

# Analiza cen akcji wybranych spółek WIG20

## Wstęp

Celem tego projektu jest analiza i prognozowanie cen akcji spółek notowanych na indeksie WIG20 na Giełdzie Papierów Wartościowych w Warszawie. WIG20 to indeks 20 największych i najbardziej płynnych spółek na GPW. Analiza tych spółek może dostarczyć cennych informacji o ogólnych trendach na polskim rynku giełdowym.

W tym projekcie skupię się na analizie historycznych danych cenowych tych spółek, które pobieram ze strony *stoq.pl*. Dane te obejmują codzienne ceny otwarcia, zamknięcia, najwyższe, najniższe oraz wolumen obrotu.

Do analizy danych i prezentacji wyników wykorzystam skrypty oraz wykresy napisane i wygenerowane przy użyciu języka programowania *Python*.

Moim celem jest dobranie odpowiednich modeli do danych, a następnie zastosowanie tych informacji do prognozowania przyszłych cen akcji. W tym celu będę korzystać z różnych technik analizy danych i modelowania predykcyjnego.

Oczywiście zdaję sobie sprawę z faktu, że na ceny akcji ma wpływ wiele czynników i nie powinno się opierać jedynie na danych historycznych, dlatego też poniższy projekt traktuję czysto hobbistycznie.

## 1. Pobieranie danych

Dane do analizy pobieram ze strony *stooq.pl*, która udostępnia historyczne ceny akcji spółek giełdowych. Pozyskuję dane przy użyciu skryptu *get\_wig20.py*, wykorzystując bibliotekę *pandas*.

```
# get_wig20.py

import pandas as pd

def download_data(symbol, start_date, end_date):
    url = f"https://stoq.pl/q/d/1/?s={symbol}&d1={start_date}&d2={end_date}&i=d"
    data = pd.read_csv(url)
    data['Symbol'] = symbol # Dodajemy kolumnę z symbolem spółki
    return data

# Zahardkodowana lista symboli spółek
symbols = ['ALR', 'CCC', 'CDR', 'CPS', 'DNP', 'JSW', 'KGH', 'LTS', 'LPP', 'MBK', 'OPL', 'PEO',
           'PGE', 'PGN', 'PKN', 'PKO', 'PLY', 'PZU', 'SPL', 'TPE']

# Pobranie danych dla każdej spółki
all_data = []
for symbol in symbols:
    data = download_data(symbol, "20230101", "20240115")
    all_data.append(data)

# Połączenie wszystkich danych w jedną ramkę danych i zapis
all_data = pd.concat(all_data)

filename = "wig20_all.csv"
all_data.to_csv(filename, index=False)
```

*Tabela 1: Pobrane dane*

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

Współpraca

## 2. Czyszczenie danych

Sprawdzam zgodność typów oraz czy w danych nie pojawiły się luki. Wśród pobranych informacji są również dane, które okazują się zbędne na potrzeby mojej prostej analizy. Przekształcam kolumny, żeby umożliwić wygodną analizę. Oczyszczone dane zapisuję w nowym pliku *cleaned\_data.csv*.

```
#clean_data.py

import pandas as pd

data = pd.read_csv("wig20_all.csv")

#Sprawdzenie typów danych i braków danych

print(data.dtypes)

data['Data'] = pd.to_datetime(data['Data'])

print(data.dtypes)

print(data.isnull().sum())

#data = data.drop(columns=['Brak danych', 'Wolumen', 'Najnizszy', 'Otwarcie', 'Najwyzszy'])

#Przekształcenie symboli spółek w kolumny

data_pivot = data.pivot(index='Data', columns='Symbol', values='Zamkniecie')

data_pivot.to_csv('cleaned data.csv')
```

Tabela 2: Dane po "oczyszczeniu"

[illegible]

### 3. Wybór słabo skorelowanych spółek

Wybieram spółki słabo skorelowane, ponieważ takie podejście pozwoli mi zdywersyfikować ewentualny portfel inwestycyjny.

Stosuję do tego algorytm, który iteracyjnie dodaje spółki do zbioru, jeśli są słabo skorelowane ze wszystkimi spółkami, które już dodaliśmy. Nie gwarantuje on co prawda znalezienia maksymalnego zbioru spółek słabo skorelowanych, ale jest za to prosty i szybki.

