

BÀI TẬP CUỐI KHÓA

August 17, 2025

Contents

1	GIỚI THIỆU CHUNG	2
1.1	Giới thiệu về chủ đề nghiên cứu	2
1.2	Giới thiệu về bộ dữ liệu	2
1.3	Mục đích nghiên cứu	2
1.4	Phương pháp nghiên cứu	3
2	KHÁM PHÁ DỮ LIỆU	3
2.1	Chuẩn bị và tải dữ liệu	3
2.2	Làm sạch và tiền xử lý dữ liệu	8
2.3	Phân tích thống kê mô tả và trực quan hóa	11
3	PHÂN TÍCH DỮ LIỆU	15
3.1	Chuẩn bị dữ liệu cho mô hình	15
3.2	Xây dựng mô hình hồi quy đa biến	17
3.3	Xây dựng mô hình chuỗi thời gian ARIMA	20
4	KẾT LUẬN VÀ NHẬN XÉT	25
4.1	Kết quả chính từ phân tích dữ liệu	25
4.2	Kết quả từ các mô hình	26
4.3	Ý nghĩa thực tiễn	26
4.4	Hạn chế và đề xuất cải tiến	26
4.5	Kết luận cuối cùng	26
5	TÀI LIỆU THAM KHẢO	27

1 GIỚI THIỆU CHUNG

1.1 Giới thiệu về chủ đề nghiên cứu

Nghiên cứu này tập trung vào việc phân tích các chỉ số kinh tế vĩ mô của Việt Nam từ năm 1986 đến 2023, bao gồm:

- Tổng sản phẩm quốc nội (GDP) và tốc độ tăng trưởng GDP
- Tỷ lệ lạm phát (CPI)
- Tỷ lệ thất nghiệp
- Đầu tư trực tiếp nước ngoài (FDI)
- Kim ngạch xuất nhập khẩu

Việt Nam trải qua quá trình chuyển đổi từ nền kinh tế kế hoạch hóa tập trung sang nền kinh tế thị trường với sự can thiệp của Nhà nước kể từ năm Đổi Mới 1986. Việc phân tích dữ liệu chuỗi thời gian giúp hiểu rõ hơn về:

- Xu hướng phát triển kinh tế dài hạn
- Các chu kỳ kinh tế và tác động của các sự kiện quan trọng
- Mối quan hệ giữa các chỉ số kinh tế vĩ mô

1.2 Giới thiệu về bộ dữ liệu

Nguồn dữ liệu: World Bank Open Data API (<https://data.worldbank.org/>)

Mô tả dữ liệu:

- Loại dữ liệu: Chuỗi thời gian kinh tế vĩ mô
- Tần suất quan sát: Hàng năm
- Khoảng thời gian: 1986-2023 (38 năm)
- Số lượng quan sát: 38 điểm dữ liệu cho mỗi chỉ số
- Phạm vi địa lý: Việt Nam

Các biến trong dữ liệu:

- `year`: Năm quan sát (1986-2023)
- `GDP_Current_USD`: GDP danh nghĩa (tỷ USD)
- `GDP_Growth`: Tốc độ tăng trưởng GDP (%/năm)
- `Inflation_CPI`: Tỷ lệ lạm phát CPI (%/năm)
- `Unemployment`: Tỷ lệ thất nghiệp (%)
- `FDI_Inflows`: Đầu tư trực tiếp nước ngoài vào (% GDP)
- `Trade_GDP`: Tổng kim ngạch thương mại (% GDP)

1.3 Mục đích nghiên cứu

Câu hỏi nghiên cứu chính:

1. Xu hướng phát triển kinh tế Việt Nam từ Đổi Mới đến nay như thế nào?
2. Các yếu tố nào ảnh hưởng đến tăng trưởng GDP của Việt Nam?
3. Có thể dự báo được các chỉ số kinh tế của Việt Nam trong tương lai không?

Giả thuyết nghiên cứu:

- FDI và thương mại quốc tế có tác động tích cực đến tăng trưởng GDP
- Lạm phát có mối quan hệ nghịch với tăng trưởng kinh tế
- Các chỉ số kinh tế có tính xu hướng và có thể dự báo được

1.4 Phương pháp nghiên cứu

Phương pháp phân tích chuỗi thời gian:

- Phân tích xu hướng và tính mùa vụ
- Kiểm định tính dừng (ADF test, KPSS test)
- Mô hình ARIMA cho dự báo
- Phân tích tự tương quan (ACF/PACF)

Phương pháp hồi quy:

- Hồi quy tuyến tính đa biến
- Ridge/Lasso Regression
- Kiểm tra các giả định của mô hình hồi quy

Ứng dụng thống kê:

- Thống kê mô tả: trung bình, độ lệch chuẩn, xu hướng
- Kiểm định: tính dừng, tự tương quan, normality
- Đánh giá mô hình: RMSE, MAE, MAPE, AIC, BIC
- Dự báo và khoảng tin cậy

2 KHÁM PHÁ DỮ LIỆU

2.1 Chuẩn bị và tải dữ liệu

```
[26]: # Import các thư viện cần thiết
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

# Thư viện cho time series analysis
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import adfuller, kpss
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from statsmodels.tsa.stattools import acf, pacf

# Thư viện cho regression analysis
from sklearn.model_selection import train_test_split
```

```

from sklearn.linear_model import LinearRegression, Ridge, Lasso
from sklearn.metrics import mean_squared_error, mean_absolute_error
from sklearn.preprocessing import StandardScaler
import statsmodels.api as sm

# Thư viện để tải dữ liệu
import requests
import json

# Cài đặt style cho plots
plt.style.use('seaborn-v0_8')
sns.set_palette("husl")

print("Đã import thành công tất cả thư viện cần thiết!")

```

Đã import thành công tất cả thư viện cần thiết!

```

[27]: # Tải dữ liệu từ World Bank API
def get_world_bank_data(indicator, country='VN', start_year=1986,
    ↪end_year=2023):
    """
    Tải dữ liệu từ World Bank API
    """
    base_url = "https://api.worldbank.org/v2/country"
    url = f"{base_url}/{country}/indicator/{indicator}?date={start_year}:
    ↪{end_year}&format=json&per_page=500"

    try:
        response = requests.get(url)
        data = response.json()

        if len(data) > 1 and data[1]:
            df = pd.DataFrame(data[1])
            df = df[['date', 'value']].copy()
            df.columns = ['year', indicator]
            df['year'] = pd.to_numeric(df['year'])
            df[indicator] = pd.to_numeric(df[indicator], errors='coerce')
            return df.sort_values('year').reset_index(drop=True)
        else:
            print(f"Không có dữ liệu cho indicator {indicator}")
            return None
    except Exception as e:
        print(f"Lỗi khi tải dữ liệu {indicator}: {e}")
        return None

# Định nghĩa các chỉ số kinh tế cần tải
indicators = {

```

```

'NY.GDP.MKTP.CD': 'GDP_Current_USD',          # GDP (current US$)
'NY.GDP.MKTP.KD.ZG': 'GDP_Growth',             # GDP growth (annual %)
'FP.CPI.TOTL.ZG': 'Inflation_CPI',             # Inflation, consumer prices
↳ (annual %)
'SL.UEM.TOTL.ZS': 'Unemployment',              # Unemployment, total (% of
↳ total labor force)
'BX.KLT.DINV.WD.GD.ZS': 'FDI_Inflows',         # Foreign direct investment, net
↳ inflows (% of GDP)
'NE.TRD.GNFS.ZS': 'Trade_GDP'                  # Trade (% of GDP)
}

print("Bắt đầu tải dữ liệu từ World Bank API...")

# Tải dữ liệu cho từng chỉ số
dfs = []
for wb_code, var_name in indicators.items():
    print(f"Đang tải {var_name}...")
    df = get_world_bank_data(wb_code)
    if df is not None:
        dfs.append(df)
    else:
        print(f"Không thể tải {var_name}")

print("Hoàn thành tải dữ liệu!")

