# **AI-Powered Data Profiling & Compliance Assistant**

- Sharanprasath S: CT (HLT-Originations)

### 1. Introduction

### 1.1 Overview

The AI-Powered Data Profiling & Compliance Assistant automates the extraction of validation rules from regulatory documents, validates datasets against these rules, detects anomalies using machine learning, and provides an interactive compliance assistant. This system is designed to streamline regulatory compliance and data accuracy in financial and enterprise datasets.

### 1.2 Problem Statement

Regulatory compliance in financial datasets is complex, requiring manual effort to extract rules, validate data, and detect inconsistencies. This leads to inefficiencies and increased risk of errors. Our system automates this process using AI and ML to improve accuracy, efficiency, and compliance.

## 1.3 Key Features

- Automated Rule Extraction: Extracts validation rules from regulatory PDFs using Generative
- **Dataset Validation**: Applies extracted rules to detect inconsistencies in uploaded datasets.
- Anomaly Detection: Identifies outliers using Isolation Forest and other ML techniques.
- Al-Powered Compliance Assistant: Provides explanations for violations and suggests rule refinements.

### 2. Technical Approach

## 2.1 Rule Extraction (Al-Powered Regulatory Rule Extraction)

## **Process:**

- Upload regulatory PDF and sample dataset (CSV).
- Text & Table Extraction: Extracts structured and unstructured data from the document.
- Chunking & Vector Storage (RAG Approach): Uses text chunking and ChromaDB for storage.
- Rule Extraction using Google Gemini AI: Al processes extracted content and generates structured validation rules.
- **JSON Rule Output:** Stores extracted rules in JSON format for dataset validation.

#### **Detailed Technical Flow:**

### 1. PDF Processing:

- Uses pdfplumber or PyPDFLoader to extract both text and table data separately.
- Pre-processing: Removes unnecessary content, identifies structured sections, and cleans extracted text.

## 2. Chunking & Vector Storage (RAG Approach):

- Splits text into semantic chunks using RecursiveCharacterTextSplitter.
- Uses GoogleGenerativeAlEmbeddings to generate vector embeddings.
- o Stores indexed chunks in **ChromaDB** for efficient retrieval.

### 3. Al-Driven Rule Extraction:

- o Sends extracted chunks to Google Gemini API.
- o Al generates **validation rules** based on dataset schema and extracted content.
- o Rules include column constraints, expected formats, data types, and business rules.

## 4. JSON Output & Storage:

- o Outputs structured rules in **JSON format**.
- Users can download extracted rules for validation purposes.

## 2.2 Dataset Validation (Automated Compliance Checking)

#### **Process:**

- Upload dataset (CSV) and validation rules (JSON).
- Batch Processing for Efficiency: Splits dataset into smaller batches for validation.
- Al-Powered Validation using Google Gemini: Checks data integrity, constraints, and expected formats.
- Validation Report Generation: Generates structured JSON & CSV reports highlighting violations.

# **Detailed Technical Flow:**

### 1. Dataset & Rule Ingestion:

- Reads dataset using Pandas.
- Loads extracted validation rules from JSON.

## 2. Batch Processing for Efficiency:

- Reduces API calls by splitting dataset into smaller batches.
- Sends batches to Google Gemini AI for rule-based validation.

# 3. Al-Powered Validation:

- o Al checks each column for format, constraints, expected values, and missing data.
- o Identifies violations and suggests fixes.
- o Ensures **logical consistency** (e.g., no negative loan amounts, proper date formats).

### 4. Real-time UI & Reporting:

- o **Progress bar updates dynamically** per batch.
- o Al-generated violations and suggestions are displayed in a structured format.
- Users can download the validation report (CSV) for further analysis.

## 2.3 Compliance Chat Assistant (Interactive AI-Powered Queries)

### **Process:**

- User Query Processing: Users ask about violations or request rule refinements.
- Retrieval of Relevant Context: Uses FAISS vector database for retrieving violations.
- Al Response Generation using Google Gemini API: Provides structured explanations and rule suggestions.
- Interactive Chat UI: Allows users to iteratively refine compliance queries.

#### **Detailed Technical Flow:**

## 1. User Query Processing:

- Accepts natural language queries from users.
- Determines intent (e.g., explanation request, rule refinement, compliance guidance).