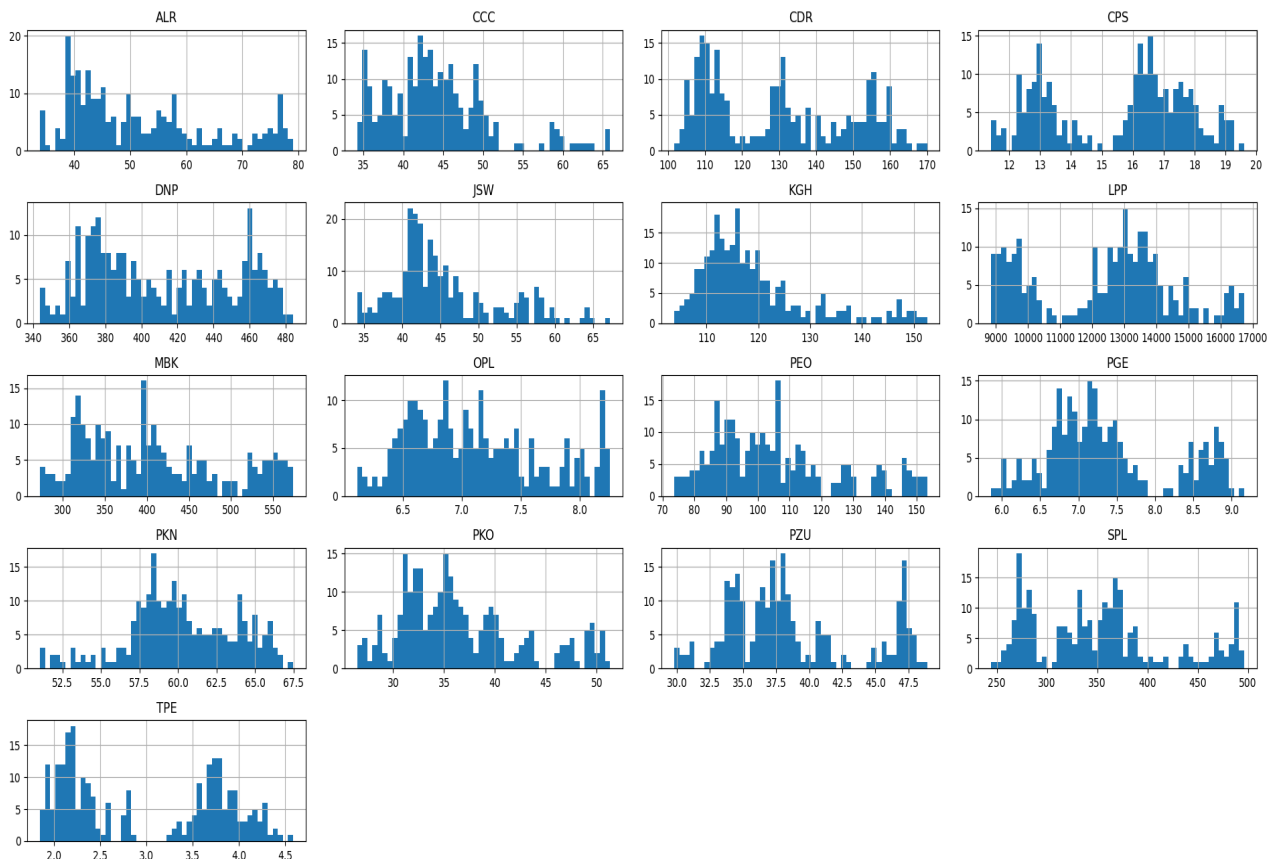
```
#eda.py

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

data = pd.read_csv("cleaned_data.csv")

data.hist(bins=50, figsize=(15,10))
plt.show()
```