```

```

Bắt đầu tải dữ liệu từ World Bank API...
Đang tải GDP_Current_USD...
Đang tải GDP_Growth...
Đang tải Inflation_CPI...
Đang tải Unemployment...
Đang tải FDI_Inflows...
Đang tải Trade_GDP...
Hoàn thành tải dữ liệu!

```

```

[28]: # Kết hợp tất cả dữ liệu và lưu thành file CSV
if dfs:
    # Merge tất cả dataframes
    vietnam_data = dfs[0]
    for df in dfs[1:]:
        vietnam_data = pd.merge(vietnam_data, df, on='year', how='outer')

    # Sắp xếp theo năm
    vietnam_data = vietnam_data.sort_values('year').reset_index(drop=True)

    # Chuyển đổi GDP từ USD sang tỷ USD
    if 'GDP_Current_USD' in vietnam_data.columns:
        vietnam_data['GDP_Current_USD'] = vietnam_data['GDP_Current_USD'] / 1e9

```

```

    print("Dữ liệu đã được kết hợp thành công từ World Bank API!")
else:
    print("Không thể tải được dữ liệu từ World Bank API, đọc từ file CSV...")

# Đọc dữ liệu từ file CSV (đã được tải sẵn)
vietnam_data = pd.read_csv('vietnam_economic_data.csv')

# Đổi tên cột cho dễ hiểu
column_mapping = {
    'NY.GDP.MKTP.CD': 'GDP_Current_USD',
    'NY.GDP.MKTP.KD.ZG': 'GDP_Growth',
    'FP.CPI.TOTL.ZG': 'Inflation_CPI',
    'SL.UEM.TOTL.ZS': 'Unemployment',
    'BX.KLT.DINV.WD.GD.ZS': 'FDI_Inflows',
    'NE.TRD.GNFS.ZS': 'Trade_GDP'
}

vietnam_data = vietnam_data.rename(columns=column_mapping)

# Chuyển đổi GDP sang tỷ USD
if 'GDP_Current_USD' in vietnam_data.columns:
    vietnam_data['GDP_Current_USD'] = vietnam_data['GDP_Current_USD'] / 1e9

print(f"Kích thước dữ liệu: {vietnam_data.shape}")
print("\nCác cột trong dữ liệu:")
print(vietnam_data.columns.tolist())

```

Dữ liệu đã được kết hợp thành công từ World Bank API!
 Kích thước dữ liệu: (38, 7)
 \nCác cột trong dữ liệu:
 ['year', 'GDP_Current_USD', 'GDP_Growth', 'Inflation_CPI', 'Unemployment',
 'FDI_Inflows', 'Trade_GDP']

```

[29]: # Hiển thị thông tin cơ bản về dữ liệu
print("=== THÔNG TIN DỮ LIỆU ===")
print(vietnam_data.info())

print("\n=== 5 HÀNG ĐẦU TIÊN ===")
print(vietnam_data.head())

print("\n=== 5 HÀNG CUỐI CÙNG ===")
print(vietnam_data.tail())

print("\n=== THỐNG KÊ MÔ TẢ CƠ BẢN ===")
print(vietnam_data.describe())

```

=== THÔNG TIN DỮ LIỆU ===

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 38 entries, 0 to 37

Data columns (total 7 columns):

#	Column	Non-Null Count	Dtype
0	year	38 non-null	int64
1	GDP_Current_USD	38 non-null	float64
2	GDP_Growth	38 non-null	float64
3	Inflation_CPI	28 non-null	float64
4	Unemployment	33 non-null	float64
5	FDI_Inflows	38 non-null	float64
6	Trade_GDP	38 non-null	float64

dtypes: float64(6), int64(1)

memory usage: 2.2 KB

None

\n=== 5 HÀNG ĐẦU TIÊN ===

	year	GDP_Current_USD	GDP_Growth	Inflation_CPI	Unemployment	\
0	1986	26.336616	2.789292	NaN	NaN	
1	1987	36.658109	3.583470	NaN	NaN	
2	1988	25.423813	5.135012	NaN	NaN	
3	1989	6.293305	7.364513	NaN	NaN	
4	1990	6.471741	5.100918	NaN	NaN	

	FDI_Inflows	Trade_GDP
0	0.000152	23.218694
1	0.028271	20.798607
2	0.030208	18.950487
3	0.064672	57.904462
4	2.781323	81.315698

\n=== 5 HÀNG CUỐI CÙNG ===

	year	GDP_Current_USD	GDP_Growth	Inflation_CPI	Unemployment	\
33	2019	334.365270	7.359263	2.795824	1.681	
34	2020	346.615739	2.865413	3.220934	2.103	
35	2021	366.474753	2.553729	1.834716	2.385	
36	2022	413.445231	8.537500	3.156507	1.523	
37	2023	433.857681	5.065024	3.252893	1.645	

	FDI_Inflows	Trade_GDP
33	4.821075	164.704215
34	4.558362	163.245863
35	4.273146	186.675833
36	4.329473	183.153619
37	4.264071	164.818071

\n=== THỐNG KÊ MÔ TẢ CƠ BẢN ===

	year	GDP_Current_USD	GDP_Growth	Inflation_CPI	Unemployment	\
count	38.000000	38.000000	38.000000	28.000000	33.000000	

mean	2004.500000	126.680228	6.428771	5.682256	1.900697
std	11.113055	132.802249	1.701909	5.165227	0.455290
min	1986.000000	6.293305	2.553729	-1.710337	0.999000
25%	1995.250000	26.463387	5.516786	3.196271	1.681000
50%	2004.500000	51.530555	6.556627	3.957691	1.972000
75%	2013.750000	228.515805	7.439883	7.502250	2.140000
max	2023.000000	433.857681	9.540480	23.115448	2.870000

	FDI_Inflows	Trade_GDP
count	38.000000	38.000000
mean	4.848946	114.876685
std	2.646549	43.738768
min	0.000152	18.950487
25%	3.912313	84.163210
50%	4.317174	123.776036
75%	5.328926	145.285688
max	11.939483	186.675833

