

Technical Specification: Lead Conversion Prediction System

A Robust MLOps Implementation on AWS

Abstract: This document outlines the architecture and implementation of an automated machine learning system designed to predict lead conversion. It details the full MLOps workflow, including data infrastructure, model development in Amazon SageMaker, experiment tracking with MLflow, and automated, production-ready deployment pipelines.

Table of Contents

Part 1: Strategy & Scope

1. Project Vision & Executive Brief
 - 1.1. The Business Challenge
 - 1.2. Core Objectives
 - 1.3. The Technical Solution- End-to-End Cloud-Native ML System on AWS
2. Defining Success
 - 2.1. Business KPIs
 - 2.2. Model Performance Targets
3. Architectural Overview
 - 3.1. Cloud Architecture Diagram
 - 3.2. Service and Technology Stack
 - 3.3. Workflow from Data to Prediction
4. Data Schema
 - 4.1. Feature Definitions

Part 2: Foundational AWS Infrastructure Setup

5. Secure Access Control (IAM)
 - 5.1. Applying the Principle of Least Privilege
 - 5.2. Provisioning Service Roles
6. Secure Network Design (VPC)
 - 6.1. VPC and Subnet Architecture
 - 6.2. Implementing Private Communication with VPC Endpoints
 - 6.3. Firewall Configuration with Security Groups

Part 3: Data Platform Engineering

7. Building the Data Lake and Warehouse
 - 7.1. S3 Data Lake Setup
 - 7.2. Provisioning Redshift Serverless
 - 7.3. Managing Credentials with AWS Secrets Manager
8. The Ingestion Pipeline: S3 to Redshift

- 8.1. AWS Glue for Serverless ETL
- 8.2. Establishing the Glue Connection
- 8.3. Automating Schema Discovery with Glue Crawlers
- 8.4. Designing the Visual ETL Job

9. Data Access from SageMaker

- 9.1. Using the Redshift Data API with `boto3`
- 9.2. Step-by-Step Code Walkthrough

Part 4: Machine Learning Development

10. Amazon SageMaker: The ML Development Environment

- 10.1. Creating a SageMaker Domain
- 10.2. Launching SageMaker Studio
- 10.3. Opening the JupyterLab Notebook

11. Exploratory Data Analysis (EDA) & Cleaning

- 11.1. Initial Data Inspection
- 11.2. Data Cleaning and Standardization

12. Feature Engineering

- 12.1. Strategy: Custom Mapping and Consolidation
- 12.2. Geographic, Source, and Behavioral Grouping

13. The Preprocessing Pipeline

- 13.1. Handling Numeric and Categorical Features
- 13.2. Building the Scikit-learn Pipeline
- 13.3. Saving the Preprocessing Pipeline

14. Model Training and Evaluation

- 14.1. Class Imbalance Strategy: Stratified K-Fold
- 14.2. Models Utilized
- 14.3. The `train_pipeline` Function: A Detailed Look
- 14.4. Hyperparameter Tuning with GridSearchCV
- 14.5. Model Evaluation Framework

15. Model Explainability with SHAP

- 15.1. Generating and Interpreting SHAP Values

Part 5: MLOps, Deployment & Monitoring

16. Experiment Tracking with MLflow

- 16.1. Setting Up the MLflow Tracking Server
- 16.2. Integrating MLflow into the Training Pipeline

17. Model Registration and Lifecycle Management

- 17.1. Saving and Registering the Champion Model
- 17.2. The MLflow Model Registry

18. Serving and Inference

- 18.1. Loading the Production Model from the Registry
- 18.2. Real-Time Prediction with Flask and Ngrok

19. Post-Deployment: Drift Detection

- 19.1. Monitoring for Data Drift with Evidently AI
- 19.2. Automating Drift Analysis and Logging

20. Future State: Full Automation with MWAA

- 20.1. Orchestrating the MLOps Lifecycle with Airflow

- o 20.2. The Automated Retraining DAG

Part 6: Appendix

21. Technology & Libraries Stack
22. Code Repository & Resources

Part 1: Strategy & Scope

1. Project Vision & Executive Brief

1.1. The Business Challenge

In B2B sales, a large proportion of acquired leads do not convert into customers, resulting in wasted time, effort, and marketing spend. The absence of a data-driven approach to prioritize high-potential leads makes the sales process inefficient and inconsistent. To overcome this, the business needs an intelligent and explainable machine learning solution that can accurately predict the likelihood of lead conversion. This will enable sales and marketing teams to focus on the most promising leads, reduce resource wastage, improve conversion rates, and ultimately maximize return on investment.