*Tabela 3: Wizualizacja rozkładu danych*



```
#Tworzenie macierzy korelacji z typów numerycznych zawartych w ramce danych
numeric_data = data.select_dtypes(include=[np.number])

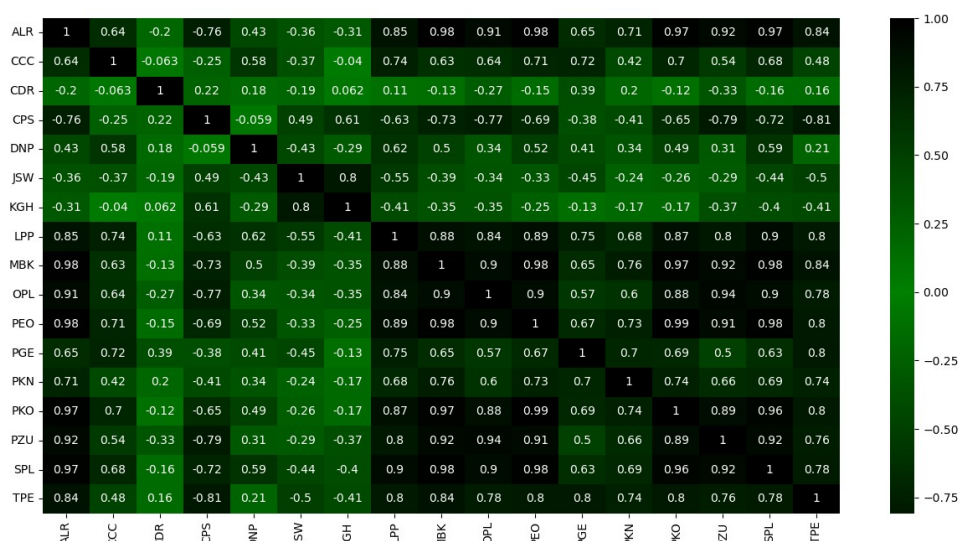
corr_matrix = numeric_data.corr()

threshold = 0.5

# Tworzenie i wyświetlanie mapy ciepła
sns.heatmap(corr_matrix, annot=True, cmap=cmap, center=0)

plt.show()
```

*Tabela 4: Macierz korelacji*



```
# Wybieranie spółek słabo skorelowanych
weak_corr = []

for i in range(len(corr_matrix.columns)):
    for j in range(i+1, len(corr_matrix.columns)):
        if abs(corr_matrix.iloc[i,j])<threshold:
            weak_corr.append((corr_matrix.columns[i],corr_matrix.columns[j]))

selected_stock = set()

for stock1, stock2 in weak_corr:
    if all(abs(corr_matrix.loc[stock, stock1]) < threshold and abs(corr_matrix.loc[stock, stock2]) <
threshold for stock in selected_stock):
        selected_stock.add(stock1)
        selected_stock.add(stock2)

# Tworzenie nowej ramki danych, zawierającej tylko wybrane kolumny i zapis nowego pliku CSV
selected_stock_df = data[["Data"]+list(selected_stock)]

selected_stock_df.to_csv("portfel.csv", index=False)
```

## 4. Dopasowanie modelu

Dopasowanie modelu zaczynam od podziału danych na dwa zbiory: treningowy (80%) i testowy (20%) (z zachowaniem porządku chronologicznego, żeby móc ocenić sprawdzalność prognozy). Następnie skorzystam z modelu *ARIMA* (AutoRegressive Integrated Moving Average) z parametrami dobranymi metodą *grid search* tak, aby zminimalizować błąd średniokwadratowy.

```
#model.py
from sklearn.metrics import mean_squared_error
import pandas as pd
from statsmodels.tsa.arima.model import ARIMA

import numpy as np

df = pd.read_csv("portfel.csv")

# Podział danych na zbiory treningowy i testowy
train_size = 0.8
train_num = int(df.shape[0] * train_size)
train_df = df.iloc[:train_num]
test_df = df.iloc[train_num:]

#Grid search
predictions = {}
tests = {}
min_mse={}
best_params={}
columns = train_df.columns[1:]

p_range = range(0,6)
d_range = range(0,6)
q_range = range(0,6)

for column in columns:
    y_train = train_df[column]
    y_test = test_df[column]

    best_params[column] = None
    lowest_mse[column] = np.inf

    for p in p_range:
        for d in d_range:
            for q in q_range:
                try:
                    model = ARIMA(y_train, order=(p,d,q))
                    model_fit = model.fit()

                    y_pred = model_fit.predict(start=len(train_df),end=len(train_df)+len(test_df)-1)

                    mse = mean_squared_error(y_test, y_pred)

                    if mse < lowest_mse[column]:
                        best_params[column] = (p, d, q)
                        lowest_mse[column] = mse
                except:
                    continue

for col in columns:
    print(f'Best params for {col}: {best_params[col]}, MSE: {lowest_mse[col]}')
```

```
# Best params for JSW: (5,1,0), MSE: 11.843689481087306 --OK
# Best params for ALR: (0, 3, 1), MSE: 6.557770249093277 --OK
# Best params for CDR: (3, 4, 2), MSE: 10.408574870811051 --OK
# Best params for DNP: (4, 3, 4), MSE: 1677.697466214512
```