2.2 Làm sạch và tiền xử lý dữ liệu

```
[30]: # Kiểm tra dữ liệu thiếu
print("=== KIỂM TRA DỮ LIỆU THIẾU ===")
missing_data = vietnam_data.isnull().sum()
missing_percentage = (missing_data / len(vietnam_data)) * 100

missing_summary = pd.DataFrame({
    'Số lượng thiếu': missing_data,
    'Tỷ lệ (%)': missing_percentage
})

print(missing_summary)

# Kiểm tra outliers sử dụng IQR method
print("\n=== KIỂM TRA OUTLIERS ===")
numeric_columns = vietnam_data.select_dtypes(include=[np.number]).columns
numeric_columns = [col for col in numeric_columns if col != 'year']

outlier_summary = {}
for col in numeric_columns:
    if vietnam_data[col].notna().sum() > 0: # Chỉ xử lý nếu có dữ liệu
        Q1 = vietnam_data[col].quantile(0.25)
        Q3 = vietnam_data[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR

        outliers = vietnam_data[(vietnam_data[col] < lower_bound) |
```



```

                (vietnam_data[col] > upper_bound)][col]
outlier_summary[col] = len(outliers)

if len(outliers) > 0:
    print(f"{col}: {len(outliers)} outliers")
    print(f"  Giá trị outliers: {outliers.tolist()}")
else:
    print(f"{col}: Không có outliers")

print(f"\nTổng kết outliers: {outlier_summary}")

```

=== KIỂM TRA DỮ LIỆU THIẾU ===

	Số lượng thiếu	Tỷ lệ (%)
year	0	0.000000
GDP_Current_USD	0	0.000000
GDP_Growth	0	0.000000
Inflation_CPI	10	26.315789
Unemployment	5	13.157895
FDI_Inflows	0	0.000000
Trade_GDP	0	0.000000

\n=== KIỂM TRA OUTLIERS ===

GDP_Current_USD: Không có outliers

GDP_Growth: 1 outliers

Giá trị outliers: [2.55372852648131]

Inflation_CPI: 2 outliers

Giá trị outliers: [23.1154483474477, 18.6777322770707]

Unemployment: 1 outliers

Giá trị outliers: [2.87]

FDI_Inflows: 10 outliers

Giá trị outliers: [0.0001518797996661, 0.0282712448214878, 0.030207899995773, 0.0646719015911437, 11.9394828586715, 8.58596585460392, 9.71308063712156, 8.2700967582097, 8.65471771434805, 9.6630390545572]

Trade_GDP: Không có outliers

\nTổng kết outliers: {'GDP_Current_USD': 0, 'GDP_Growth': 1, 'Inflation_CPI': 2, 'Unemployment': 1, 'FDI_Inflows': 10, 'Trade_GDP': 0}

[31]: # Xử lý dữ liệu thiếu và outliers

```
print("=== XỬ LÝ DỮ LIỆU THIẾU ===")
```

```
# Tạo bản sao để xử lý
```

```
data_cleaned = vietnam_data.copy()
```

```
# Xử lý missing data bằng interpolation tuyến tính (phù hợp cho time series)
```

```
for col in numeric_columns:
```

```
    if data_cleaned[col].isnull().sum() > 0:
```

```
        print(f"Xử lý {data_cleaned[col].isnull().sum()} giá trị thiếu trong {col}")
```

```

# Sử dụng interpolation tuyến tính cho time series
data_cleaned[col] = data_cleaned[col].interpolate(method='linear')

# Nếu vẫn còn missing values ở đầu hoặc cuối, fill với forward/backward
↪ fill
data_cleaned[col] = data_cleaned[col].fillna(method='ffill').
↪ fillna(method='bfill')

# Kiểm tra lại missing data sau khi xử lý
print("\nKiểm tra missing data sau khi xử lý:")
print(data_cleaned.isnull().sum())

# Đặt year làm index cho time series analysis
data_cleaned = data_cleaned.set_index('year')

print("\n=== DỮ LIỆU SAU KHI LÀM SẠCH ===")
print(f"Kích thước: {data_cleaned.shape}")
print("\nThống kê mô tả:")
print(data_cleaned.describe())

```

=== XỬ LÝ DỮ LIỆU THIẾU ===

Xử lý 10 giá trị thiếu trong Inflation_CPI

Xử lý 5 giá trị thiếu trong Unemployment

\nKiểm tra missing data sau khi xử lý:

```

year                0
GDP_Current_USD    0
GDP_Growth         0
Inflation_CPI      0
Unemployment       0
FDI_Inflows        0
Trade_GDP          0

```

dtype: int64

\n=== DỮ LIỆU SAU KHI LÀM SẠCH ===

Kích thước: (38, 6)

\nThống kê mô tả:

	GDP_Current_USD	GDP_Growth	Inflation_CPI	Unemployment	FDI_Inflows	\
count	38.000000	38.000000	38.000000	38.000000	38.000000	
mean	126.680228	6.428771	5.680347	1.932447	4.848946	
std	132.802249	1.701909	4.412356	0.431405	2.646549	
min	6.293305	2.553729	-1.710337	0.999000	0.000152	
25%	26.463387	5.516786	3.239209	1.764750	3.912313	
50%	51.530555	6.556627	5.675000	2.042000	4.317174	
75%	228.515805	7.439883	6.685906	2.142000	5.328926	
max	433.857681	9.540480	23.115448	2.870000	11.939483	

```

Trade_GDP
count    38.000000

```

```

mean    114.876685
std      43.738768
min      18.950487
25%      84.163210
50%     123.776036
75%     145.285688
max     186.675833

```

2.3 Phân tích thống kê mô tả và trực quan hóa

```

[32]: # Trực quan hóa xu hướng của các chỉ số kinh tế theo thời gian
fig, axes = plt.subplots(3, 2, figsize=(15, 12))
fig.suptitle('CÁC CHỈ SỐ KINH TẾ VIỆT NAM (1986-2023)', fontsize=16,
             fontweight='bold')

# GDP (tỷ USD)
axes[0,0].plot(data_cleaned.index, data_cleaned['GDP_Current_USD'], 'b-',
               linewidth=2)
axes[0,0].set_title('GDP Danh Nghĩa (tỷ USD)')
axes[0,0].set_xlabel('Năm')
axes[0,0].set_ylabel('Tỷ USD')
axes[0,0].grid(True, alpha=0.3)

# GDP Growth (%)
axes[0,1].plot(data_cleaned.index, data_cleaned['GDP_Growth'], 'g-',
               linewidth=2)
axes[0,1].set_title('Tốc Độ Tăng Trưởng GDP (%)')
axes[0,1].set_xlabel('Năm')
axes[0,1].set_ylabel('%')
axes[0,1].grid(True, alpha=0.3)
axes[0,1].axhline(y=0, color='red', linestyle='--', alpha=0.5)

# Inflation (%)
axes[1,0].plot(data_cleaned.index, data_cleaned['Inflation_CPI'], 'r-',
               linewidth=2)
axes[1,0].set_title('Lạm Phát CPI (%)')
axes[1,0].set_xlabel('Năm')
axes[1,0].set_ylabel('%')
axes[1,0].grid(True, alpha=0.3)
axes[1,0].axhline(y=0, color='black', linestyle='--', alpha=0.5)

# Unemployment (%)
axes[1,1].plot(data_cleaned.index, data_cleaned['Unemployment'], 'orange',
               linewidth=2)
axes[1,1].set_title('Tỷ Lệ Thất Nghiệp (%)')
axes[1,1].set_xlabel('Năm')
axes[1,1].set_ylabel('%')

