

**Міністерство освіти і науки України
Національний технічний університет України «КПІ» імені Ігоря Сікорського
Кафедра обчислювальної техніки ФІОТ**

**ЗВІТ
з лабораторної роботи №2
з навчальної дисципліни «Вступ до технології Data Science»**

Тема:

СТАТИСТИЧНЕ НАВЧАННЯ З ПОЛІНОМІАЛЬНОЮ РЕГРЕСІЄЮ

Виконав:

Студент 3 курсу кафедри ІПІ ФІОТ,
Навчальної групи ІП-11
Лошак В.І.

Перевірив:

Професор кафедри ОТ ФІОТ
Писарчук О.О.

Київ 2023

I. Мета:

Виявити дослідити та узагальнити особливості реалізації процесів статистичного навчання із застосуванням методів обробки Big Data масивів та калмановської рекурентної фільтрації з використанням можливостей мови програмування Python.

II. Завдання:

Реалізація проекту триває та спрямовано на збільшення функціональності програмної компоненти

Лабораторія провідної IT-компанії реалізує масштабний проект розробки універсальної платформи з обробки Big Data масиву статистичних даних поточного спостереження для виявлення закономірностей і прогнозування розвитку контрольованого процесу. Платформа передбачає розташування back-end компоненти на власному хмарному сервері з наданням повноважень користувачам заздалегідь адаптованого front-end функціоналу універсальної платформи.

III. Завдання IV рівня(максимум 15 балів):

Реалізувати групу вимог 1 та (або) 2 з імплементацією однієї з групи вимог 3.

Докладно опитати отримані R&D рішення

Група вимог_1:

1. Отримання вхідних даних із властивостями, заданими в Лр_1;
2. Модель вхідних даних із аномальними вимірами;
3. Очищення вхідних даних від аномальних вимірів. Спосіб виявлення аномалій та очищення обрати самостійно;
4. Визначення показників якості та оптимізація моделі (вибір моделі залежно від значення показника якості). Показник якості та спосіб оптимізації обрати самостійно.
5. Статистичне навчання поліноміальної моделі за методом найменших квадратів (МНК – LSM) – поліноміальна регресія для вхідних даних, отриманих в п.1,2. Спосіб реалізації МНК обрати самостійно;
6. Прогнозування (екстраполяцію) параметрів досліджуваного процесу за «навченою» у п.5 моделлю на 0,5 інтервалу спостереження (об'єму вибірки);
7. Провести аналіз отриманих результатів та верифікацію розробленого скрипта.

Група вимог 3:

3.1. Здійснити розробку власного алгоритму виявлення аномальних вимірів та / або «навчання» параметрів відомих алгоритмів «бачити» властивості статистичної вибірки.

IV. Результати виконання лабораторної роботи.

1. Отримання вхідних даних із властивостями, заданими в Лр_1

Для отримання даних попередньої лабораторної використано бібліотеку pandas. Загружені дані про тренд та залишки скомбіновано в один data frame. В лабораторній 1 ці дані вже було очищено від аномалій.