Błąd średniokwadratowy dla spółki DNP jest zbyt wysoki w stosunku do wariancji, więc spróbowałem użyć modelu *SARIMAX*, który pozwolił obniżyć *MSE*, ale nie sądzę, żeby ceny akcji DNP wykazywały się miesięczną sezonowością, więc myślę że nie warto brać ich prognoz pod uwagę podczas dalszej analizy.

```
#Sprawdzanie wariancję, żeby ocenić MSE
```

```
variance={}
```

```
for column in train_df.columns[1:]:
```

```
    variance[column]= train_df[column].var()
```

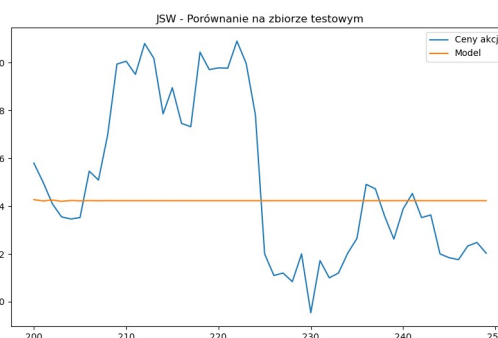
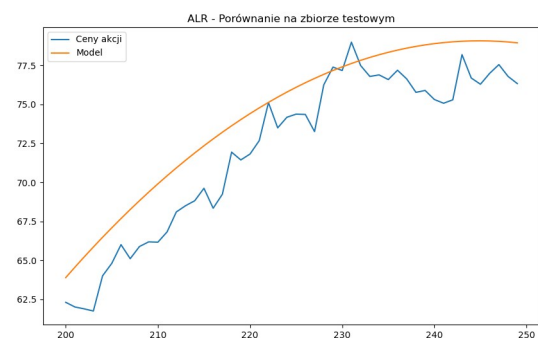
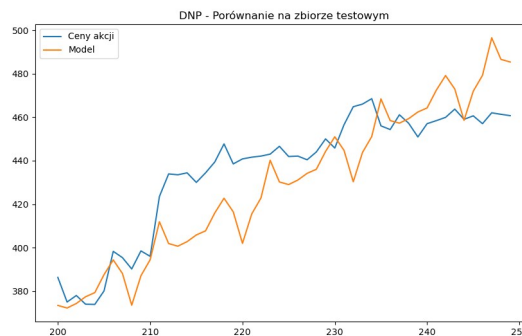
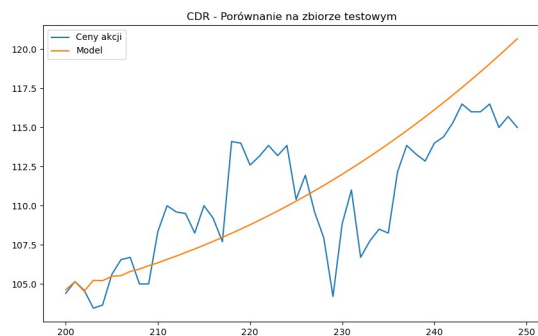
```
    print(f'Wariancja dla kolumny {column}: {variance[column]}, MSE: {min_mse[column]}, BEST PARAMS: {best_params[column]}')
```

```
# # Wariancja dla kolumny ALR: 51.0976731155779, MSE: 6.557770249093277, BEST PARAMS: (0, 3, 1)
```

```
# # Wariancja dla kolumny JSW: 56.38201414824117, MSE: 11.843689481087306, BEST PARAMS: (5, 1, 0)
```

```
# # Wariancja dla kolumny CDR: 341.07862969846747, MSE: 10.408574870811051, BEST PARAMS: (3, 4, 2)
```

```
# # Wariancja dla kolumny DNP: 1510.9146530150758, MSE: 329.1053547512303, BEST PARAMS: (1, 0, 0, 1, 0, 0, 12)
```