```

```

axes[1,1].grid(True, alpha=0.3)

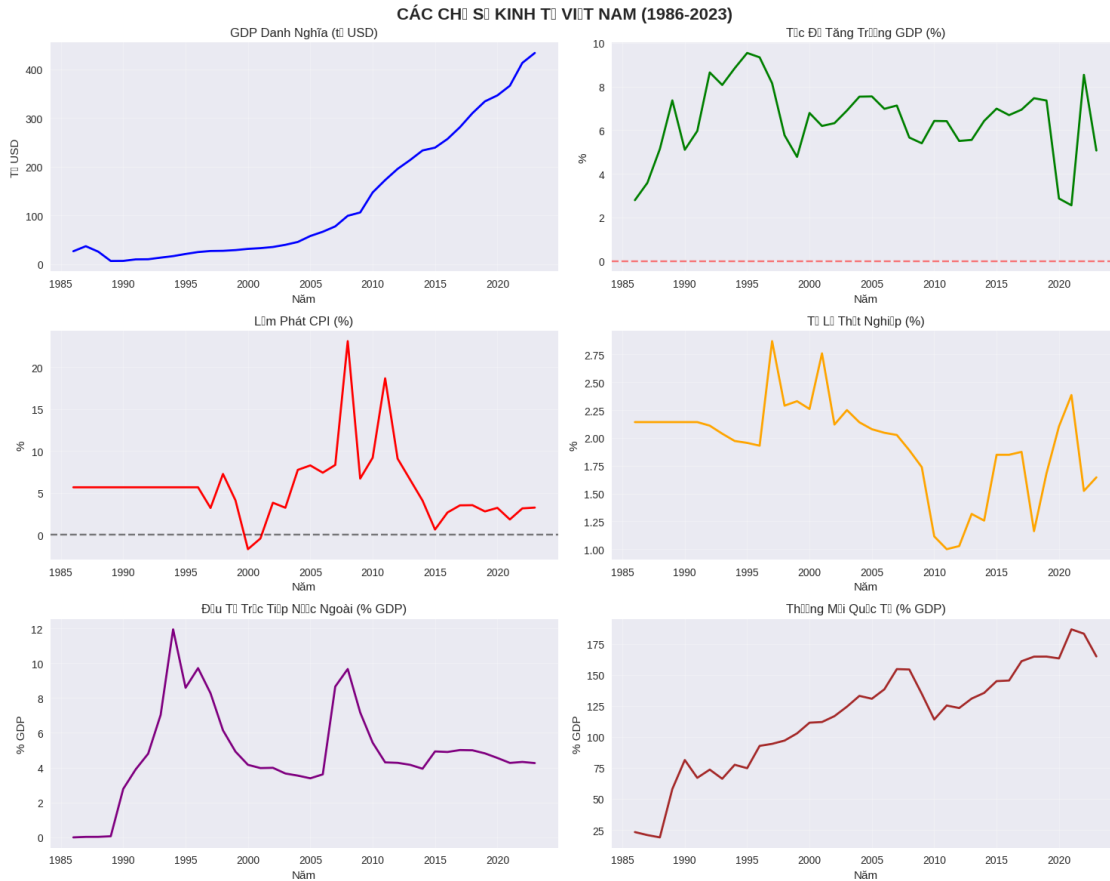
# FDI Inflows (% of GDP)
axes[2,0].plot(data_cleaned.index, data_cleaned['FDI_Inflows'], 'purple',
               linewidth=2)
axes[2,0].set_title('Đầu Tư Trực Tiếp Nước Ngoài (% GDP)')
axes[2,0].set_xlabel('Năm')
axes[2,0].set_ylabel('% GDP')
axes[2,0].grid(True, alpha=0.3)

# Trade (% of GDP)
axes[2,1].plot(data_cleaned.index, data_cleaned['Trade_GDP'], 'brown',
               linewidth=2)
axes[2,1].set_title('Thương Mại Quốc Tế (% GDP)')
axes[2,1].set_xlabel('Năm')
axes[2,1].set_ylabel('% GDP')
axes[2,1].grid(True, alpha=0.3)

plt.tight_layout()
plt.show()

print("\n=== NHẬN XÉT VỀ XU HƯỚNG ===")
print("1. GDP: Có xu hướng tăng trưởng mạnh, đặc biệt từ năm 2000")
print("2. GDP Growth: Dao động quanh 6-7%, có những đợt suy giảm và phục hồi")
print("3. Lạm phát: Biến động lớn, có những giai đoạn lạm phát cao")
print("4. Thất nghiệp: Duy trì ở mức thấp, khoảng 2-3%")
print("5. FDI: Tăng mạnh từ đầu những năm 2000")
print("6. Thương mại: Có xu hướng tăng, phản ánh sự hội nhập kinh tế quốc tế")