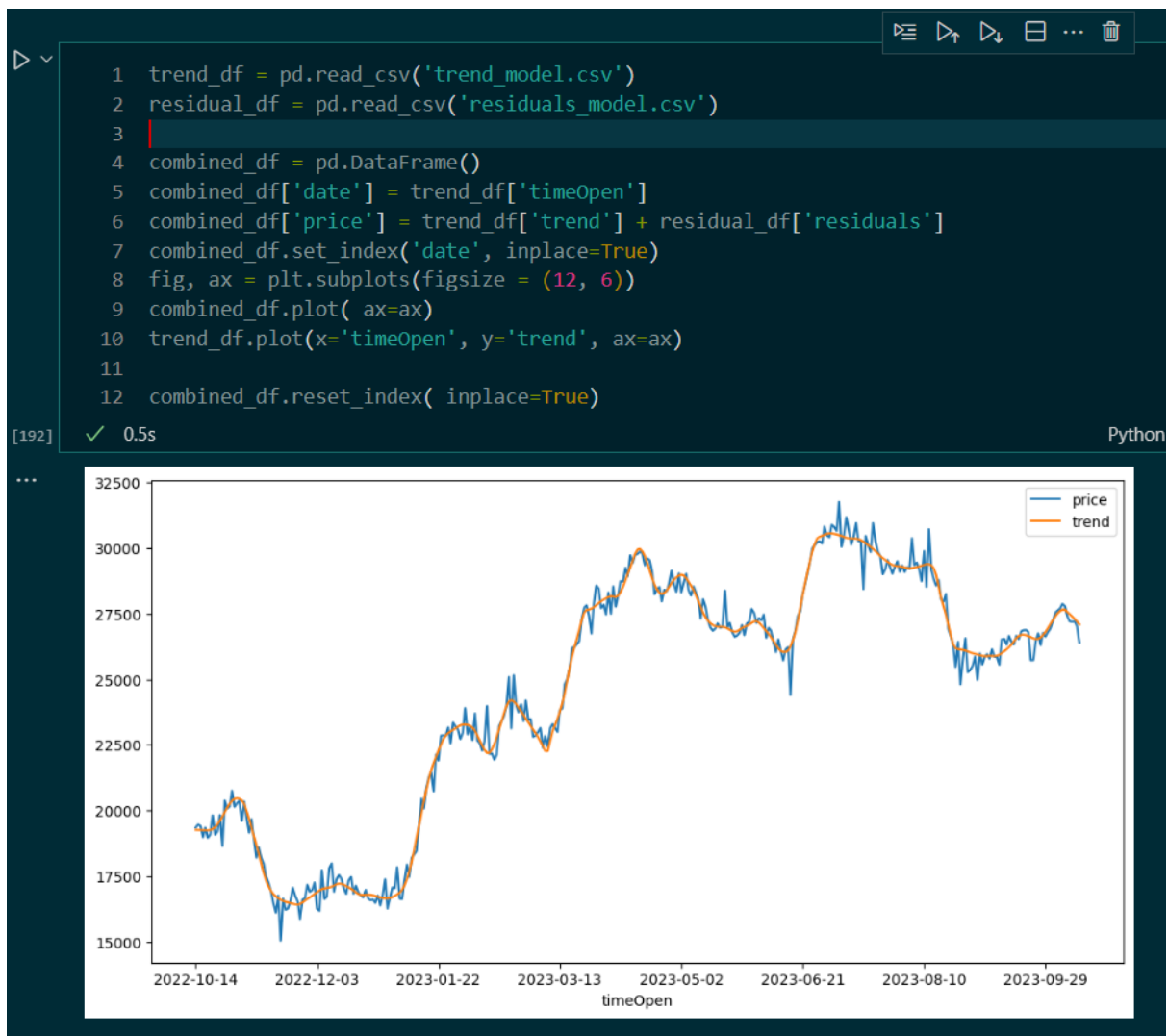


Рис. 1— Візуалізація тренду та комбінованої моделі даних.

2. Модель вхідних даних із аномальними вимірами;

Для того щоб додати аномалії до даних використано метод додавання кратної певному мультиплікатору кількості стандартних відхилень. Згенеровані аномалії розподілені рівномірно.

```
2 threshold = 5
3 anomaly_percentage = 0.1 # 10% of total data
4 price_std = np.std(residual_df['residuals'])
5
6 total_data = len(combined_df)
7 num_anomalies = int(total_data * anomaly_percentage)
8 anomaly_indices = np.random.choice(np.arange(total_data), size=num_anomalies, replace=False)
9
10 combined_df_with_anomalies = combined_df.copy()
11 for index in anomaly_indices:
12     sign = np.random.choice([-1, 1])
13     combined_df_with_anomalies.loc[index, 'price'] += sign * threshold * price_std
14
15 combined_df_with_anomalies.plot()
16 combined_df.plot()
```

