red')
7
8 ax.set_title('SARIMA Forecast (Monthly Data with Seasonality)',
9              fontsize=14, fontweight='bold')
10 ax.set_xlabel('Date', fontsize=12)
11 ax.set_ylabel('Sales', fontsize=12)
12 ax.legend(fontsize=11)
13 ax.grid(True, alpha=0.3)
14 plt.tight_layout()
15 plt.show()
```

SARIMA Forecast (Monthly Data with Seasonality)

# 6. Model Selection

## ⌄ Information Criteria

### AIC (Akaike Information Criterion)

**AIC = 2k - 2ln(L)**

Where:

- k = number of parameters
- L = maximum likelihood

### BIC (Bayesian Information Criterion)

**BIC = k·ln(n) - 2ln(L)**

Where:

- n = number of observations

**Rule**: Lower AIC/BIC = better model (balances fit and complexity)

## Grid Search for Best Model

```python
1 def evaluate_arima_models(data, p_range, d_range, q_range):
2     """
3     Evaluate different ARIMA configurations
4     """
5     results = []
6
7     for p in p_range:
8         for d in d_range:
9             for q in q_range:
10                 try:
11                     model = ARIMA(data, order=(p, d, q))
12                     fitted = model.fit()
13
14                     results.append({
15                         'order': (p, d, q),
16                         'AIC': fitted.aic,
17                         'BIC': fitted.bic
18                     })
19                 except:
20                     continue
21
22     return pd.DataFrame(results)
23
24 # Search for best model
25 print("Searching for best ARIMA model...")
26 results_df = evaluate_arima_models(
27     train,
28     p_range=range(0, 3),
29     d_range=range(0, 2),
30     q_range=range(0, 3)
31 )
32
33 # Sort and display
34 results_df = results_df.sort_values('AIC')
35 print("\nTop 5 Models by AIC:")
36 print(results_df.head())
37
38 best_order = results_df.iloc[0]['order']
39 print(f"\nBest Model: ARIMA{best_order}")
```

```
Searching for best ARIMA model...

Top 5 Models by AIC:
       order          AIC          BIC
6   (1, 0, 0)   697.563191   708.005108
13  (2, 0, 1)   698.028252   715.431446
```

```
10  (1, 1, 1)  698.670134  709.099525
12  (2, 0, 0)  699.335230  713.257786
 7  (1, 0, 1)  699.341125  713.263681

Best Model: ARIMA(1, 0, 0)
```

# 7. Model Diagnostics

## Why Check Diagnostics?

After fitting a model, we need to verify that:

1. **Residuals are white noise** (random, no pattern)
2. **Residuals are normally distributed**
3. **No autocorrelation in residuals**
4. **Model assumptions are satisfied**

## Theory: Good Residuals

If the model is appropriate, residuals should have:

- Mean ≈ 0
- Constant variance (homoscedastic)
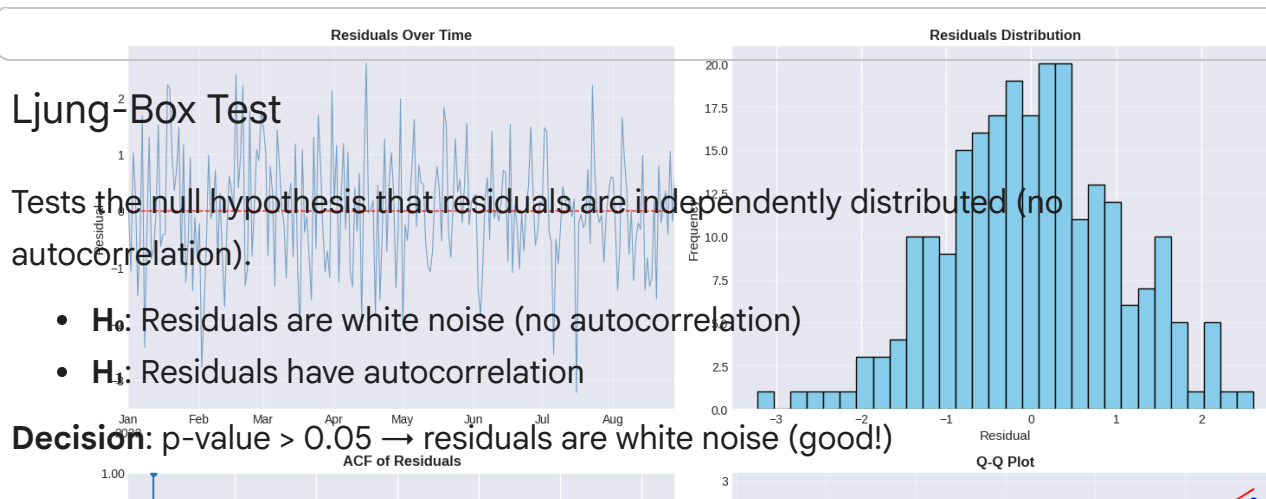- No autocorrelation
- Normal distribution

## Diagnostic Plots