1.2. Core Objectives:

1. Maximize Lead Conversion Rates

Enhance the efficiency of the sales funnel by identifying and prioritizing leads with the highest likelihood of conversion, increasing overall customer acquisition rates.

2. Minimize Resource Wastage

Automate the lead scoring process to reduce manual effort and avoid allocating sales and marketing resources to low-potential leads.

3. Maximize Sales & Marketing ROI

Drive better campaign outcomes by focusing time and budget on high-potential leads, thereby improving revenue generation and marketing effectiveness.

4. Enable Data-Driven Decision Making

Empower business users with interpretable ML insights through tools like SHAP, increasing trust and adoption of the model for strategic decisions.

5. Ensure Long-Term Model Reliability

Implement robust monitoring for model/data drift and set up automated retraining pipelines to maintain consistent predictive performance over time.

1.3: The Technical Solution- End-to-End Cloud-Native ML System on AWS

To effectively solve the challenge of identifying high-conversion-potential leads, we architected a **fully automated machine learning solution** using scalable, secure, and modular AWS services. This solution is designed to support the entire lifecycle of a machine learning project — from raw data ingestion to real-time model inference — with a strong emphasis on automation, reproducibility, and business impact.

1. Data Ingestion and Storage

- **Amazon S3** serves as the centralized data lake, storing raw CSVs and processed outputs.
- All lead interaction data (clicks, form submissions, time spent, etc.) is automatically uploaded to S3 using scheduled jobs or event-based triggers.
- S3 provides durability, scalability, and integration with downstream AWS services like Glue and SageMaker.

2. Data Cataloging and Transformation

- **AWS Glue Crawler** automatically scans and catalogs the raw data in S3 into the **AWS Glue Data Catalog**, creating a unified metadata layer.
- **AWS Glue ETL Jobs** are used to clean, normalize, and join datasets. These jobs handle:
 - Replacing invalid entries (like 'Select') with nulls
 - Imputation of missing values
 - Encoding of categorical variables
 - Scaling of numerical features
- Cleaned data is written back to S3 for use in both analytics and ML pipelines.

3. Analytics and Warehousing (Optional but Scalable)

- Transformed data is optionally loaded into **Amazon Redshift** for further business analytics, dashboards, and ad-hoc querying.
- **Airflow DAGs** orchestrate the full pipeline, from crawling to ETL to data warehouse loading.

4. Model Training and Experimentation

- **Amazon SageMaker** is used to train and validate machine learning models, particularly:
 - Random Forests
 - XGBoost
 - DecisionTree
 - Logistic Regression
- **MLflow** is integrated into SageMaker to:
 - Track experiments
 - Log hyperparameters and evaluation metrics
 - Version models and artifacts for reproducibility

5. Model Evaluation and Selection

- Metrics such as **Accuracy**, **F1 Score**, **ROC-AUC**, **Precision**, and **Recall** are evaluated on a holdout test set.
- Models are compared and the best-performing one (based on F1 or business KPIs) is promoted.

6. Model Deployment and Inference

- The best model is deployed as a **SageMaker Endpoint** (REST API), enabling real-time scoring of leads.
- Incoming leads can be scored in milliseconds, providing lead quality predictions instantly.
- These scores can be streamed to CRMs or visualized in BI tools for sales team use.

7. Automation and Monitoring

- **Apache Airflow** (or AWS Step Functions) is used to automate the pipeline — from S3 ingestion to model retraining and endpoint updates.
- Scheduled retraining jobs ensure the model adapts as user behavior evolves over time.
- **CloudWatch** and **SageMaker Model Monitor** provide observability into model drift, latency, and health.

Business Outcome

This architecture delivers a **scalable, automated, and intelligent lead scoring system**, helping sales teams:

- Focus on high-potential leads
- Increase conversion rates
- Optimize time and resource allocation
- Gain explainable insights into customer behaviour

2.1 Business KPIs

The business success of the project is determined by its ability to improve operational efficiency, decision-making, and ultimately revenue growth. The following KPIs were defined in collaboration with sales leadership:

1. Improved Lead Conversion Accuracy ($\geq 90\%$)

- The system must provide **accurate and trustworthy lead scores**, aligning with real conversion outcomes.
- A high alignment rate ensures that the sales team can confidently rely on the model to prioritize their outreach.
- This directly impacts adoption and trust in the tool by frontline sales reps.