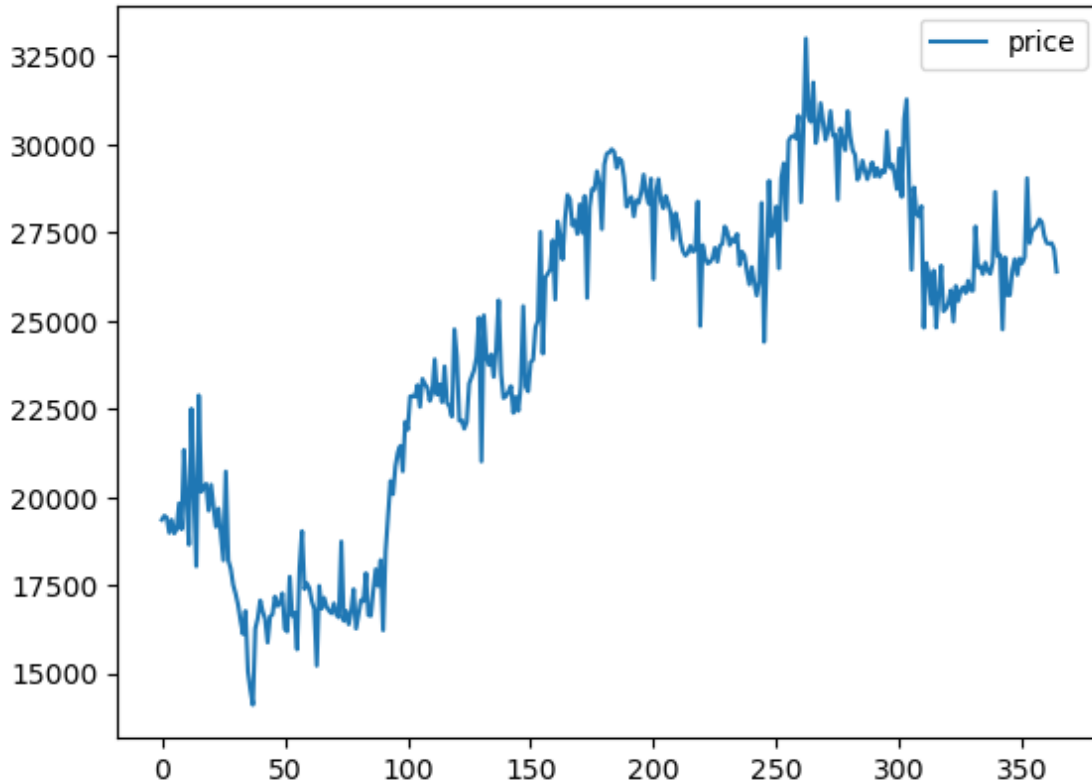


Рис. 2— Візуалізація часового ряду з доданими аномаліями що відхиляються на 5 стандартних відхилень.

3. Очищення вхідних даних від аномальних вимірів. Спосіб виявлення аномалій та очищення обрати самостійно;

3.1. Здійснити розробку власного алгоритму виявлення аномальних вимірів та / або «навчання» параметрів відомих алгоритмів «бачити» властивості статистичної вибірки.

Для розробки власного алгоритму для виявлення аномалій було вирішено модифікувати алгоритм `Sliding_Window_AV_Detect_medium` наведений Писарчук О.О. в файлі `L_1_3_Statistical_learnin`. В ході модифікації було створено алгоритм для підбору оптимальних параметрів для заданого часового ряду, а також модифіковано спосіб збільшення вікна. Оновлений алгоритм `medium` має вигляд:

Крок 1: Створити вікно даних з заданою параметром розмірністю

Крок 2: Визначити статистичні характеристики поточного вікна

Крок 3: Зсунути вікно вздовж інтервалу на одну поділку

Крок 4: Визначити статистичні характеристики нового вікна

Крок 5: Порівняти стат. Характеристики нового вікна з даними попереднього вікна і замінити включене до вікна значення на медіанне значення по поточному вікну.

Крок 6: встановити характеристики поточного вікна замість минулих характеристик

Крок 7: Якщо не досягнуто кінця даних — повернутися на крок 3.

```

1 def sliding_window_anomaly_detection(data, window_size, threshold):
2     result = data.copy()
3     # standard deviation of the initial window
4     previous_std = np.std(result[:window_size])
5
6     # slide the window across the data
7     for i in range(len(result) - window_size + 1):
8         window = result[i : i + window_size]
9         window_std = np.std(window)
10
11         if window_std > threshold * previous_std:
12             result[i + window_size - 1] = np.median(window)
13             previous_std = window_std
14
15     return result

```

✓ 0.0s

Рис. 3— Оновлений алгоритм детекції аномалій.

Для підбору параметрів моделі було розроблено алгоритм що проходиться через усі значення розміру вікна та трешхолду в певному діапазоні значень, та для кожної комбінації оцінює якість очистки використовуючи MSE.

```

1 def gridsearch(func_name ):
2     # Define the ranges for threshold and window_size
3     thresholds = np.linspace(0.1, 2.0, 20) # Change these values as needed
4     window_sizes = range(1, 15) # Change these values as needed
5
6     # Initialize variables to store the best parameters and the smallest error
7     best_threshold = None
8     best_window_size = None
9     smallest_error = np.inf
10
11     # Grid search over all combinations of threshold and window_size
12     for threshold in thresholds:
13         for window_size in window_sizes:
14             cleaned_data = sliding_window_anomaly_detection(combined_df_with_anomalies['price'].values, window_size, threshold)
15             error = mean_squared_error(combined_df['price'].values, cleaned_data)
16
17             if error < smallest_error:
18                 best_threshold, best_window_size, smallest_error = threshold, window_size, error
19                 # print(best_threshold, best_window_size, smallest_error)
20
21     return best_threshold, best_window_size, smallest_error