## 2. Retrieval-Augmented Search:

- Stores violations in FAISS (vector database).
- o Retrieves most relevant validation errors for the query.

# 3. Al Response Generation:

- Al processes retrieved violations and generates structured explanations.
- o If rule refinement is requested, AI suggests modifications to the extracted rules.

## 4. Interactive Chat UI:

- Displays user messages in chat bubbles.
- Dynamically updates responses, allowing iterative compliance improvements.

# 2.4 Anomaly Detection (Machine Learning-Based Outlier Detection)

### **Process:**

- Feature Selection & Preprocessing: Selects numeric columns and standardizes features.
- Model-Based Anomaly Detection: Uses Isolation Forest for anomaly detection.
- Auto-Calculated Contamination Rate: All estimates the proportion of anomalies dynamically.
- Interactive UI & Visualization: Scatter plots highlight anomalies; users can download flagged data.

#### **Detailed Technical Flow:**

# 1. Feature Selection & Preprocessing:

- Selects **only numeric columns** for anomaly detection.
- o Standardizes features using **StandardScaler**.

# 2. Model-Based Anomaly Detection:

- o Uses **Isolation Forest** for anomaly detection (default model).
- Allows switching to:
  - Local Outlier Factor (LOF) → Density-based anomaly detection.
  - **DBSCAN** → Clustering-based anomaly detection.

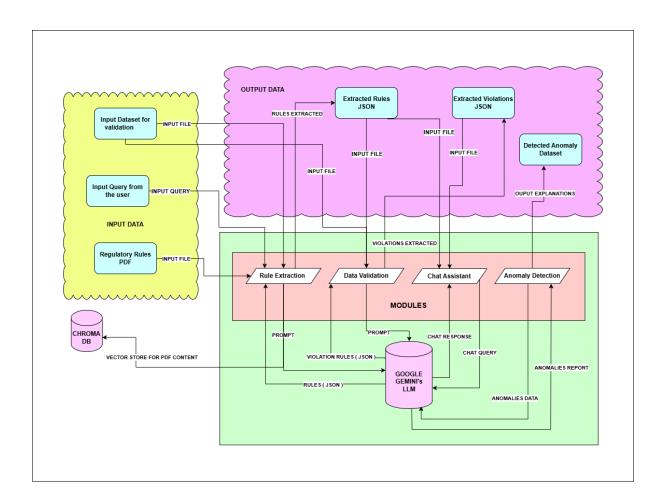
## 3. Auto-Calculated Contamination Rate:

- Instead of manually setting the contamination rate, it is estimated using a hybrid approach:
- o Isolation Forest-based Scoring: Model initially assumes a higher contamination rate.
- Z-Score Method: Identifies anomalies based on statistical deviations (2 standard deviations below mean).
- o Percentile-Based Approach: Flags the lowest 2% of anomaly scores as outliers.
- Final Contamination Rate: The maximum value between the Z-Score and Percentile-Based Approach is chosen.
- Ensures Stability: The contamination rate is clamped between 0.01 and 0.15 to avoid overfitting or under-detection.

### 4. Interactive UI & Visualization:

- o Generates **scatter plot with color-coded anomalies** (red for outliers).
- Users can download flagged anomalies (CSV) for further analysis.

# 2.5 Architecture Diagram



# 3. Implementation Details

# 3.1 Libraries & Technologies Used

- Backend: Python, LangChain, Pandas
- Frontend: Streamlit, Plotly, Altair
- Al Models: Google Gemini (LLM), Isolation Forest (ML)
- Database: ChromaDB (Vector Storage)
- Deployment: Streamlit Cloud, GitHub Actions

# 3.2 Optimization Strategies

- Batch Processing: Reduces API requests and improves validation efficiency.
- Vector Search (FAISS): Ensures fast and relevant retrieval of validation rules.
- Dynamic Contamination Rate: Improves anomaly detection accuracy.
- **Streamlit UI Enhancements:** Enhances user interaction with collapsible reports and real-time updates.

## 4. Challenges & Solutions

# 4.1 Handling Large PDFs

- Challenge: Extracting meaningful rules from complex regulatory documents.
- Solution: Used ChromaDB for vector storage and chunk-based processing.

# 4.2 Reducing API Usage & Cost

- Challenge: Generative AI models can be costly with frequent API calls.
- Solution: Batch validation & caching reduced unnecessary API calls.