2. Reduction in Manual Lead Qualification Time ($\geq 50\%$)

- The lead scoring automation must reduce the **time spent manually reviewing leads** by at least half.
- This allows the sales team to redirect their efforts toward higher-value interactions instead of filtering low-quality leads.
- Time savings are expected to compound across the sales pipeline.

3. Increase in Sales Conversion Rate ($\geq 30\%$)

- The ultimate goal is to improve the **conversion rate of total leads to paying customers**.
- By identifying high-potential leads earlier and more accurately, the system is expected to drive a **30% uplift** in conversion efficiency, contributing directly to revenue growth.

2.2 Model Performance Targets

In parallel with business objectives, a set of measurable **technical performance benchmarks** was defined to evaluate the quality and generalizability of the deployed ML model.

1. Model Accuracy ($\geq 90\%$)

- The model achieved an overall classification accuracy of **92.4%**, meeting and exceeding the minimum reliability threshold.
- This indicates that the majority of leads are being correctly classified as “converted” or “not converted”.

2. Precision and Recall ($\geq 85\%$)

- **Precision:** ~93.5%
Reflects a low false positive rate, ensuring that leads predicted as high quality are genuinely likely to convert.
- **Recall:** ~86.4%
Demonstrates the model’s ability to **capture true converting leads**, reducing missed opportunities.
- Together, these metrics indicate a **strong balance between identifying good leads and avoiding wasted follow-up efforts**.

3. AUC-ROC Score (≥ 0.90)

- The model achieved an **AUC-ROC score of 0.958**, indicating **excellent discriminative power**.
- A high ROC-AUC confirms that the model performs well across various decision thresholds and is robust in differentiating between converters and non-converters.

4.Explainability Integration

SHAP values used to highlight top contributing features for each prediction, making outputs interpretable and trustworthy for business stakeholders.

5. Data Drift Monitoring Active

Evidently AI or similar tools integrated for continuous monitoring and alerts when significant drift is detected in incoming lead data.