```

[197] ✓ 0.0s Python

Рис. 4— Алгоритм підбору параметрів для очищення вибірки від аномалій способом sliding_window_medium.

```

1 print('Best threshold: ', sw1_threshold)
2 print('Best window size: ', sw1_window_size)
3 print('Smallest error: ', sw1_smallest_error)

```

[10] ✓ 0.0s

... Best threshold: 1.8
Best window size: 5
Smallest error: 164207.3610226861

Рис. 5— Значення підібраних параметрів алгоритму sliding_window_medium

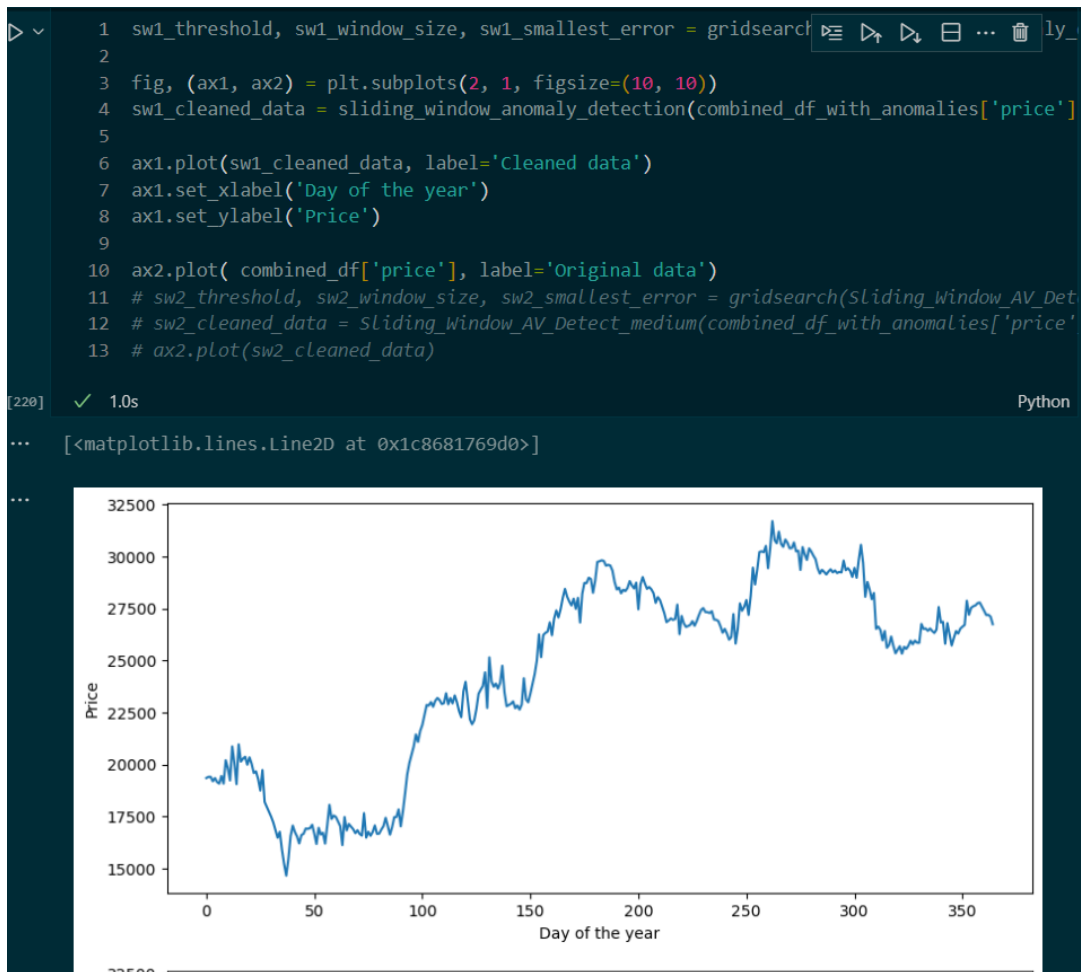


Рис. 6—Ряд очищений від аномалій з використанням оптимальних параметрів очищення.