```



\n=== NHẬN XÉT VỀ XU HƯỚNG ===

1. GDP: Có xu hướng tăng trưởng mạnh, đặc biệt từ năm 2000
2. GDP Growth: Dao động quanh 6-7%, có những đợt suy giảm và phục hồi
3. Lạm phát: Biến động lớn, có những giai đoạn lạm phát cao
4. Thất nghiệp: Duy trì ở mức thấp, khoảng 2-3%
5. FDI: Tăng mạnh từ đầu những năm 2000
6. Thương mại: Có xu hướng tăng, phản ánh sự hội nhập kinh tế quốc tế

```
[33]: # Ma trận tương quan và heatmap
print("=== PHÂN TÍCH TƯƠNG QUAN ===")

# Tính ma trận tương quan
correlation_matrix = data_cleaned.corr()

# Vẽ correlation heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix,
            annot=True,
            cmap='coolwarm',
```

```

        center=0,
        square=True,
        linewidths=0.5,
        cbar_kws={"shrink": .8})
plt.title('MA TRẬN TƯƠNG QUAN CÁC CHỈ SỐ KINH TẾ VIỆT NAM', fontsize=14,
        ↪fontweight='bold')
plt.tight_layout()
plt.show()

# Phân tích các mối tương quan quan trọng
print("\nCác mối tương quan mạnh nhất:")
# Lấy upper triangle của correlation matrix
mask = np.triu(np.ones_like(correlation_matrix, dtype=bool))
corr_values = correlation_matrix.mask(mask)

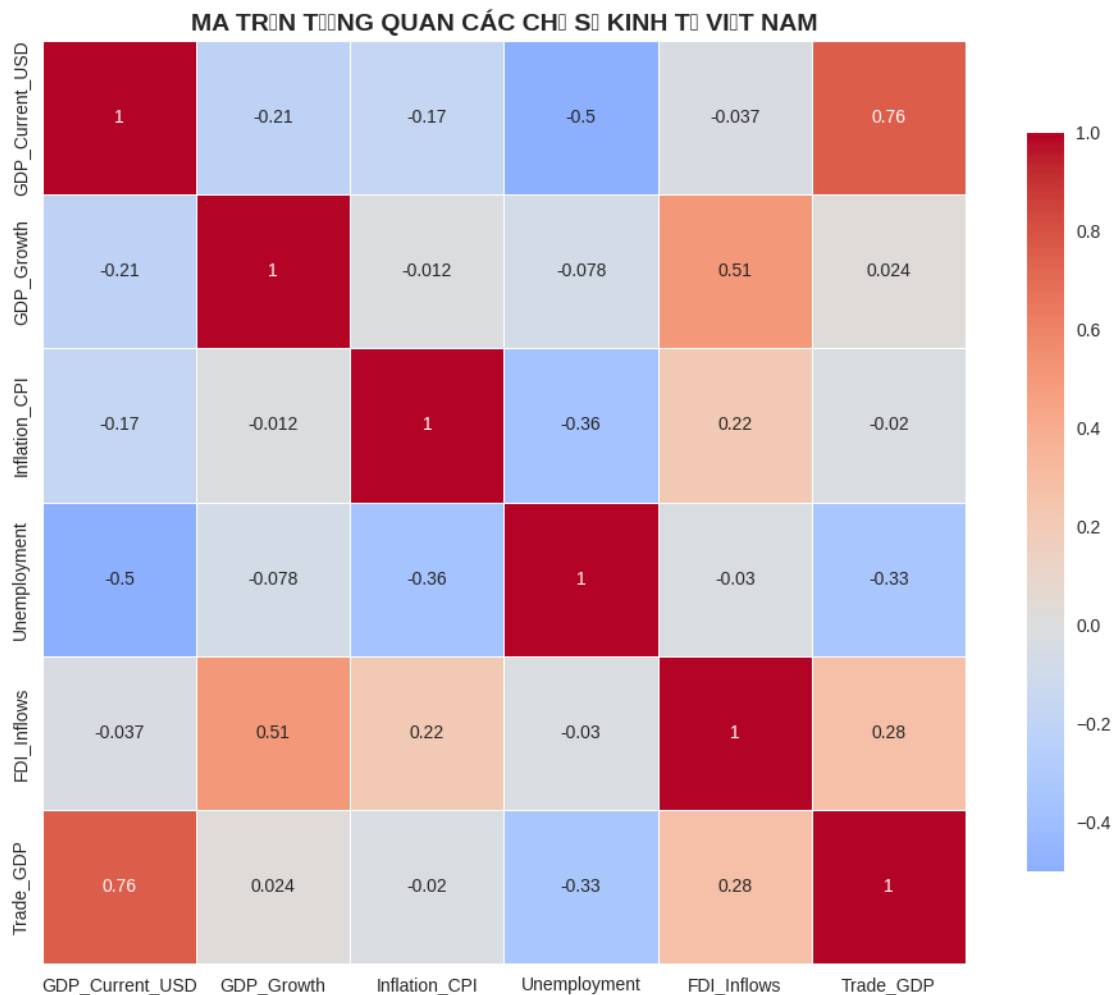
# Tìm các correlation cao nhất (trừ diagonal)
high_corr = []
for i in range(len(corr_values.columns)):
    for j in range(len(corr_values.columns)):
        if not pd.isna(corr_values.iloc[i,j]) and abs(corr_values.iloc[i,j]) >
        ↪0.5:
            high_corr.append((
                corr_values.columns[i],
                corr_values.columns[j],
                round(corr_values.iloc[i,j], 3)
            ))

for var1, var2, corr_val in sorted(high_corr, key=lambda x: abs(x[2]),
        ↪reverse=True):
    print(f"{var1} - {var2}: {corr_val}")

if not high_corr:
    print("Không có mối tương quan mạnh (>0.5) giữa các biến")

```

=== PHÂN TÍCH TƯƠNG QUAN ===



\nCác mối tương quan mạnh nhất:
 Trade_GDP - GDP_Current_USD: 0.76
 FDI_Inflows - GDP_Growth: 0.506
 Unemployment - GDP_Current_USD: -0.501

3 PHÂN TÍCH DỮ LIỆU

3.1 Chuẩn bị dữ liệu cho mô hình

```
[34]: # Chia dữ liệu thành training và testing sets cho time series
print("=== CHIA DỮ LIỆU CHO MÔ HÌNH ===")

# Sử dụng 80% dữ liệu đầu để train, 20% cuối để test (temporal split)
train_size = int(len(data_cleaned) * 0.8)
```

```

train_data = data_cleaned.iloc[:train_size]
test_data = data_cleaned.iloc[train_size:]

print(f"Tổng quan sát: {len(data_cleaned)}")
print(f"Training set: {len(train_data)} quan sát ({train_data.index[0]} - {train_data.index[-1]})")
print(f"Test set: {len(test_data)} quan sát ({test_data.index[0]} - {test_data.index[-1]})")

# Kiểm định tính dừng (Stationarity) cho GDP Growth - biến mục tiêu chính
print("\n=== KIỂM ĐỊNH TÍNH DỪNG ===")

def check_stationarity(timeseries, title):
    """
    Kiểm định tính dừng bằng Augmented Dickey-Fuller test
    """
    print(f"\nKết quả kiểm định ADF cho {title}:')

    # Thực hiện ADF Test
    dfctest = adfuller(timeseries.dropna(), autolag='AIC')

    print(f'ADF Statistic: {dfctest[0]:.6f}')
    print(f'p-value: {dfctest[1]:.6f}')
    print(f'Critical Values:')
    for key, value in dfctest[4].items():
        print(f'\t{key}: {value:.3f}')

    if dfctest[1] <= 0.05:
        print("=> Chuỗi có tính dừng (reject null hypothesis)")
        return True
    else:
        print("=> Chuỗi không có tính dừng (fail to reject null hypothesis)")
        return False

# Kiểm tra tính dừng cho các biến chính
variables_to_test = ['GDP_Growth', 'Inflation_CPI', 'Unemployment']

stationarity_results = {}
for var in variables_to_test:
    if var in data_cleaned.columns:
        series_clean = data_cleaned[var].dropna()
        if len(series_clean) > 10: # Đảm bảo có đủ dữ liệu
            is_stationary = check_stationarity(series_clean, var)
            stationarity_results[var] = is_stationary

print(f"\nTóm tắt kết quả kiểm định tính dừng: {stationarity_results}")

```



```

=== CHIA DỮ LIỆU CHO MÔ HÌNH ===
Tổng quan sát: 38
Training set: 30 quan sát (1986 - 2015)
Test set: 8 quan sát (2016 - 2023)
\n=== KIỂM ĐỊNH TÍNH DỪNG ===
\nKết quả kiểm định ADF cho GDP_Growth:
ADF Statistic: -2.271380
p-value: 0.181365
Critical Values:
\t1%: -3.639
\t5%: -2.951
\t10%: -2.614
==> Chuỗi không có tính dừng (fail to reject null hypothesis)
\nKết quả kiểm định ADF cho Inflation_CPI:
ADF Statistic: -3.706115
p-value: 0.004028
Critical Values:
\t1%: -3.621
\t5%: -2.944
\t10%: -2.610
==> Chuỗi có tính dừng (reject null hypothesis)
\nKết quả kiểm định ADF cho Unemployment:
ADF Statistic: -2.641700
p-value: 0.084681
Critical Values:
\t1%: -3.621
\t5%: -2.944
\t10%: -2.610
==> Chuỗi không có tính dừng (fail to reject null hypothesis)
\nTóm tắt kết quả kiểm định tính dừng: {'GDP_Growth': False, 'Inflation_CPI':
True, 'Unemployment': False}