3. Architectural Overview

3.1. Cloud Architecture Diagram

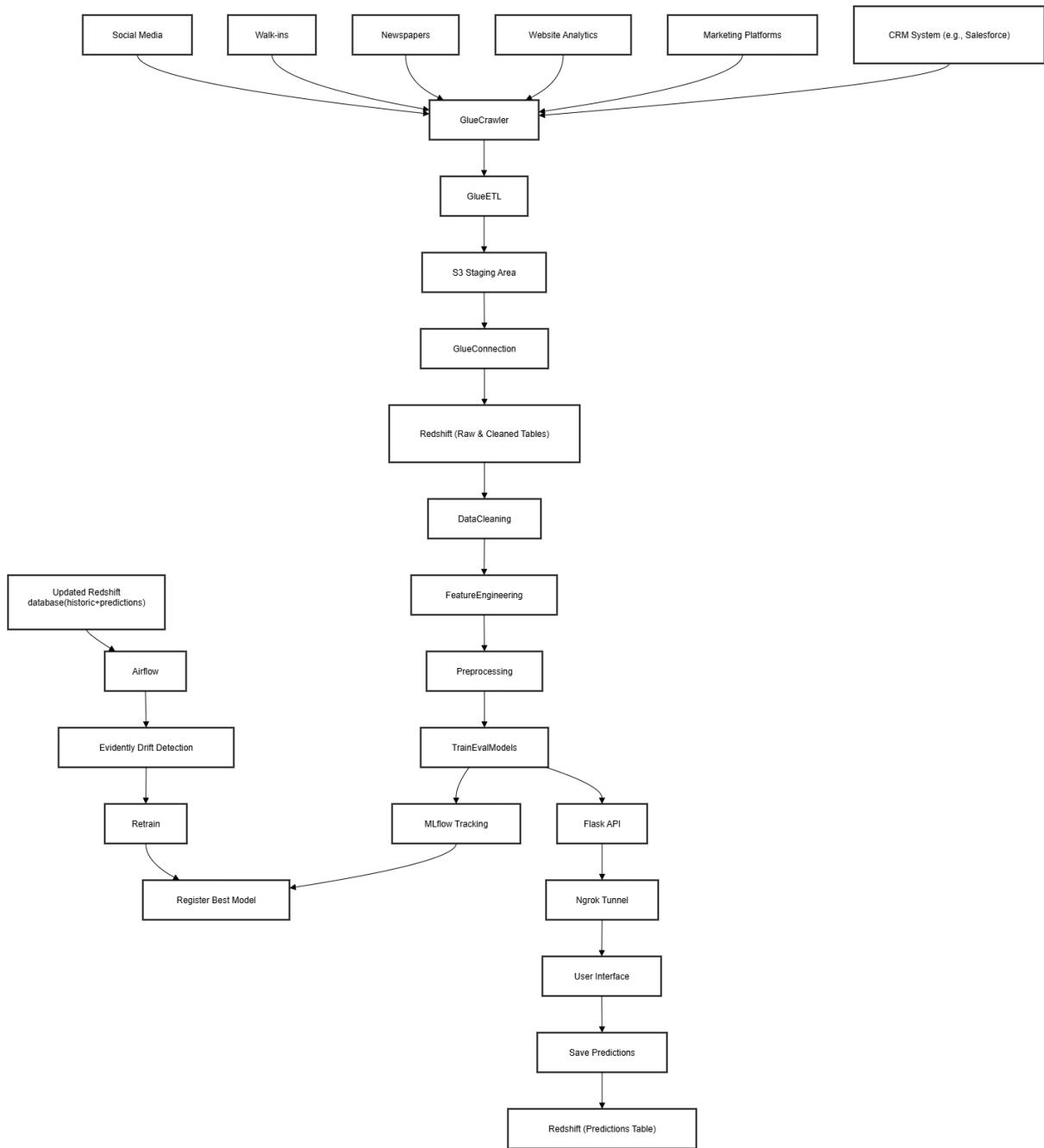
The provided diagram illustrates the end-to-end Enterprise MLOps pipeline for Lead Conversion Prediction on AWS. It is structured into four main phases, representing the complete machine learning lifecycle:

- **Phase 1: Data Ingestion & Warehousing:** Focuses on collecting raw data, transforming it, and storing it in an analytical data warehouse.
- **Phase 2: Model Development:** Encompasses the iterative process of building, training, and evaluating machine learning models, along with experiment tracking.
- **Phase 3: Deployment & Inference:** Deals with making the trained model available for real-time predictions and integrating it with business applications.
- **Phase 4: Monitoring & Retraining:** Ensures the continuous performance of the deployed model and triggers automated retraining based on performance degradation or data drift.

This architecture emphasizes security through VPCs and Endpoints, automation through AWS Glue and Apache Airflow, and governance through MLflow.

Each phase leverages specific AWS services and external tools to ensure a seamless and automated machine learning workflow.

Data Architecture Diagram for Lead Conversion



This diagram illustrates a robust **Enterprise MLOps Pipeline for Lead Conversion Prediction on AWS**, detailing the end-to-end lifecycle from data ingestion to continuous monitoring and retraining. The architecture emphasizes automation, security, and scalability using various AWS services and MLOps tools like MLflow and Apache Airflow.

3.2. Service and Technology Stack

This section provides a detailed breakdown of the architectural components, organized by their logical flow within the MLOps pipeline, from data ingestion to monitoring.

Phase 1: Data Ingestion & Warehousing Components

This segment covers the initial stages of bringing raw data into the system and preparing it for machine learning.