4. Визначення показників якості та оптимізація моделі (вибір моделі залежно від значення показника якості). Показник якості та спосіб оптимізації обрати самостійно.

Для створення прогнозу було обрано побудувати поліноміальну модель. Для оцінки якості моделі вибрано показник R2. Для підбору оптимального значення параметру поліному використано алгоритм пошуку що використовує R2. Пошук параметра здійснено в діапазоні 1—10.

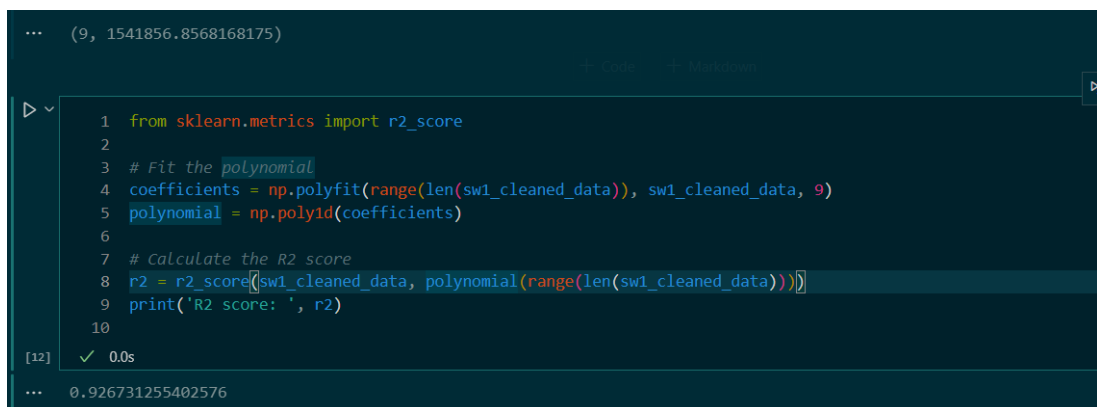


Рис. 7— Метрика R2 для оптимальної поліноміальної моделі.

5. Статистичне навчання поліноміальної моделі за методом найменших квадратів (МНК – LSM) – поліноміальна регресія для вхідних даних, отриманих в п.1,2. Спосіб реалізації МНК обрати самостійно;

```
1 degrees = range(1, 10) # Change the range of degrees as needed
2
3 best_degree = None
4 smallest_error = np.inf
5
6 for degree in degrees:
7     # Fit the polynomial
8     coefficients = np.polyfit(range(len(sw1_cleaned_data)), sw1_cleaned_data, degree)
9     polynomial = np.poly1d(coefficients)
10
11     # Calculate the mean squared error
12     error = mean_squared_error(sw1_cleaned_data, polynomial(range(len(sw1_cleaned_data))))
13
14     # Update the best degree and smallest error if necessary
15     if error < smallest_error:
16         best_degree = degree
17         smallest_error = error
18
19 # Plot the polynomial fit
20 plt.plot(sw1_cleaned_data)
21 plt.plot(polynomial(range(len(sw1_cleaned_data))), label=f'Degree {best_degree} Polynomial Fit')
22 plt.legend()
23 plt.show()
```

Рис. 8— Алгоритм підбору параметрів поліному.

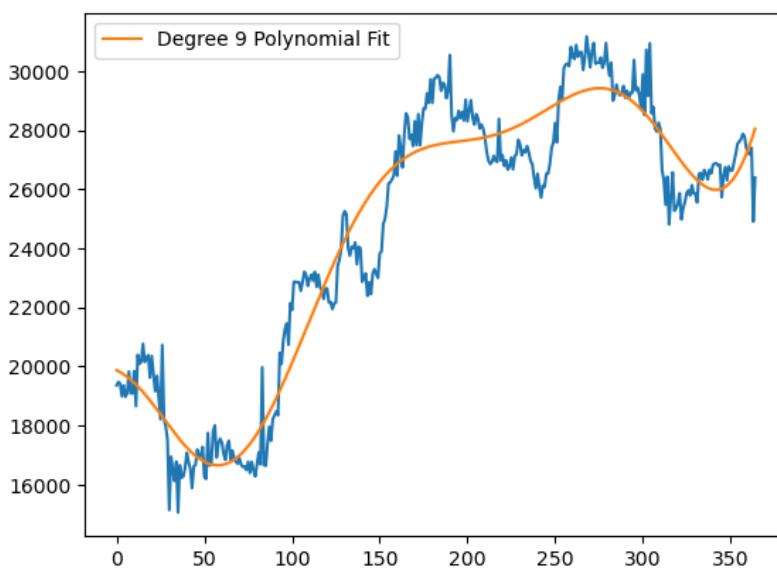


Рис. 9— Поліноміальна модель натренована за МНК після проведення підбору параметрів