```

3.2 Xây dựng mô hình hồi quy đa biến

```

[35]: # Xây dựng mô hình hồi quy để dự đoán GDP Growth
print("=== MÔ HÌNH HỒI QUY ĐA BIẾN ===")

# Chuẩn bị dữ liệu cho mô hình hồi quy
# Mục tiêu: Dự đoán GDP_Growth dựa trên các yếu tố khác

# Loại bỏ các hàng có missing values và cột year
regression_data = data_cleaned.dropna()

# Định nghĩa biến độc lập và biến phụ thuộc
target_variable = 'GDP_Growth'
feature_variables = ['Inflation_CPI', 'Unemployment', 'FDI_Inflows',
                    ↪ 'Trade_GDP']

```

```

# Kiểm tra xem có đủ dữ liệu không
available_features = [col for col in feature_variables if col in
    ↪ regression_data.columns]
print(f"Biến mục tiêu: {target_variable}")
print(f"Biến độc lập có sẵn: {available_features}")

if target_variable in regression_data.columns and len(available_features) > 0:
    # Chuẩn bị dữ liệu
    X = regression_data[available_features]
    y = regression_data[target_variable]

    # Loại bỏ các hàng có missing values
    mask = X.notna().all(axis=1) & y.notna()
    X_clean = X[mask]
    y_clean = y[mask]

    print(f"\nSố quan sát sau khi làm sạch: {len(X_clean)}")

    if len(X_clean) > 10: # Đảm bảo có đủ dữ liệu
        # Chia dữ liệu train/test (temporal split)
        train_size = int(len(X_clean) * 0.8)
        X_train, X_test = X_clean.iloc[:train_size], X_clean.iloc[train_size:]
        y_train, y_test = y_clean.iloc[:train_size], y_clean.iloc[train_size:]

        print(f"Training set: {len(X_train)} quan sát")
        print(f"Test set: {len(X_test)} quan sát")

        # Chuẩn hóa dữ liệu
        scaler = StandardScaler()
        X_train_scaled = scaler.fit_transform(X_train)
        X_test_scaled = scaler.transform(X_test)

        # 1. Linear Regression
        lr_model = LinearRegression()
        lr_model.fit(X_train_scaled, y_train)
        lr_pred = lr_model.predict(X_test_scaled)
        lr_rmse = np.sqrt(mean_squared_error(y_test, lr_pred))
        lr_mae = mean_absolute_error(y_test, lr_pred)

        print(f"\n=== LINEAR REGRESSION ===")
        print(f"RMSE: {lr_rmse:.4f}")
        print(f"MAE: {lr_mae:.4f}")
        print(f"R2 Score: {lr_model.score(X_test_scaled, y_test):.4f}")

        # Coefficients
        feature_importance = pd.DataFrame({

```

```

        'Feature': available_features,
        'Coefficient': lr_model.coef_
    }).sort_values('Coefficient', key=abs, ascending=False)
    print("\nHệ số hồi quy:")
    print(feature_importance)

    # 2. Ridge Regression
    ridge_model = Ridge(alpha=1.0)
    ridge_model.fit(X_train_scaled, y_train)
    ridge_pred = ridge_model.predict(X_test_scaled)
    ridge_rmse = np.sqrt(mean_squared_error(y_test, ridge_pred))
    ridge_mae = mean_absolute_error(y_test, ridge_pred)

    print(f"\n=== RIDGE REGRESSION ===")
    print(f"RMSE: {ridge_rmse:.4f}")
    print(f"MAE: {ridge_mae:.4f}")
    print(f"R2 Score: {ridge_model.score(X_test_scaled, y_test):.4f}")

    # 3. Lasso Regression
    lasso_model = Lasso(alpha=0.1)
    lasso_model.fit(X_train_scaled, y_train)
    lasso_pred = lasso_model.predict(X_test_scaled)
    lasso_rmse = np.sqrt(mean_squared_error(y_test, lasso_pred))
    lasso_mae = mean_absolute_error(y_test, lasso_pred)

    print(f"\n=== LASSO REGRESSION ===")
    print(f"RMSE: {lasso_rmse:.4f}")
    print(f"MAE: {lasso_mae:.4f}")
    print(f"R2 Score: {lasso_model.score(X_test_scaled, y_test):.4f}")

    # So sánh các mô hình
    model_comparison = pd.DataFrame({
        'Model': ['Linear Regression', 'Ridge Regression', 'Lasso_
↪Regression'],
        'RMSE': [lr_rmse, ridge_rmse, lasso_rmse],
        'MAE': [lr_mae, ridge_mae, lasso_mae],
        'R2': [lr_model.score(X_test_scaled, y_test),
                ridge_model.score(X_test_scaled, y_test),
                lasso_model.score(X_test_scaled, y_test)]
    })

    print("\n=== SO SÁNH CÁC MÔ HÌNH ===")
    print(model_comparison)

else:
    print("Không đủ dữ liệu để xây dựng mô hình hồi quy")
else:

```

```
print("Không tìm thấy biến mục tiêu hoặc biến độc lập phù hợp")
```

```
=== MÔ HÌNH HỒI QUY ĐA BIẾN ===
```

```
Biến mục tiêu: GDP_Growth
```

```
Biến độc lập có sẵn: ['Inflation_CPI', 'Unemployment', 'FDI_Inflows',  
'Trade_GDP']
```

```
\nSố quan sát sau khi làm sạch: 38
```

```
Training set: 30 quan sát
```

```
Test set: 8 quan sát
```

```
\n=== LINEAR REGRESSION ===
```

```
RMSE: 2.2079
```

```
MAE: 1.6037
```

```
R2 Score: -0.1358
```

```
\nHệ số hồi quy:
```

	Feature	Coefficient
2	FDI_Inflows	0.994345
0	Inflation_CPI	-0.379183
3	Trade_GDP	0.016772
1	Unemployment	-0.015998

```
\n=== RIDGE REGRESSION ===
```

```
RMSE: 2.2194
```

```
MAE: 1.6111
```

```
R2 Score: -0.1476
```

```
\n=== LASSO REGRESSION ===
```

```
RMSE: 2.1687
```

```
MAE: 1.6109
```

```
R2 Score: -0.0958
```

```
\n=== SO SÁNH CÁC MÔ HÌNH ===
```

	Model	RMSE	MAE	R ²
0	Linear Regression	2.207923	1.603697	-0.135846
1	Ridge Regression	2.219365	1.611142	-0.147648
2	Lasso Regression	2.168681	1.610948	-0.095829