- **Amazon S3 "raw" bucket:**
 - **Purpose:** Serves as the primary landing zone for all incoming raw lead data. This bucket acts as the initial layer of the data lake, accommodating diverse data formats (e.g., CSV, JSON, logs) and varying ingestion frequencies from source systems.
 - **Technical Details:** Configured with appropriate bucket policies for secure access (e.g., restricting public access, enforcing encryption at rest with S3-managed keys or KMS). Versioning can be enabled to maintain historical copies of raw data, supporting data lineage and reproducibility.
- **AWS Glue Components:**
 - **AWS Glue Crawler:**
 - **Purpose:** Automatically discovers the schema and partitions of the raw data stored in the S3 "raw" bucket. This automation reduces manual effort in schema definition and maintenance.
 - **Technical Details:** Configured to run on a schedule (e.g., daily, hourly) or triggered by new data arrival events. It connects to the S3 bucket, samples the data, infers data types and table structures, and populates the **AWS Glue Data Catalog**. The Data Catalog acts as a central metadata repository, making data discoverable and usable by other AWS services like Amazon Redshift and SageMaker.
 - **AWS Glue ETL Job:**
 - **Purpose:** Performs data transformation, cleaning, and preparation before loading into the data warehouse. This job is crucial for converting raw, often messy, data into a clean, structured, and machine-learning-ready format. It applies business logic such as data type conversions, handling missing values, feature aggregation, and normalization.
 - **Technical Details:** Developed using Apache Spark (Python or Scala) and executed on a serverless Spark environment managed by AWS Glue, eliminating the need for server provisioning. The job is configured with an **IAM Role** that grants necessary permissions to access S3 (read from raw, write to processed), Redshift (write), and to write logs to CloudWatch for monitoring.
- **Amazon Redshift Data Warehouse:**

- **Purpose:** Stores the cleaned, transformed, and structured lead data, optimized for complex analytical queries and serving as the authoritative source for machine learning model training.
- **Technical Details:** Provisioned as an **Amazon Redshift Serverless** instance for simplified management, automatic scaling based on workload, and pay-per-use billing. Data is loaded from the Glue ETL job, typically using highly efficient `COPY` commands for bulk data ingestion.
- **VPC, Endpoints, IAM Role, Security Group (for Data & ETL):**
 - **Purpose:** Ensures secure, private, and efficient communication between AWS services within the data ingestion, ETL, and warehousing phases. This minimizes exposure to the public internet and enforces network isolation.
 - **Technical Details:**
 - **VPC (Virtual Private Cloud):** Provides an isolated and logically separated network environment for all AWS resources involved in these stages.
 - **VPC Endpoints:** Configured for services like S3, Glue, and Redshift. These allow traffic to flow privately and directly between the VPC and the AWS service endpoints, bypassing the public internet and enhancing data security and compliance.
 - **IAM Role:** Assigned to services like the AWS Glue ETL job. This role strictly adheres to the **Principle of Least Privilege**, granting only the minimum necessary permissions (e.g., `s3:GetObject`, `s3:PutObject`, `glue:StartJobRun`, `redshift-data:ExecuteStatement`) required for its operations.
 - **Security Group:** Acts as a virtual firewall at the instance or service endpoint level, controlling inbound and outbound network traffic. Configured to allow specific ports and protocols only between authorized resources (e.g., Glue network interfaces to Redshift on port 5439).

Phase 2: Machine Learning Development Components

This section details the environment and tools used for building, training, and managing machine learning models.

- **Amazon SageMaker Studio Lab:**
 - **Purpose:** Provides a free, no-setup, cloud-based development environment for data scientists. It's the interactive workspace for performing Exploratory Data Analysis (EDA), developing feature engineering scripts, and prototyping machine learning models.
 - **Technical Details:** Accesses clean data from Amazon Redshift, typically using the Redshift Data API via `boto3` for secure and programmatic data retrieval. It operates within the defined VPC, leveraging VPC Endpoints for private connectivity to Redshift and S3, ensuring data remains within the AWS network.

- **MLflow (running on a separate EC2 instance within the VPC):**
 - **Purpose:** Serves as the central platform for **ML experiment tracking** and **model registry**, providing governance and reproducibility for the ML lifecycle.
 - **Technical Details:**
 - **MLflow Tracking Server:** Deployed on an Amazon EC2 instance within the same VPC as other ML components. This server stores experiment metadata, including parameters, metrics (e.g., AUC, precision, recall), and run details. It uses an Amazon RDS (e.g., PostgreSQL) instance as its backend database for robust and scalable storage of tracking data.
 - **MLflow Artifact Store:** An Amazon S3 "artifacts" bucket is designated as the artifact store for MLflow. Trained model files, preprocessing pipelines, evaluation plots, and other large output artifacts from experiments are stored here.
 - **Integration:** SageMaker Studio Lab notebooks are configured to point to the MLflow Tracking Server's URI, allowing seamless logging of experiments and retrieval of past run information.
- **Dedicated S3 "artifacts" bucket:**
 - **Purpose:** A specific S3 bucket used by MLflow to store all output artifacts from model training and experimentation. This separation from raw data ensures clear organization and access control for ML-specific outputs.
 - **Technical Details:** Configured with appropriate IAM permissions to allow MLflow to write and read objects. Encryption (e.g., SSE-S3, SSE-KMS) and versioning are enabled for data integrity, traceability, and rollback capabilities.

