

Міністерство освіти та науки України
Національний університет «Львівська політехніка»



Звіт до лабораторної роботи
“Аналітичні та нереляційні бази даних”

Виконав:
Студент ІР-32
Лис Ярослав

Прийняв:
Верес З. Є.

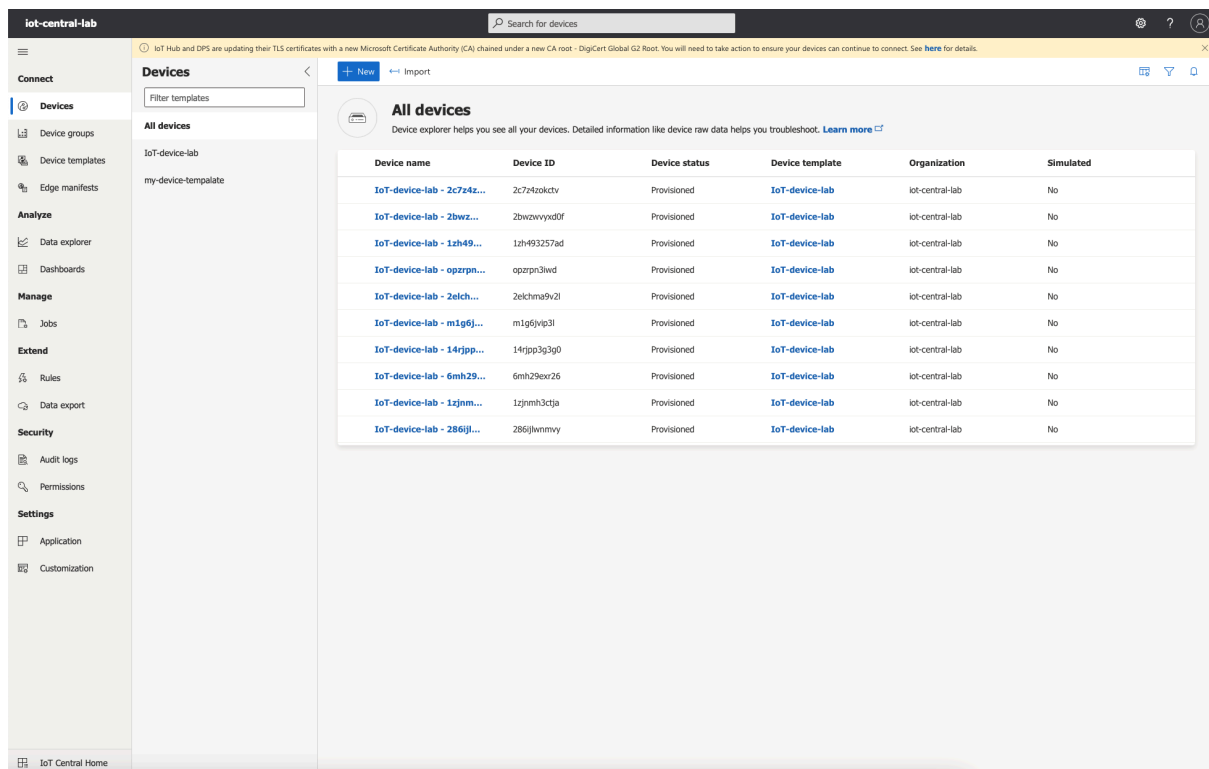
Львів – 2023

Завдання

1. Реєстрація 10 аплікацій у Azure IoT Core
2. Передача даних по MQTT протоколу у Azure IoT Core та подальше їх збереження у Apache Kafka (дані будемо брати з відкритих джерел та з Kaggle, тобто у вас буде 10 клієнтських апок, які відрізняються лише даними)
3. Реалізація data enrichment засобами Spark та збереження даних у Delta Lake
4. Налаштування AML воркспейсу
5. Реєстрація дата лейку як джерела даних у AML
6. Налаштування автомл на даних, (задача класифікації)