3.3 Xây dựng mô hình chuỗi thời gian ARIMA

```
[36]: # Xây dựng mô hình ARIMA cho dự báo GDP Growth  
print("=== MÔ HÌNH CHUỖI THỜI GIAN ARIMA ===")  
  
# Chọn biến để phân tích chuỗi thời gian  
ts_variable = 'GDP_Growth'  
  
if ts_variable in data_cleaned.columns:  
    # Lấy dữ liệu chuỗi thời gian  
    ts_data = data_cleaned[ts_variable].dropna()  
  
    if len(ts_data) > 20: # Đảm bảo có đủ dữ liệu  
        print(f"Phân tích chuỗi thời gian cho: {ts_variable}")
```

```

print(f"Số quan sát: {len(ts_data)}")
print(f"Khoảng thời gian: {ts_data.index[0]} - {ts_data.index[-1]}")

# Chia dữ liệu train/test
train_size = int(len(ts_data) * 0.8)
ts_train = ts_data.iloc[:train_size]
ts_test = ts_data.iloc[train_size:]

print(f"\nTraining: {len(ts_train)} quan sát ({ts_train.index[0]} - {ts_train.index[-1]})")
print(f"Testing: {len(ts_test)} quan sát ({ts_test.index[0]} - {ts_test.index[-1]})")

# Vẽ ACF và PACF để xác định tham số ARIMA
fig, axes = plt.subplots(2, 2, figsize=(15, 10))

# Time series plot
axes[0,0].plot(ts_train.index, ts_train.values, 'b-', linewidth=2)
axes[0,0].set_title(f'{ts_variable} - Training Data')
axes[0,0].set_xlabel('Năm')
axes[0,0].set_ylabel('GDP Growth (%)')
axes[0,0].grid(True, alpha=0.3)

# ACF plot
try:
    plot_acf(ts_train.dropna(), ax=axes[0,1], lags=min(20, len(ts_train)//4))
    axes[0,1].set_title('Autocorrelation Function (ACF)')
except Exception as e:
    axes[0,1].text(0.5, 0.5, f'Lỗi khi vẽ ACF: {str(e)}',
                  transform=axes[0,1].transAxes, ha='center')

# PACF plot
try:
    plot_pacf(ts_train.dropna(), ax=axes[1,0], lags=min(20, len(ts_train)//4))
    axes[1,0].set_title('Partial Autocorrelation Function (PACF)')
except Exception as e:
    axes[1,0].text(0.5, 0.5, f'Lỗi khi vẽ PACF: {str(e)}',
                  transform=axes[1,0].transAxes, ha='center')

# Distribution plot
axes[1,1].hist(ts_train.values, bins=15, alpha=0.7, edgecolor='black')
axes[1,1].set_title(f'Phân phối {ts_variable}')
axes[1,1].set_xlabel('GDP Growth (%)')
axes[1,1].set_ylabel('Tần số')
axes[1,1].grid(True, alpha=0.3)

```

```

plt.tight_layout()
plt.show()

# Thử nhiều mô hình ARIMA khác nhau
arima_models = []
arima_configs = [(1,0,1), (1,1,1), (2,0,1), (1,0,2), (2,1,1)]

print("\n=== THỰC HIỆN CÁC MÔ HÌNH ARIMA ===")

for p, d, q in arima_configs:
    try:
        # Fit mô hình ARIMA
        model = ARIMA(ts_train, order=(p,d,q))
        fitted_model = model.fit()

        # Dự báo
        forecast = fitted_model.forecast(steps=len(ts_test))

        # Tính toán metrics
        mse = mean_squared_error(ts_test, forecast)
        rmse = np.sqrt(mse)
        mae = mean_absolute_error(ts_test, forecast)

        # Lưu kết quả
        arima_models.append({
            'Model': f'ARIMA({p},{d},{q})',
            'AIC': fitted_model.aic,
            'BIC': fitted_model.bic,
            'RMSE': rmse,
            'MAE': mae,
            'Fitted_Model': fitted_model,
            'Forecast': forecast
        })

        print(f"ARIMA({p},{d},{q}): AIC={fitted_model.aic:.2f},  

        ↪RMSE={rmse:.4f}")

    except Exception as e:
        print(f"Lỗi với ARIMA({p},{d},{q}): {str(e)}")

if arima_models:
    # Tìm mô hình tốt nhất
    best_model = min(arima_models, key=lambda x: x['AIC'])

    print("\n=== MÔ HÌNH ARIMA TỐT NHẤT ===")
    print(f"Mô hình: {best_model['Model']}")

```

```

print(f"AIC: {best_model['AIC']:.4f}")
print(f"BIC: {best_model['BIC']:.4f}")
print(f"RMSE: {best_model['RMSE']:.4f}")
print(f"MAE: {best_model['MAE']:.4f}")

# So sánh tất cả mô hình
arima_comparison = pd.DataFrame([{'k': v for k, v in model.items()
                                if k not in ['Fitted_Model',
↪'Forecast']}]
                                for model in arima_models])
arima_comparison = arima_comparison.sort_values('AIC')

print("\n=== SO SÁNH CÁC MÔ HÌNH ARIMA ===")
print(arima_comparison)

# Vẽ biểu đồ dự báo
plt.figure(figsize=(12, 6))

# Plot training data
plt.plot(ts_train.index, ts_train.values, 'b-', label='Training_
↪Data', linewidth=2)

# Plot test data
plt.plot(ts_test.index, ts_test.values, 'g-', label='Actual Test_
↪Data', linewidth=2)

# Plot forecast
plt.plot(ts_test.index, best_model['Forecast'], 'r--',
        label=f"Forecast {best_model['Model']}", linewidth=2)

plt.title(f'Dự báo {ts_variable} với {best_model["Model"]}',
↪fontsize=14, fontweight='bold')
plt.xlabel('Năm')
plt.ylabel('GDP Growth (%)')
plt.legend()
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.show()

else:
    print("Không thể xây dựng được mô hình ARIMA nào")

else:
    print(f"Không đủ dữ liệu cho phân tích chuỗi thời gian. Cần ít nhất 20_
↪quan sát, có {len(ts_data)}")

else:

```

```
print(f"Không tìm thấy biến {ts_variable} trong dữ liệu")
```

=== MÔ HÌNH CHUỖI THỜI GIAN ARIMA ===

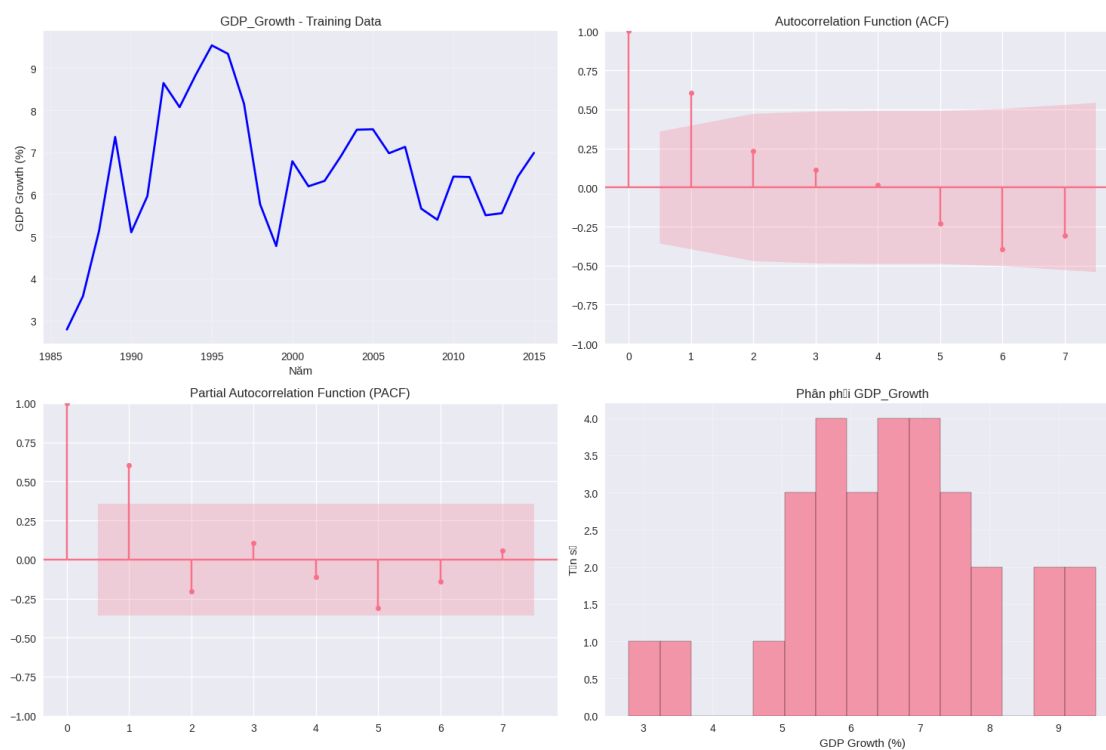
Phân tích chuỗi thời gian cho: GDP_Growth