6. Прогнозування (екстраполяцію) параметрів досліджуваного процесу за «навченою» у п.5 моделлю на 0,5 інтервалу спостереження (об'єму вибірки);

З вище наведених метрик зроблено висновок що модель достатньо добре описує дані одже ми можемо використати її для прогнозування.

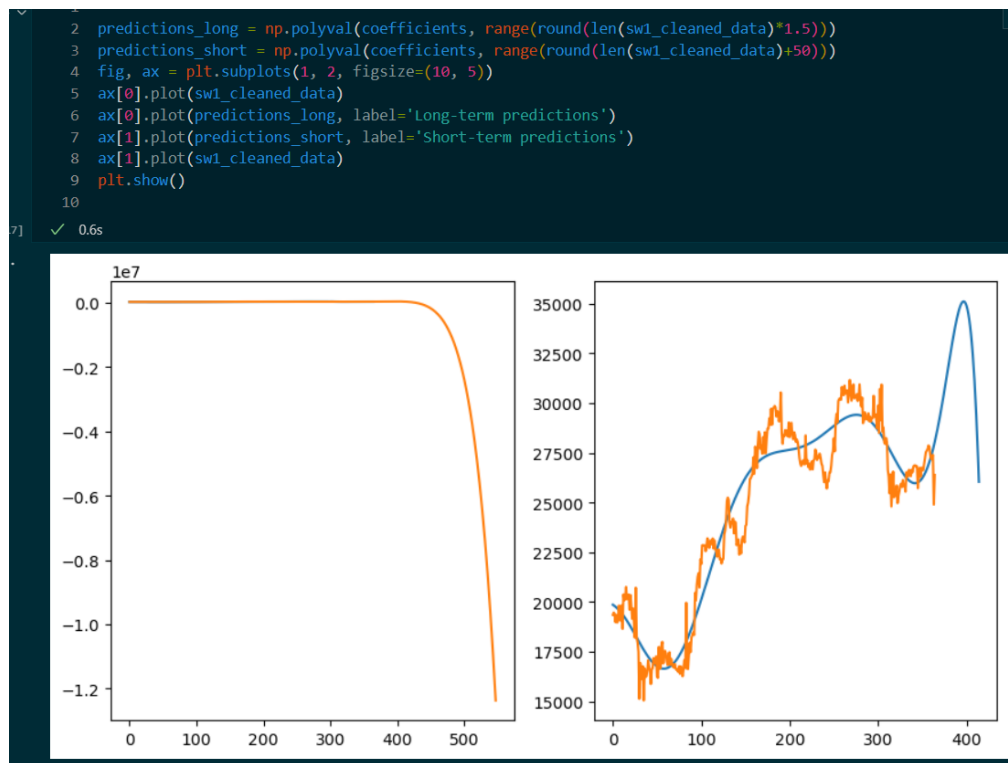


Рис. 10— Прогноз з використанням поліноміальної моделі на 1.5 та +10 відповідно до розміру вибірки.

7. Провести аналіз отриманих результатів та верифікацію розробленого скрипта.

З графіків видно що в короткостроковій перспективі модель прогнозує різкий стрибок вгору ціни Bitcoin з подальшим його стрімким падінням, що приводить нас до висновку що така волатильність може бути використана для спекуляцій на ринку. В той же час різке падіння відображає факт того що поліноміальна модель містить в собі певні надоліки в плані прогнозування оскільки не враховує історичні дані про діапазон в якому ціни біткоїн зазвичай перебувають що і приводить до появи від'ємних значень ціни і безпрецедентного прогнозу падіння. Незважаючи на вище сказане, в короткостроковій перспективі модель показує себе задовільно.

IV. Висновки.

У ході цієї роботи, я практикував побудову поліноміальної моделі для прогнозування часових рядів. Мною було створено/модифіковано алгоритм sliding_window_medium очищення даних від аномальних вимірів що адаптується до наданих даних за своїми параметрами і якісно відтворює оригінальні дані.

Виконав: студент ФІОТ Лошак В.І. ІІ-11