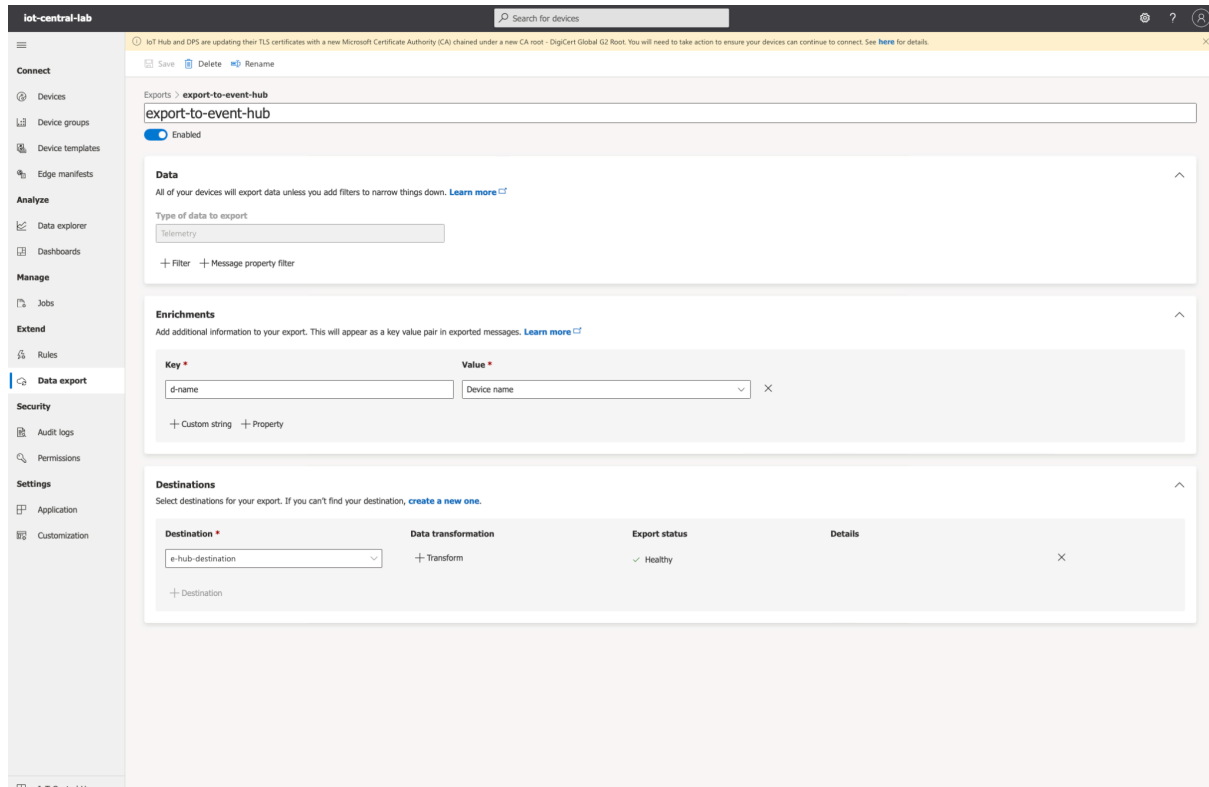
Хід роботи

1. Створюємо 10 скриптів на пайтоні. Створюємо 10 сенсорів у IoT Central та вписуємо connection string у скрипти. Результат виконання кроку:

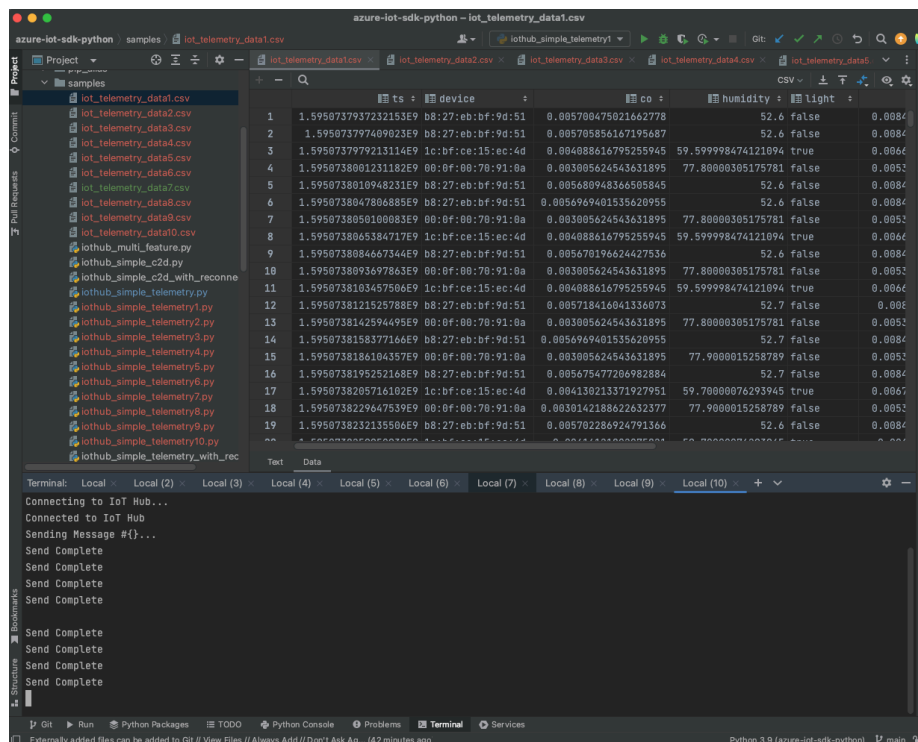


2. Налаштовуємо експортування даних до event hub та створюємо сам event hub.

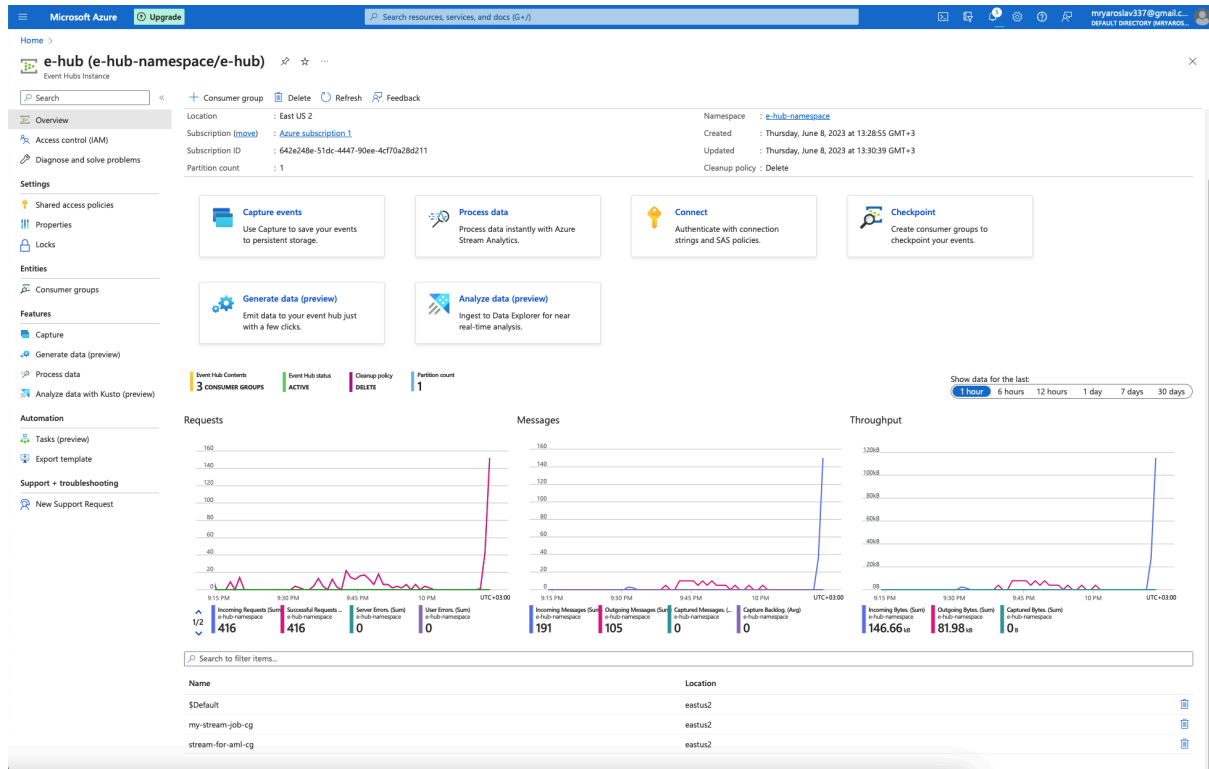
Data Export до event hub:



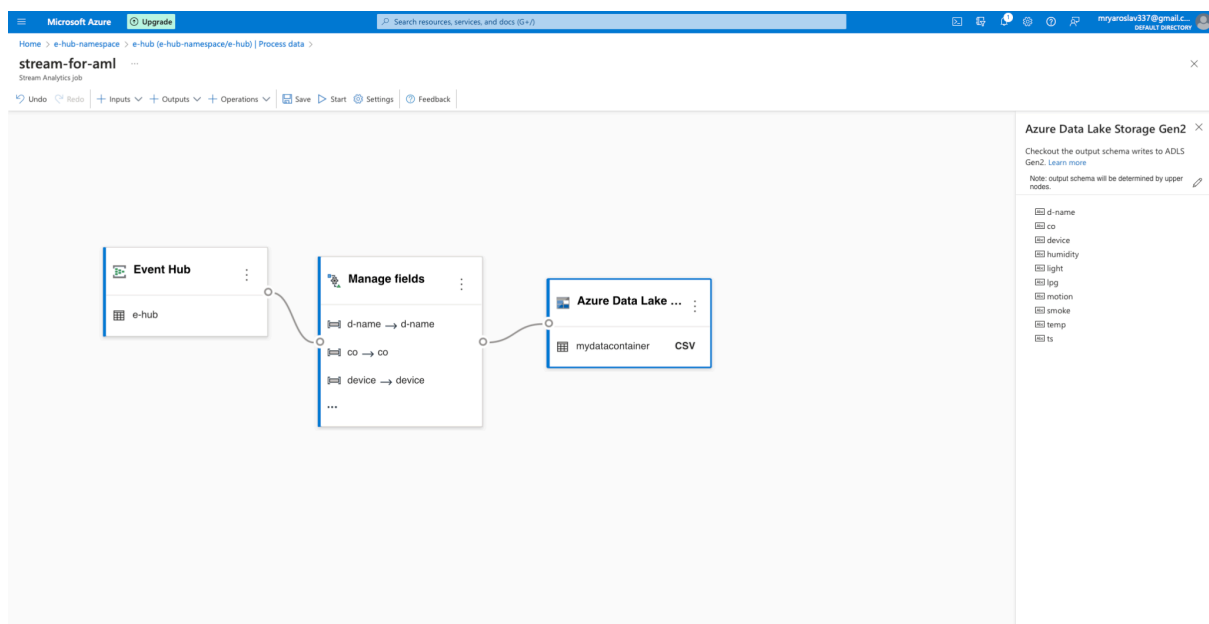
3. Запускаємо скрипти для відправки даних на event hub.



На цьому скріні можна побачити, що дані успішно передалися з IoT Central до event hub:



4. Налаштування stream job analytics для витягнення даних з event hub до azure data lake.



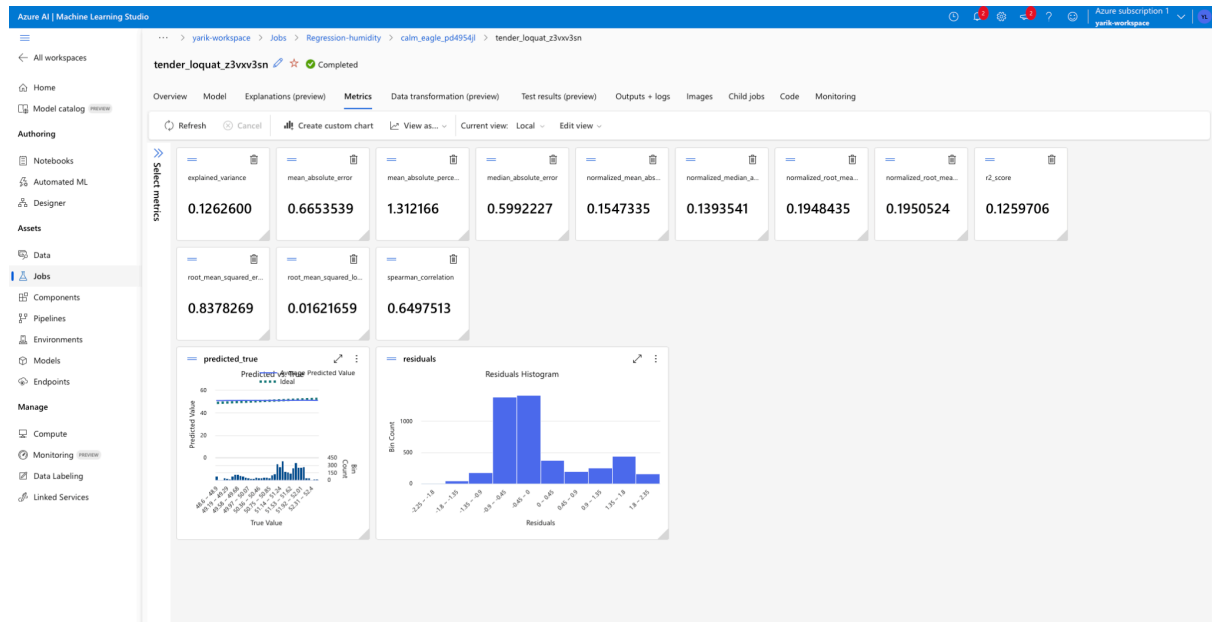
5. Налаштування вокрспейсу та старт тренування моделі:

The screenshot shows the Azure ML Studio interface. The left sidebar contains navigation links: All workspaces, Home, Model catalog, Authoring, Notebooks, Automated ML, Designer, Assets, Data, Jobs, Components, Pipelines, Environments, Models, Endpoints, Manage, Compute, Monitoring, Data Labeling, and Linked Services. The main panel displays the 'salmon_card_rx83h8jp' job, which is in a 'Running' state. The 'Overview' tab is active, showing job properties: Status (Running), Model training, Created on (Jun 8, 2023 10:30 PM), Start time (Jun 8, 2023 10:30 PM), Compute target (mryaroslav3371), Name (AutoML_91687098-58ab-425a-922c-0ab1c550858f), and Script name (--). The 'Inputs' section shows 'Input name: training_data' and 'Dataset: for-aml:1'. The 'Best model summary' section indicates 'No data'. The 'Run summary' section shows 'Task type: Regression' and 'Primary metric: Normalized root mean squared error'. The 'Tags' section shows 'No tags'.

6. Успішне закінчення тренування моделі

The screenshot shows the Azure ML Studio interface with the 'calm_eagle_pd4954jl' job, which is now in a 'Completed' state. The 'Overview' tab is active, showing job properties: Status (Completed), Script name (--), Created by (Yaroslav Lys), Job type (Automated ML), Experiment (Regression-humidity), Arguments (None), See all properties, Raw JSON, and See YAML job definition. The 'Inputs' section shows 'Input name: training_data' and 'Dataset: for-aml-one:1'. The 'Outputs' section shows 'Output name: best_model', 'Model: azureml_AutoML_74586d47-07f2-4413-a08a-fef81887e3d2_0_output_mfllow_log_model_1908001935:1', 'Output name: full_training_dataset', and 'Dataset: 50fb1344-6bc4-40de-973b-7603fe8e2b8'. The 'Best model summary' section shows 'Algorithm name: MaxAbsScaler, LightGBM', 'Hyperparameters: View hyperparameters', 'Normalized root mean squared error: 0.05596', and 'Sampling: 100.00 %'. The 'Run summary' section shows 'Task type: Regression' and 'Featureization'. The 'Tags' section shows a long list of tags, including 'dynamic_allowlisting_iterations' and 'fit_time_000'.

Результати тренування моделі:



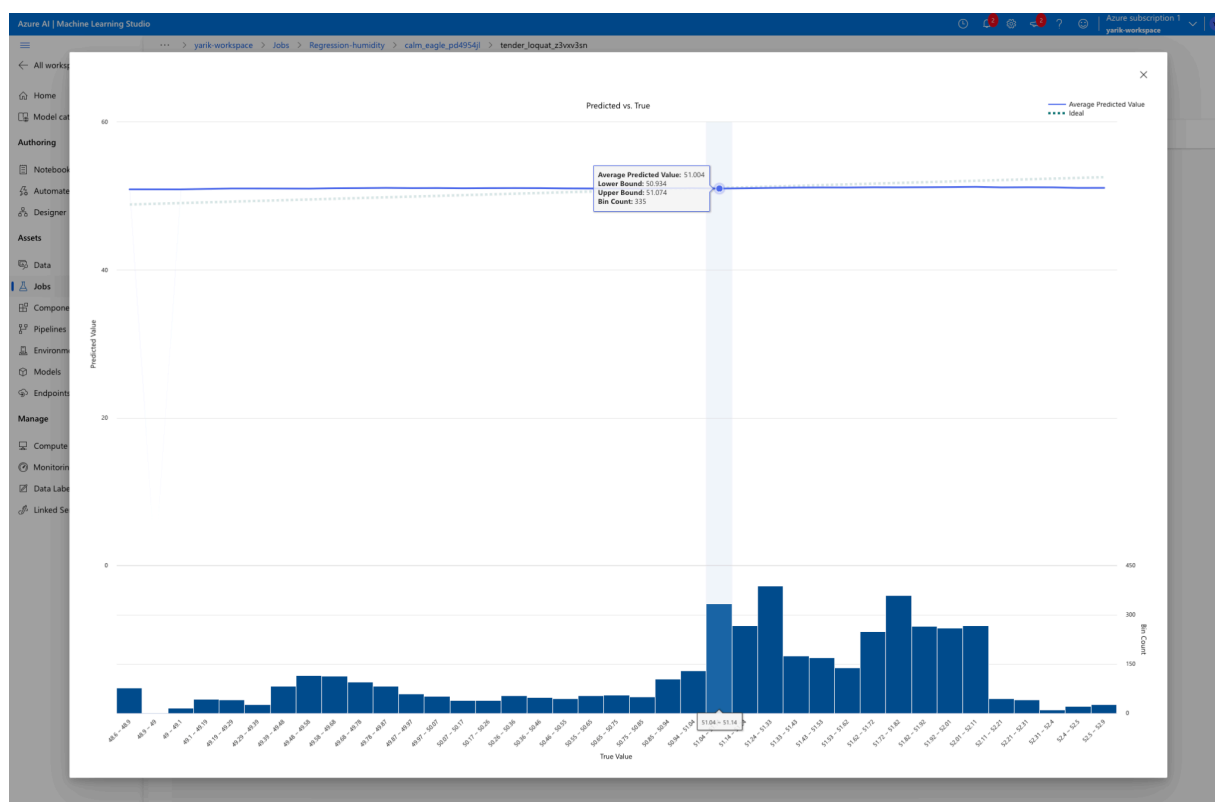
Пояснення метрик:

Метрика	Опис
explained_variance	Пояснена дисперсія вимірює, якою мірою модель враховує варіацію цільової змінної. Це відсоткове зменшення дисперсії вихідних даних до дисперсії помилок. Коли середнє значення помилок дорівнює 0, воно дорівнює коефіцієнту детермінації (див. r2_score нижче).
mean_absolute_error	Середня абсолютна похибка – це очікуване значення абсолютної різниці між цільовим показником і прогнозом.
mean_absolute_percentage_error	Середня абсолютна відсоткова похибка (MAPE) — це міра середньої різниці між

	прогнозованим і фактичним значенням.
median_absolute_error	Середня абсолютна похибка – це медіана всіх абсолютних різниць між цільовим показником і прогнозом. Ця втрата стійка до викидів.
r2_score	R2 (the coefficient of determination) measures the proportional reduction in mean squared error (MSE) relative to the total variance of the observed data.
root_mean_squared_error	Середньоквадратична помилка (RMSE) – це квадратний корінь із очікуваної квадратичної різниці між цільовим показником і прогнозом. Для неупередженого оцінювача RMSE дорівнює стандартному відхиленню.
root_mean_squared_log_error	Середньоквадратична логарифмічна помилка – це квадратний корінь із очікуваної квадратичної логарифмічної помилки.
spearman_correlation	Кореляція Спірмена є непараметричною мірою монотонності зв'язку між двома наборами даних. На відміну від кореляції Пірсона, кореляція Спірмена не передбачає, що обидва набори даних розподілені нормально. Як і інші коефіцієнти кореляції, Спірмена коливається від -1 до 1, де 0 означає відсутність кореляції. Кореляції -1 або 1 означають точне монотонне співвідношення.

Діаграма Predicted vs True

Для регресії та експерименту з прогнозуванням прогнозована та істинна діаграма відображає зв'язок між цільовою ознакою (справжні/фактичні значення) та прогнозами моделі. Істинні значення розміщуються вздовж осі x, і для кожного біну середнє прогнозоване значення наноситься на графік зі смугами помилок. Це дає змогу побачити, чи є модель упередженою щодо передбачення певних значень. Рядок відображає середній прогноз, а заштрихована область вказує на дисперсію прогнозів навколо цього середнього значення.

