# Anomaly Detection at Pension Architects

## Summary

**Bachelor's degree in Computer Science field AI**

**Aynur Guliyeva**

Academic year 2025-2026

Campus Geel, Kleinhoefstraat 4, BE-2440 Geel

THOMAS MORE

LID VAN ASSOCIATIE KU LEUVEN

# PoC Anomaly Detection in Pension Data

This internship was carried out at **Pension Architects**, a Belgian fintech company specializing in pension administration systems and actuarial calculations for occupational pension plans. The internship took place in a professional environment where large volumes of financial and employment-related data are processed under strict regulatory and operational constraints. The project was embedded within the company's broader objective to improve data quality assurance and risk detection in pension plan administration.

## Context and Problem Statement

Pension contribution and reserve calculations are inherently complex. They depend on heterogeneous data sources, plan-specific rules, temporal patterns, and employer-reported information that may contain inconsistencies or errors. Traditional validation approaches rely heavily on predefined business rules and manual controls. While these methods are interpretable, they are difficult to maintain, scale poorly across plans, and are vulnerable to edge cases and reporting irregularities.

As data volumes continue to grow, there is a clear need for **automated anomaly detection mechanisms** that can assist experts by flagging potentially incorrect records, while still allowing human judgment to remain decisive. The central challenge addressed during this internship was therefore how to design an AI based system that can detect meaningful anomalies in pension data without embedding fragile business logic, and that can continuously improve through expert feedback.

## Objective of the Internship

The main objective of the internship was to design and implement a **scalable, plan anomaly detection framework** for pension contribution and reserve data. The system needed to:

- Operate on large, heterogeneous datasets.
- Avoid reliance on hard-coded actuarial formulas.
- Control false positives to limit unnecessary manual review.
- Integrate human expertise through a structured feedback loop.
- Be deployable in a production-ready cloud environment.

## Approach

The project followed an **iterative and exploratory approach**, evolving over time as insights were gained from different pension funds (PF1, PF2, PF3, and PF4). Initially, the work focused on **unsupervised anomaly detection**, using models such as Isolation Forests, density-based clustering techniques, and Autoencoders. These methods were well suited to situations where no labelled anomalies were available and helped uncover extreme outliers and structural inconsistencies. However, purely unsupervised approaches showed limitations, particularly in distinguishing true anomalies from rare but valid cases.

To overcome this, the project transitioned towards a **data-driven and eventually supervised approach**. Realistic anomalies were synthetically injected into cleaned datasets to create controlled training labels. This enabled the use of supervised models while preserving domain realism. In parallel, feature engineering focused exclusively on **generic, plan-independent characteristics**, such as ratios, temporal patterns, volatility measures, and distributional properties, rather than explicit business rules.

A key methodological innovation was the integration of a **human-in-the-loop (HITL)** mechanism. Model outputs get reviewed by domain experts, who label flagged cases as true anomalies or false positives. These labels then merged back into the training data, overriding synthetic labels where applicable. This allows the system to learn directly from expert judgment while retaining exposure to genuinely anomalous patterns.

## Technical Realisation

The final solution was implemented as a set of **AWS SageMaker pipelines**, covering data preprocessing, model training, batch scoring, and retraining. Large input files stored in Amazon S3 were processed using SageMaker Processing jobs, and models were trained and evaluated using SageMaker Training jobs. Batch Transform was used for scalable, non-real-time scoring.

## Results and Outcomes

The project resulted in a **fully automated, cloud-based anomaly detection framework** capable of processing large pension datasets and producing interpretable anomaly scores. Key outcomes include:

- A reusable pipeline architecture applicable across multiple pension plans.
- Improved precision through supervised learning and threshold optimisation.
- A functioning human-in-the-loop retraining mechanism.
- Reduced dependence on brittle, rule-based validation logic.
- A clear separation between detection, review, and retraining stages.

## Learning Outcomes and Relevance

From an academic and professional perspective, the internship provided hands-on experience at the intersection of **information management, machine learning, and organisational decision-making**. Key learning outcomes include:

- Applying machine learning techniques in a real-world, regulated domain.
- Understanding the trade-offs between automation and human oversight.
- Designing scalable data pipelines in a cloud environment.
- Translating abstract ML concepts into operational systems.

The internship aligns closely with the **Information Management** programme by combining data engineering, analytics, governance considerations, and human-centred system design. It demonstrates how AI systems can support, rather than replace, expert decision-making in complex organisational contexts.