# 4.3 Ensuring JSON Response Integrity

- Challenge: Al-generated responses sometimes contained incomplete JSON.
- Solution: Implemented parsing checks & regex fixes to ensure valid JSON.

# 4.4 Aligning Rule Extraction with Validation Logic

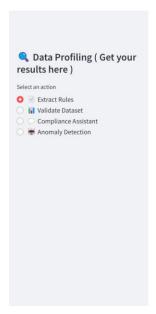
- Challenge: Extracted rules sometimes lacked dataset alignment.
- Solution: Fine-tuned AI prompts to improve rule relevance.

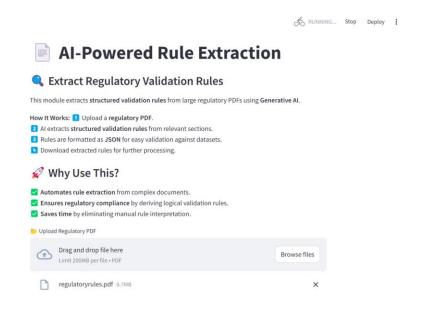
## 4.5 Making UI More Interactive

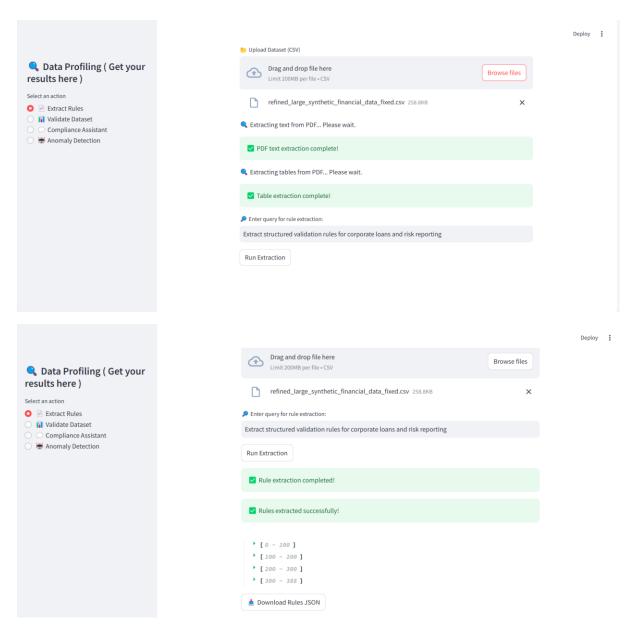
- Challenge: Users needed a clearer way to track validation & anomalies.
- Solution: Used dynamic progress bars, collapsible reports, and color-coded chat bubbles.

#### 5. Demo

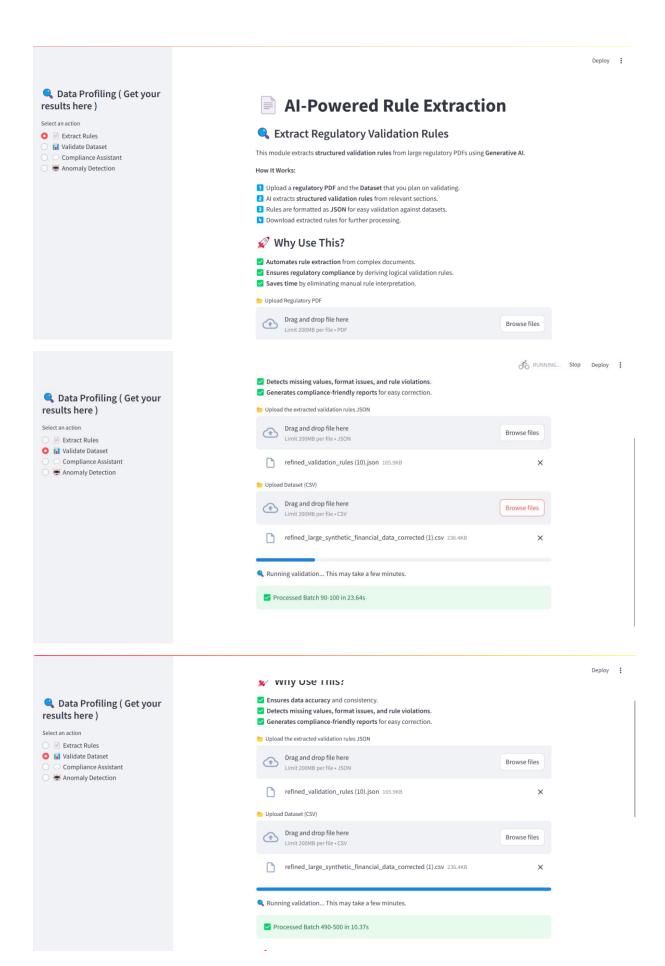
## 5.1 Rule Extraction (Al-Powered Regulatory Rule Extraction)