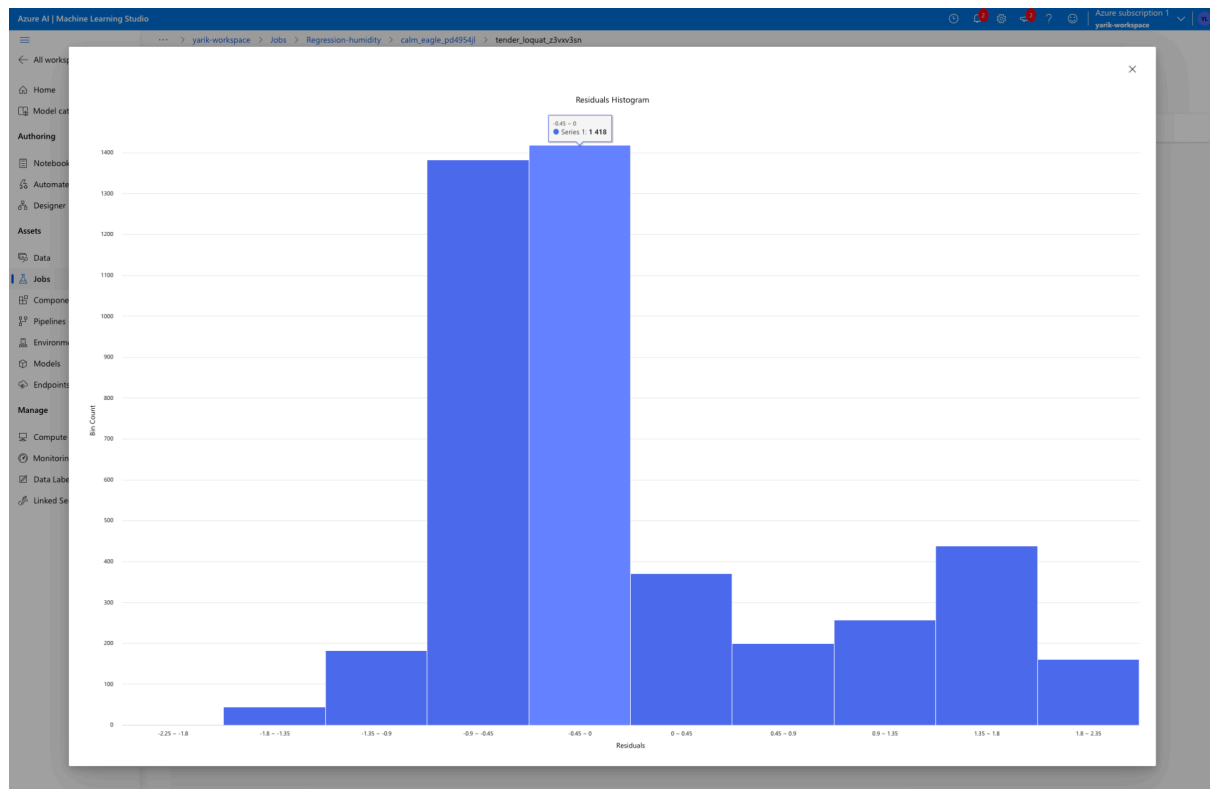


2

Діаграма залишків

Діаграма залишків — це гістограма помилок передбачення (залишків), згенерованих для експериментів регресії та прогнозування.

Залишки обчислюються як $y_{\text{predicted}} - y_{\text{true}}$ для всіх зразків, а потім відображаються як гістограма, щоб показати зміщення моделі.



Висновок: У цій лабораторній роботі я використовував сервіс IoT Central, щоб підключити пристрої до хмари та збирати дані з них. IoT Central надає простий спосіб налаштування та керування пристроями, а також забезпечує безпечний обмін даними з сервісами Azure. Event Hub дозволив мені збирати та обробляти великі обсяги подій з пристроїв IoT Central. Також я налаштував хмарне сховище для зберігання подій та надсилав їх до моделі машинного навчання для подальшого аналізу. Data Lake був використаний для зберігання, управління та аналізу великих обсягів даних.