## 5. Prognoza

Na podstawie wygenerowanych modeli tworzę prognozy na okres 01.01.2024 – 31.01.2024 i porównuję je z rzeczywistymi cenami akcji.

```
#prognoza.py

import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.tsa.statespace.sarimax import SARIMAX

df = pd.read_csv("portfel.csv", index_col='Data', parse_dates=True)

# Definiowanie parametrów modelu
best_params = {
    'ALR': (0, 3, 1),
    'JSW': (5, 1, 0),
    'CDR': (3, 4, 2),
    'DNP': (1, 0, 0, 1, 0, 0, 12)
}

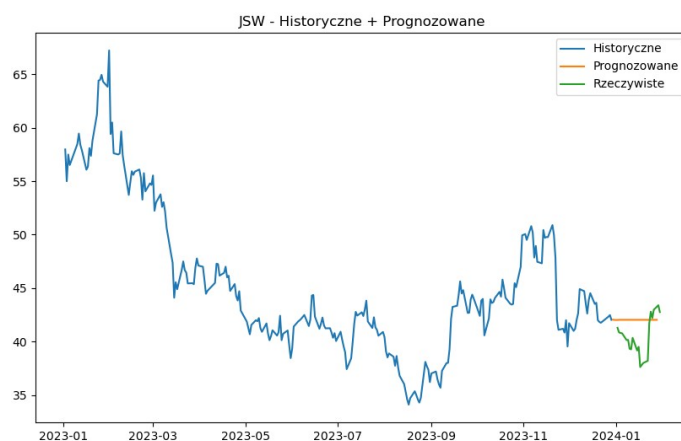
# Przygotowanie słownika do przechowywania prognoz
predictions = {}

from pandas.tseries.offsets import DateOffset

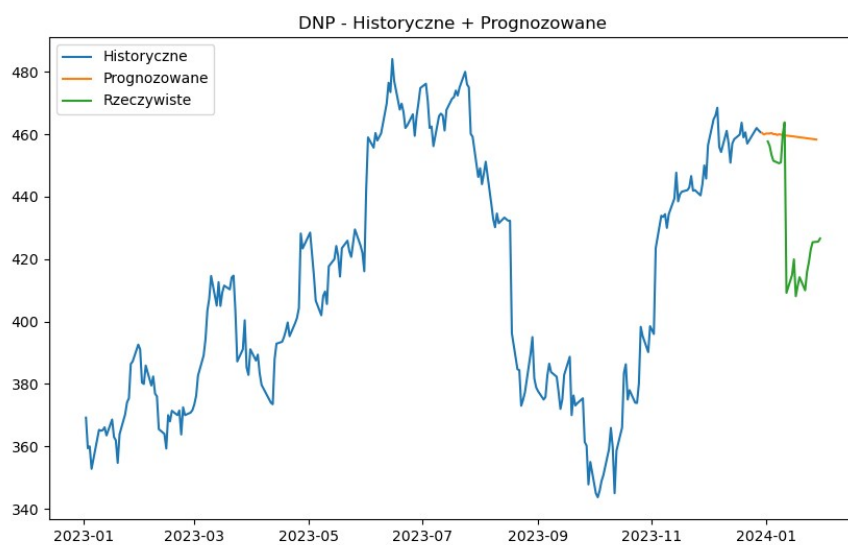
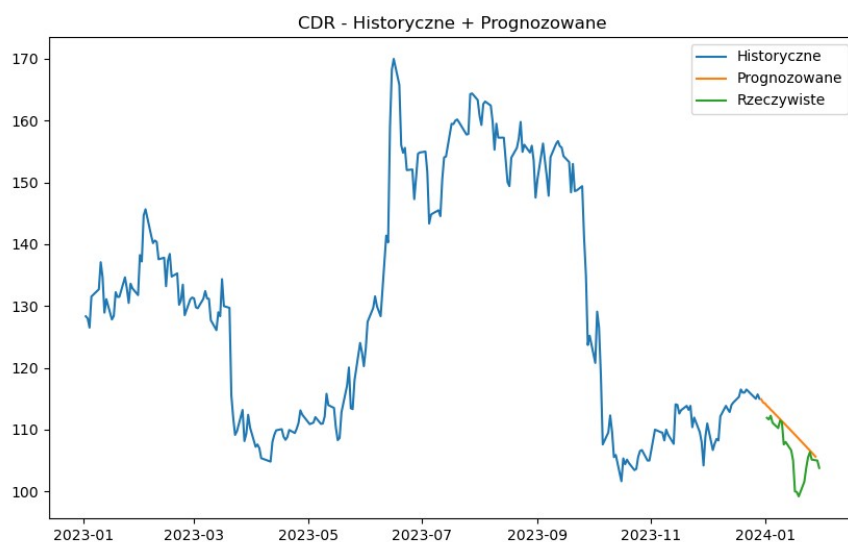
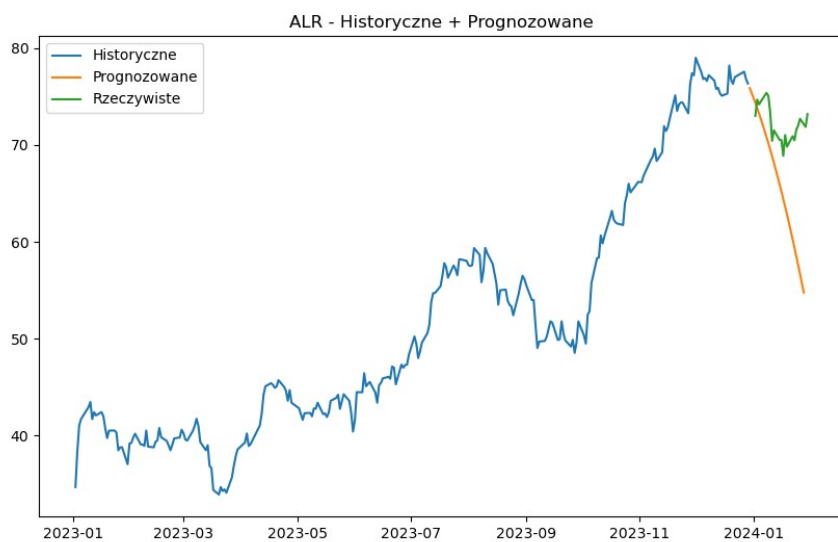
# Trenowanie modelu i generowanie prognozy dla każdej kolumny
for column in df.columns:
    if column == 'DNP':
        model = SARIMAX(df[column], order=(1, 0, 0), seasonal_order=(1, 0, 0, 12))
    else:
        model = ARIMA(df[column], order=best_params[column])
    model_fit = model.fit()
    forecast = model_fit.forecast(steps=30) # prognoza na następne 30 dni
    predictions[column] = forecast

future_data=pd.read_csv('cleaned_future_stocks.csv', index_col='Data', parse_dates=True)

# Wyświetlenie wykresów prognoz
for column in predictions.keys():
    plt.figure(figsize=(10,6))
    plt.plot(df[column], label='Historyczne')
    plt.plot(predictions[column], label='Prognozowane')
    plt.plot(future_data[column], label='Rzeczywiste')
    plt.title(f'{column} - Historyczne + Prognozowane')
    plt.legend()
    plt.show()
```







