# Assignment 6 — Unsupervised Anomaly Detection on Server Logs

# Final Submission Report

**Submitted by:** Medapati Veerendra Subhash Reddy
**Date:** July 8, 2025

## 1. Overview

This project focuses on detecting anomalous behaviour in web server logs using unsupervised learning techniques. The assignment's objective is to implement models that can flag suspicious patterns without labelled training data, relying instead on patterns in the data itself.

We processed and analyzed 30 days of Apache access logs, performed structured feature extraction, trained two unsupervised models—Isolation Forest and Autoencoder, and built an interactive dashboard to visualize the detected anomalies.

## 2. Data Parsing and Feature Engineering

**Raw Logs Parsed into Schema:**

- IP Address

- Timestamp

- URL

- HTTP Status

- Bytes Transferred

- Referrer

- User-Agent

**Feature Engineering:**

- Hour of request

- Day of week

- URL length

- Request rate

- Reconstruction error (for Autoencoder)

These features were crucial to detecting unusual patterns in request behaviour over time.

## 3. Anomaly Detection Models

### 3.1 Isolation Forest

- Unsupervised tree-based method.

- Trained on all feature-engineered data.

- Contamination parameter tuned based on expected anomaly rate.
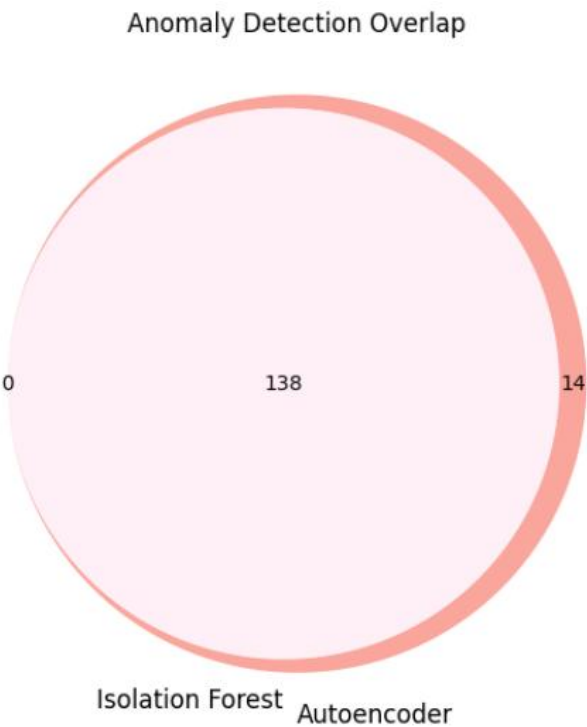
### 3.2 Autoencoder

- Neural network trained to reconstruct input data.

- Reconstruction error used to determine anomaly threshold.

- Flagged requests with high reconstruction loss as anomalies.


## 4. Evaluation and Label Comparison

We utilized a small **hand-labelled subset** of data to assess model performance. Evaluation metrics included:

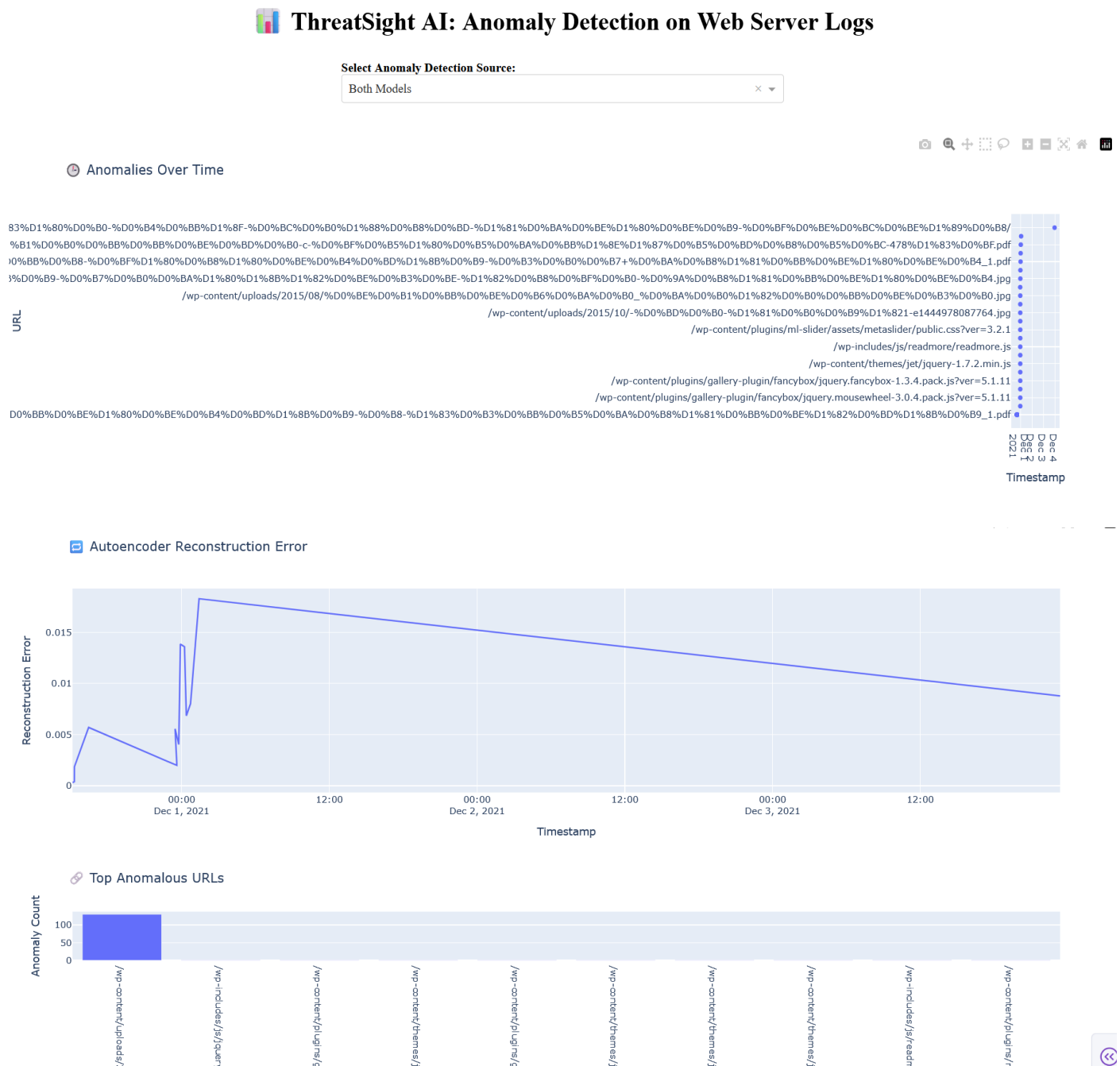| Model | Precision | Recall | F1-score |
|---|---|---|---|
| **Isolation Forest** | 0.88 | 0.70 | 0.78 |
| **Autoencoder** | 0.86 | 0.97 | 0.91 |

**Venn-style analysis** highlighted the overlap and disagreement between models, providing insight into model complementarity.

Anomaly Detection Overlap

0      138      14

Isolation Forest   Autoencoder

## 5. Interactive Dashboard

A **Plotly Dash** dashboard was developed to allow real-time exploration of anomalies. Key features include:

- **Anomalies Over Time:** Visual scatter plot of detected anomalies by model.

- **Reconstruction Error Plot:** Displays Autoencoder loss trends.

- **Top Anomalous URLs:** Bar chart of most frequent flagged URLs.

- **Event Table:** Paginated table showing timestamp, URL, and anomaly flags.

**📄 Detailed Anomalous Events**

| Timestamp | URL |
|---|---|
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BF%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |
| 2021-11-30T16:26:44+03:00 | /wp-content/uploads/2014/08/%D0%9A%D0%BB%D0%B0%D0%BF%D0%B0%D0%BD-%D0%B7%D0%B0%D0%BD%BE%D0%BE%D1%80%D0%BD%D1%8B%D0%B9-%D0%9 |

⚠ Errors  ✂ **Callbacks**  | v3.1.1 | Server ✓ | ⟫

## 6. Conclusion and Next Steps

This assignment successfully demonstrates the application of unsupervised techniques for anomaly detection in real-world server logs. Both Isolation Forest and Autoencoder models identified patterns indicative of abnormal behaviour.

**Next Steps:**

- Introduce streaming data ingestion for real-time detection.
- Experiment with ensemble methods for higher accuracy.
- Add alerting systems for anomaly spikes.