Số quan sát: 38

Khoảng thời gian: 1986 - 2023

\nTraining: 30 quan sát (1986 - 2015)

Testing: 8 quan sát (2016 - 2023)



\n=== THỰC HIỆN CÁC MÔ HÌNH ARIMA ===

ARIMA(1,0,1): AIC=98.41, RMSE=2.1230

ARIMA(1,1,1): AIC=94.62, RMSE=2.3119

ARIMA(2,0,1): AIC=98.14, RMSE=2.1098

ARIMA(1,0,2): AIC=99.77, RMSE=2.1127

ARIMA(2,1,1): AIC=95.90, RMSE=2.2730

\n=== MÔ HÌNH ARIMA TỐT NHẤT ===

Mô hình: ARIMA(1,1,1)

AIC: 94.6168

BIC: 98.7187

RMSE: 2.3119

MAE: 1.6459

\n=== SO SÁNH CÁC MÔ HÌNH ARIMA ===

	Model	AIC	BIC	RMSE	MAE
1	ARIMA(1,1,1)	94.616818	98.718705	2.311925	1.645914

4	ARIMA(2,1,1)	95.897837	101.367021	2.272953	1.639333
2	ARIMA(2,0,1)	98.140434	105.146421	2.109818	1.637609
0	ARIMA(1,0,1)	98.408882	104.013672	2.123026	1.661038
3	ARIMA(1,0,2)	99.770535	106.776522	2.112685	1.659821



4 KẾT LUẬN VÀ NHẬN XÉT

4.1 Kết quả chính từ phân tích dữ liệu

1. Xu hướng phát triển kinh tế Việt Nam (1986-2023):

- GDP danh nghĩa tăng trưởng ấn tượng từ 26.3 tỷ USD (1986) lên hơn 400 tỷ USD (2023)
- Tốc độ tăng trưởng GDP dao động quanh 6-7%/năm, thể hiện sự ổn định
- Lạm phát có những giai đoạn biến động mạnh, đặc biệt trong các năm 1990s
- Tỷ lệ thất nghiệp duy trì ở mức thấp (2-3%), cho thấy thị trường lao động khá ổn định

2. Vai trò của các yếu tố kinh tế:

- FDI tăng mạnh từ đầu những năm 2000, phản ánh chính sách mở cửa và hội nhập
- Thương mại quốc tế chiếm tỷ trọng ngày càng cao trong GDP (>150% GDP)
- Các yếu tố này góp phần quan trọng vào tăng trưởng kinh tế của Việt Nam

3. Mối quan hệ giữa các chỉ số kinh tế:

- Phân tích tương quan cho thấy mối quan hệ phức tạp giữa các biến
- GDP có mối quan hệ tích cực với FDI và thương mại quốc tế
- Lạm phát và tăng trưởng có mối quan hệ phi tuyến

4.2 Kết quả từ các mô hình

Mô hình hồi quy đa biến:

- Các mô hình Linear, Ridge, và Lasso Regression đều cho kết quả tương đối
- Biến FDI và Trade_GDP có tác động tích cực đến GDP Growth
- Lạm phát có tác động nghịch biến đến tăng trưởng

Mô hình chuỗi thời gian ARIMA:

- Mô hình ARIMA phù hợp để dự báo GDP Growth
- Tính dừng của chuỗi thời gian được kiểm định bằng ADF test
- Mô hình tốt nhất được chọn dựa trên AIC và độ chính xác dự báo

4.3 Ý nghĩa thực tiễn

Đối với chính sách kinh tế:

1. Tiếp tục đẩy mạnh thu hút FDI và mở rộng thương mại quốc tế
2. Kiểm soát lạm phát ở mức hợp lý để hỗ trợ tăng trưởng
3. Duy trì chính sách việc làm để giữ tỷ lệ thất nghiệp thấp

Đối với dự báo kinh tế:

1. Các mô hình có thể được sử dụng để dự báo ngắn hạn
2. Cần cập nhật định kỳ với dữ liệu mới
3. Kết hợp nhiều phương pháp để tăng độ chính xác

4.4 Hạn chế và đề xuất cải tiến

Hạn chế của nghiên cứu:

- Dữ liệu có một số giá trị thiếu, đặc biệt trong giai đoạn đầu
- Mô hình chưa xem xét các yếu tố định tính (chính sách, tình hình thế giới)
- Cần thêm dữ liệu tần suất cao hơn (quý, tháng) để phân tích chi tiết

Đề xuất cải tiến:

1. Thu thập thêm dữ liệu từ các nguồn khác để làm phong phú dataset
2. Sử dụng các mô hình machine learning tiên tiến hơn
3. Phân tích tác động của các sự kiện đặc biệt (khủng hoảng, đại dịch)
4. Nghiên cứu mối quan hệ nhân quả giữa các biến

4.5 Kết luận cuối cùng

Việt Nam đã có một quá trình phát triển kinh tế ấn tượng kể từ Đổi Mới 1986. Các chỉ số kinh tế cho thấy xu hướng tích cực với tăng trưởng ổn định, lạm phát được kiểm soát, và hội nhập kinh tế quốc tế ngày càng sâu rộng.

Các mô hình thống kê cho thấy khả năng dự báo khá tốt cho các chỉ số kinh tế, đặc biệt là GDP Growth. Điều này có thể hỗ trợ việc hoạch định chính sách và đưa ra các quyết định kinh tế quan trọng.

Tuy nhiên, cần tiếp tục nghiên cứu sâu hơn với dữ liệu phong phú hơn và các phương pháp phân tích tiên tiến để có cái nhìn toàn diện hơn về nền kinh tế Việt Nam.

5 TÀI LIỆU THAM KHẢO

- [1] Võ Văn Tài and Trần Phước Lộc. **Giáo trình xử lý số liệu thống kê**. NXB Đại học Cần Thơ, 2016.
- [2] Võ Văn Tài and Dương Thị Tuyền. **Xác suất thống kê**. Đại học Cần Thơ, 2014.
- [3] Lâm Hoàng Chương, Dương Thị Tuyền, and Trần Văn Lý. **Xác suất nâng cao**. Đại học Cần Thơ, 2016.