## 6. Interpretacja i wnioski

Na podstawie wykresów można zauważyć, że prognozy generowane przez modele ARIMA i SARIMAX nie odzwierciedlają w pełni dynamiki obserwowanej w danych empirycznych. W szczególności, prognozy nie uwzględniają znaczących fluktuacji cen, które są typowe dla rynków finansowych.

Ponadto, prognozy nie odzwierciedlają trendów obserwowanych w danych empirycznych. Na przykład w przypadku JSW, w danych empirycznych widzimy, że cena spada od 41.29 do 37.62, a następnie wzrasta do 43.40. Tymczasem prognozy utrzymują się na stałym poziomie około 42.03, nie odzwierciedlając tych trendów.

Jednym z możliwych wyjaśnień takiego zachowania modeli ARIMA i SARIMAX może być fakt, że są one modelami liniowymi i mogą mieć trudności z modelowaniem nieliniowych, długoterminowych zależności obserwowanych w danych finansowych.

W związku z powyższym, prognozy generowane przez te modele mogą nie być użyteczne do celów inwestycyjnych. Wysokie ryzyko błędu prognozy, jak wskazuje brak zgodności prognoz z obserwowanymi trendami i fluktuacjami, może prowadzić do błędnych decyzji inwestycyjnych i potencjalnych strat.

## Podsumowanie

Rynki finansowe są złożone i ma na nie wpływ wiele różnych czynników, takich jak zmiany w gospodarce, polityce, technologii, czy też nastroje inwestorów. Te czynniki są trudne do przewidzenia i modelowania, co sprawia, że prognozy oparte na samych danych historycznych są często niedokładne.

Dlatego, choć te modele mogą być użyteczne jako część strategii inwestycyjnej, zawsze powinny być używane z ostrożnością i w połączeniu z innymi narzędziami i informacjami. Zawsze istnieje ryzyko, że rzeczywiste wyniki będą różnić się od prognoz.