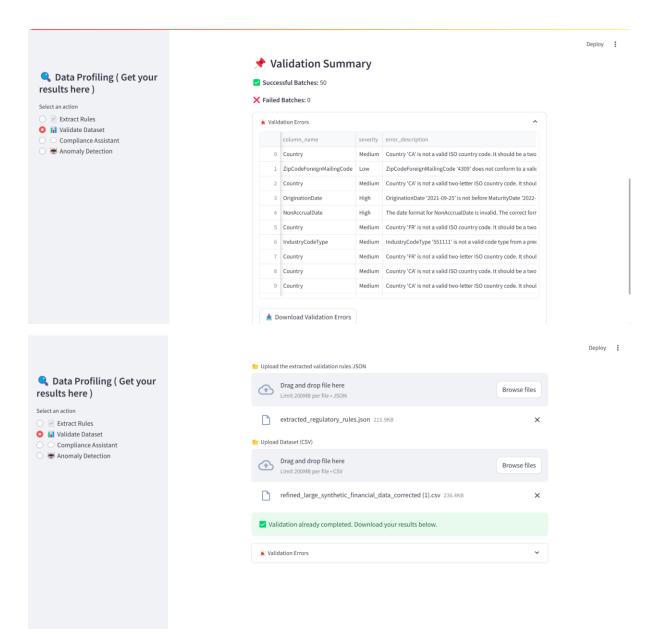




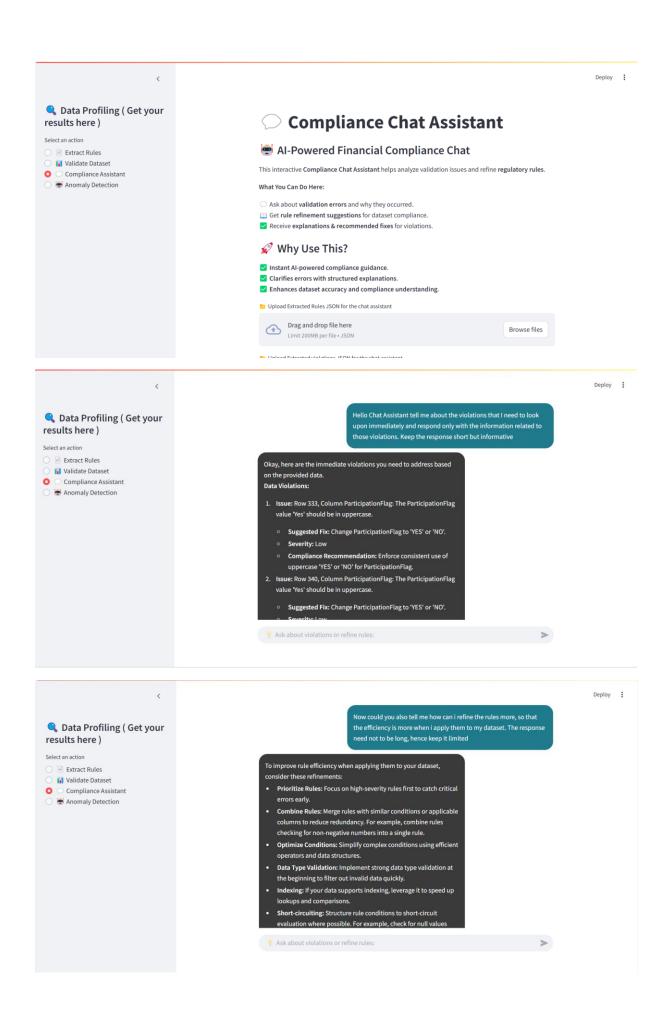


5.2 Dataset Validation (Automated Compliance Checking)

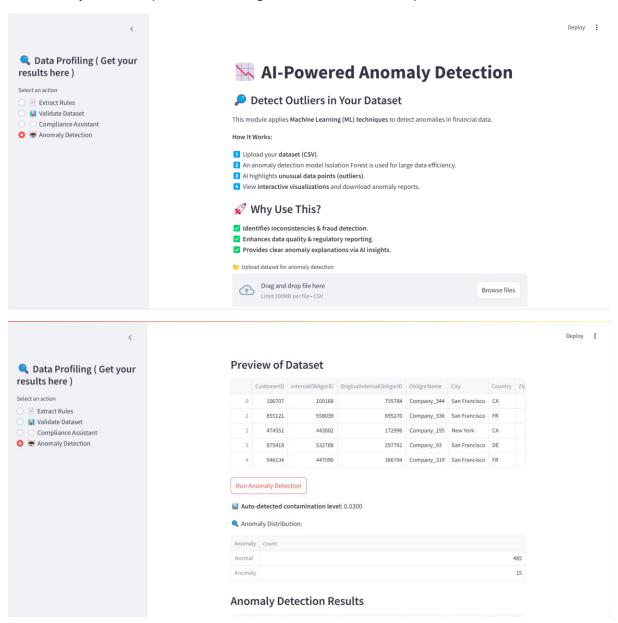


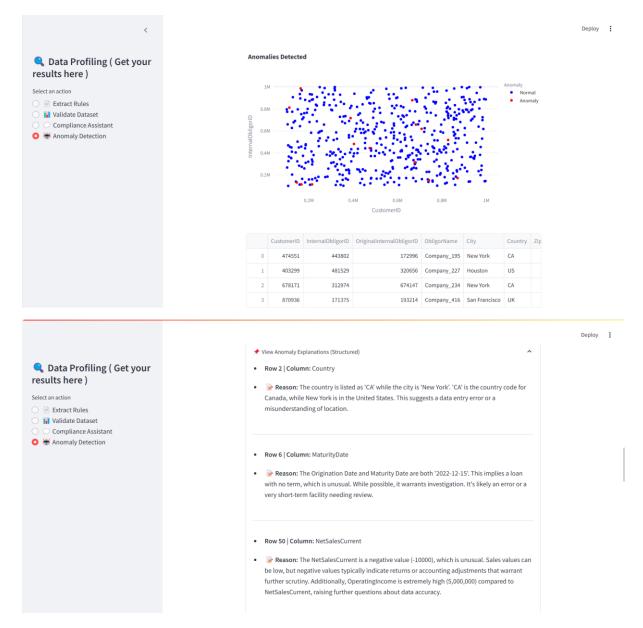


5.3 Compliance Chat Assistant (Interactive AI-Powered Queries)



# 5.4 Anomaly Detection (Machine Learning-Based Outlier Detection)





## 5. Future Scope

- Advanced Machine Learning for Anomaly Detection: Incorporate deep learning-based anomaly detection (e.g., Autoencoders, GANs). Implement ensemble models that combine multiple anomaly detection techniques for better accuracy
- Real-Time Data Validation & Streaming Support: Integrate with real-time data sources (Kafka, Spark Streaming). Perform continuous validation instead of batch processing.
- Automated Data Correction & Compliance Reports: Instead of just detecting violations, AI
  can suggest & auto-correct issues in datasets. Generate automated compliance reports for
  auditors & regulators.
- Adaptive AI-Powered Rule Learning: Enable AI to dynamically learn new validation rules over time based on past violations. Use reinforcement learning to improve compliance recommendations

 API Integration & Enterprise Deployment: Develop REST APIs for easy integration with existing banking & enterprise systems. Deploy as a cloud-based SaaS solution for scalability & security..

### 6. Conclusion

The Al-Powered Data Profiling & Compliance Assistant automates regulatory validation, improves dataset accuracy, and reduces compliance risks.

# **Key Takeaways:**

- Al-Driven Rule Extraction saves time & ensures compliance.
- Batch-Based Validation optimizes performance & reduces API costs.
- Unsupervised Anomaly Detection identifies potential data inconsistencies.
- Interactive AI Chat Assistant enhances compliance understanding.
- Scalable & Future-Ready with potential for real-time integration.

This system provides a **powerful**, **automated approach** to data validation, making regulatory compliance **faster**, **smarter**, **and more efficient**.