```python
1 # Get residuals from best model
2 best_model = ARIMA(train, order=best_order).fit()
3 residuals = best_model.resid
4
5 # Create diagnostic plots
6 fig, axes = plt.subplots(2, 2, figsize=(14, 10))
7
8 # 1. Residuals over time
9 residuals.plot(ax=axes[0, 0], linewidth=0.8, color='steelblue', alpha=0.7
10 axes[0, 0].axhline(y=0, color='red', linestyle='--', linewidth=1)
11 axes[0, 0].set_title('Residuals Over Time', fontsize=12, fontweight='bold
12 axes[0, 0].set_ylabel('Residual', fontsize=11)
13 axes[0, 0].grid(True, alpha=0.3)
14
15 # 2. Histogram
16 residuals.hist(ax=axes[0, 1], bins=30, edgecolor='black', color='skyblue'
17 axes[0, 1].set_title('Residuals Distribution', fontsize=12, fontweight='b
```

```
18 axes[0, 1].set_xlabel('Residual', fontsize=11)
19 axes[0, 1].set_ylabel('Frequency', fontsize=11)
20 axes[0, 1].grid(True, alpha=0.3, axis='y')
21
22 # 3. ACF of residuals
23 plot_acf(residuals, lags=30, ax=axes[1, 0])
24 axes[1, 0].set_title('ACF of Residuals', fontsize=12, fontweight='bold')
25
26 # 4. Q-Q plot
27 stats.probplot(residuals, dist="norm", plot=axes[1, 1])
28 axes[1, 1].set_title('Q-Q Plot', fontsize=12, fontweight='bold')
29 axes[1, 1].grid(True, alpha=0.3)
30
31 plt.tight_layout()
32 plt.show()
```

# Ljung-Box Test

Tests the null hypothesis that residuals are independently distributed (no autocorrelation).

- **H$_0$**: Residuals are white noise (no autocorrelation)
- **H$_1$**: Residuals have autocorrelation

**Decision**: p-value > 0.05 → residuals are white noise (good!)

```
1 # Perform Ljung-Box test
2 lb_test = acorr_ljungbox(residuals, lags=[10, 20, 30], return_df=True)
3
4 print("\nLjung-Box Test Results:")
5 print(lb_test)
6 print("\nInterpretation:")
7 print("If p-value > 0.05: Residuals are white noise ✓")
8 print("If p-value < 0.05: Residuals have autocorrelation ✗")
```

```
Ljung-Box Test Results:
       lb_stat   lb_pvalue
10    7.147022   0.711496
20   31.838671   0.045048
30   41.570079   0.077860

Interpretation:
If p-value > 0.05: Residuals are white noise ✓
If p-value < 0.05: Residuals have autocorrelation ✗
```

# Residual Statistics

```
1 print("\nResidual Summary Statistics:")
2 print("="*50)
3 print(f"Mean:      {residuals.mean():>10.6f}  (should be ≈ 0)")
4 print(f"Std Dev:   {residuals.std():>10.6f}")
5 print(f"Min:       {residuals.min():>10.6f}")
6 print(f"Max:       {residuals.max():>10.6f}")
7 print(f"Skewness:  {residuals.skew():>10.6f}  (should be ≈ 0)")
8 print(f"Kurtosis:  {residuals.kurtosis():>10.6f}  (should be ≈ 0)")
```

```
Residual Summary Statistics:
==================================================
Mean:        -0.000090  (should be ≈ 0)
Std Dev:      1.022687
Min:         -3.222888
```

```
Max:           2.614799
Skewness:    -0.045184  (should be ≈ 0)
Kurtosis:     0.071627  (should be ≈ 0)
```

## 8. Model Evaluation

⌄  Common Forecasting Metrics

### 1. MAE (Mean Absolute Error)

**MAE = (1/n) $\Sigma|y_i - \hat{y}_i|$**

- Average absolute difference
- Same units as original data
- Easy to interpret

### 2. RMSE (Root Mean Squared Error)

**RMSE = $\sqrt{[(1/n)\,\Sigma(y_i - \hat{y}_i)^2]}$**

- Penalizes large errors more
- Same units as original data
- More sensitive to outliers than MAE

### 3. MAPE (Mean Absolute Percentage Error)

**MAPE = (100/n) $\Sigma|((y_i - \hat{y}_i)/y_i)|$**

- Expressed as percentage
- Scale-independent
- Undefined when $y_i = 0$

## Computing Metrics

```python
1 def calculate_metrics(actual, predicted):
2     """Calculate forecasting metrics"""
3     mae = mean_absolute_error(actual, predicted)
4     rmse = np.sqrt(mean_squared_error(actual, predicted))
5     mape = np.mean(np.abs((actual - predicted) / actual)) * 100
6
7     return {'MAE': mae, 'RMSE': rmse, 'MAPE': mape}
8
9 # Evaluate forecast
10 forecast = best_model.forecast(steps=len(test))
11 metrics = calculate_metrics(test, forecast)
```

```
12
13 print("\nForecast Evaluation Metrics:")
14 print("="*50)
15 print(f"MAE:  {metrics['MAE']:.4f}")
16 print(f"RMSE: {metrics['RMSE']:.4f}")
17 print(f"MAPE: {metrics['MAPE']:.2f}%")
```

```
Forecast Evaluation Metrics:
==================================================
MAE:  1.2387
RMSE: 1.5522
MAPE: 102.35%
```

## Summary and Key Takeaways