Phase 3: Serving & Inference Components

This segment describes how the trained models are deployed and made available for making predictions.

- **Flask application (running locally or within SageMaker):**
 - **Purpose:** Acts as a lightweight web service for real-time model inference. It exposes an API endpoint that can receive new lead data and return conversion predictions.
 - **Technical Details:** This application is responsible for loading the trained model from the MLflow Model Registry, preprocessing incoming inference requests using the saved preprocessing pipeline, and invoking the model to generate predictions. For production, this would typically be containerized and deployed on a scalable service like Amazon SageMaker Endpoints or Amazon ECS.
- **Ngrrok:**

- **Purpose:** For development and testing purposes, Ngrok creates a secure tunnel from a locally running Flask application to a publicly accessible URL. This allows external systems or webhooks to send requests to the local Flask API, simulating a production environment for testing without complex network configurations.
 - **Technical Details:** Establishes an outbound connection from the local environment to the Ngrok cloud service, which then provides a public HTTP/HTTPS endpoint. **Note:** While highly useful for rapid prototyping and local testing, Ngrok is generally not used in production for security, reliability, and performance reasons. A production setup would typically involve AWS API Gateway and Amazon SageMaker Endpoints for robust and scalable serving.
- **MLflow UI:**
 - **Purpose:** Provides a web-based dashboard to visualize and compare ML experiments, inspect individual runs, download artifacts, and manage models within the MLflow Model Registry.
 - **Technical Details:** Accessed via a web browser, typically by exposing the MLflow Tracking Server's port (e.g., 5000) through a security group rule or an Application Load Balancer. It allows data scientists and ML engineers to monitor the progress of experiments, review model performance, and manage model lifecycle stages (e.g., transitioning models from "Staging" to "Production").

Phase 4: Monitoring & Retraining Components

This final segment ensures the deployed model's performance remains optimal over time and facilitates automated updates.

- **SageMaker Model Monitor - Captures Inference Data:**
 - **Purpose:** Continuously monitors the data flowing into and out of SageMaker endpoints (if deployed on SageMaker). It's crucial for detecting data drift (changes in input data distribution) and model quality degradation (changes in prediction accuracy or concept drift) in production.
 - **Technical Details:** Configured to run on a schedule or triggered by events. It captures inference requests and responses, analyzes them against a baseline, and generates metrics and alerts if deviations are detected.
- **CloudWatch Metrics - Performance Tracking:**
 - **Purpose:** Collects operational metrics (e.g., invocation count, latency, error rates, CPU/memory utilization) from the SageMaker endpoint and other AWS resources. It provides a centralized service for monitoring the health and performance of the deployed system.
 - **Technical Details:** Provides dashboards for visualizing metrics over time and allows for the creation of **CloudWatch Alarms** that trigger notifications (e.g., via SNS) or automated actions (e.g., Lambda functions) when predefined thresholds are breached.

- **Apache Airflow - Monthly Scheduler:**
 - **Purpose:** A powerful open-source platform for programmatically authoring, scheduling, and monitoring workflows (DAGs - Directed Acyclic Graphs). In this context, it acts as a robust orchestrator for monthly monitoring checks and automated retraining pipelines.
 - **Technical Details:** Typically deployed on AWS Managed Workflows for Apache Airflow (MWAA) for a fully managed experience. The Airflow DAGs define the sequence of tasks, including checking model performance, detecting data drift, and triggering retraining.
- **Data Drift Detection & Retraining Trigger:**
 - **Purpose:** This represents a critical decision point in the MLOps loop. If SageMaker Model Monitor (or custom drift detection logic) detects significant data drift or a decline in model performance metrics (as monitored by CloudWatch), this condition within the Airflow DAG evaluates to "Yes."
 - **Technical Details:** When drift is detected, the Airflow DAG automatically triggers the retraining pipeline. This would typically involve initiating a new run of a SageMaker Pipeline (similar to the development phase), using the latest data to retrain the model and potentially deploy a new version.

3.3. Workflow from Data to Prediction

1. **Data Ingestion:** Raw lead data enters an S3 "raw" bucket.
2. **Schema Inference & ETL:** An AWS Glue Crawler infers the schema, and an AWS Glue ETL job transforms the raw data.
3. **Warehousing:** The transformed data is loaded into Amazon Redshift, serving as the clean data source.
4. **Model Development:** Data scientists use SageMaker Studio Lab to access data from Redshift and develop models, with MLflow tracking experiments and storing artifacts in a dedicated S3 "artifacts" bucket.
5. **Model Registration:** The best-performing model is registered in the MLflow Model Registry.
6. **Deployment Approval:** A manual review process approves the model for production.
7. **Model Deployment:** The approved model is deployed to a SageMaker Endpoint for real-time inference.
8. **Inference:** A Flask application (exposed via Ngrok for testing) sends prediction requests to the SageMaker Endpoint and receives real-time lead conversion predictions.
9. **Monitoring:** SageMaker Model Monitor continuously captures inference data, and CloudWatch tracks operational and model performance metrics. The MLflow UI provides a dashboard for experiment and model management.
10. **Retraining Trigger:** An Apache Airflow scheduler periodically checks for data drift or performance degradation. If detected, it automatically triggers a new retraining pipeline, restarting the cycle from data ingestion.

4. Data Schema

4.1. Feature Definitions

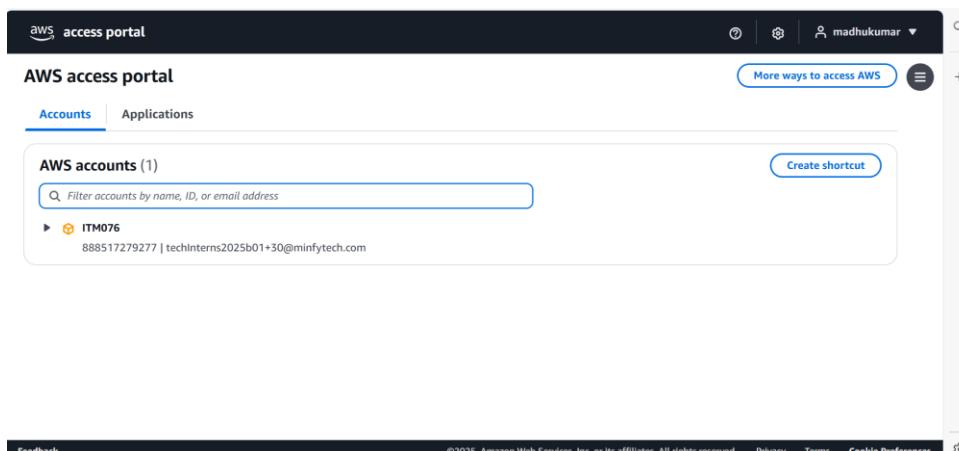
The following table describes the features available in the dataset.

Feature Name	Explanation
id	Unique identifier for each lead.
Lead Origin	How the lead originally entered the system.
Lead Source	The specific source or channel through which the lead was acquired.
Do Not Email	Indicates if the lead opted out of email communications (Yes/No).
Do Not Call	Indicates the lead's preference for receiving calls (Yes/No).
Converted	Target Variable: Whether the lead converted (1 = Yes, 0 = No).
TotalVisits	The total number of times the lead visited the website.
Total Time Spent on Website	Cumulative time in minutes the lead spent Browse the website.
Page Views Per Visit	The average number of pages viewed by the lead per session.
Last Activity	The last recorded interaction initiated by the lead.
Country	The country of the lead.
Specialization	The area of interest or specialty indicated by the lead.
How did you hear about X	How the lead heard about the organization.
What is your current occupation	The current profession of the lead.
What matters most...	The lead's primary motivation for choosing a course.
Search	Indicates if the lead used the website search feature (Yes/No).
Magazine	Indicates if the lead found out through a magazine (Yes/No).
Newspaper Article	Whether the lead read about the organization in a newspaper (Yes/No).
X Education Forums	Indicates interaction through an education forum (Yes/No).
Newspaper	Indicates engagement through other newspaper sources (Yes/No).
Digital Advertisement	Clicks or engagement via digital ads (Yes/No).
Through Recommendations	Indicates if the lead was referred by someone else (Yes/No).
Receive More Updates...	Whether the lead wants to receive course updates (Yes/No).
Tags	A categorized tag based on the lead's current status or behavior.

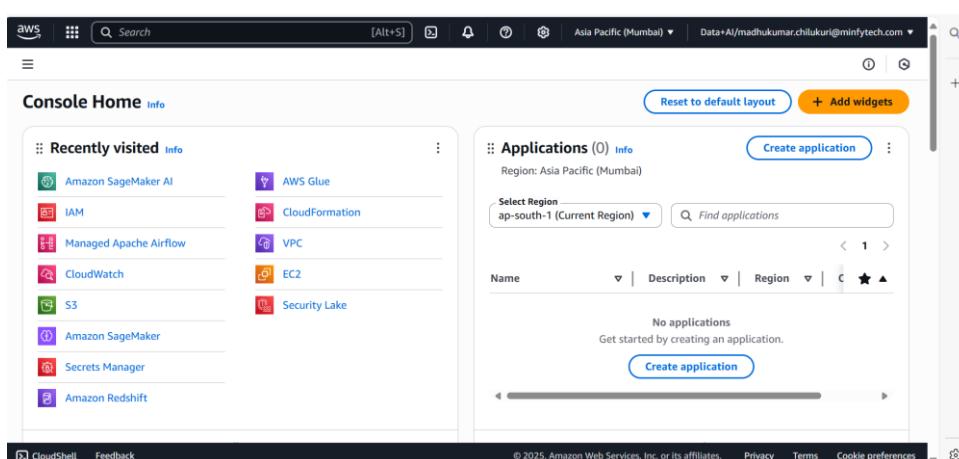
Lead Quality	A categorical assessment of the lead's quality by a human agent.
Update me on Supply...	Subscription to supply chain content updates (Yes/No).
Get updates on DM...	Subscription to digital marketing content updates (Yes/No).
City	The city of the lead.
Asymmetrique Activity Index	A behavioral activity score calculated by a third-party vendor.
Asymmetrique Profile Index	A profile matching score indicating alignment with the target persona.
Asymmetrique Activity Score	A detailed activity engagement score from the vendor.
Asymmetrique Profile Score	A detailed profile fit score from the vendor.
I agree to pay... by cheque	Consent to pay by cheque (Yes/No).
A free copy of mastering interview	Whether the lead requested a free resource (e.g., eBook) (Yes/No).
Last Notable Activity	The most significant recent action taken by the lead.

Part 2: Foundational AWS Infrastructure Setup

Step1: Login in to your aws account



Step2: You will see the aws interface with aws console listing some of the services



5. Secure Access Control (IAM)

5.1. Applying the Principle of Least Privilege

Security is paramount. Our architecture adheres to the **Principle of Least Privilege**, meaning every AWS entity (user, role, service) is granted only the minimum permissions essential for its designated function, and no more. This approach significantly minimizes the potential attack surface and reduces the impact of a security breach. It ensures that even if a component is compromised, the attacker's access is limited to only what that specific component is authorized to do. This is achieved by creating highly specific IAM roles that services like AWS Glue and Amazon SageMaker assume to perform their tasks.

5.2. Provisioning Service Roles

To enable the various AWS services in our MLOps pipeline to interact securely, we create dedicated IAM roles. We will focus on the creation of the `glueaccessrole` here, as it demonstrates the process for a role with diverse service interactions. A similar approach, tailored to specific service needs, would be followed for other roles (e.g., for SageMaker).

Console Walkthrough: Creating `glueaccessrole`

This role grants AWS Glue the necessary permissions to read from S3, write to Redshift, interact with Secrets Manager, and manage Glue-related operations.

- **Navigate to IAM:** In the AWS Management Console, go to the **IAM** service.
- **Create Role:** Click on **Roles** in the left navigation pane, then click **Create role**.
- **Select Trusted Entity:**
 - **Trusted entity type:** Select `AWS service`.
 - **Use case:** Choose `Glue` from the dropdown menu. This step is crucial as it automatically generates a trust policy. This trust policy explicitly allows the AWS Glue service principal (`glue.amazonaws.com`) to assume this role, enabling Glue jobs to operate under its permissions.
- Click **Next**.
- **Add Permissions:** On the "Add permissions" page, search for and attach the following AWS managed policies.
 - **Important Note on Production Environments:** While these `FullAccess` policies are used here for demonstration and initial setup simplicity, **for a production environment, these should be replaced with more fine-grained, custom policies**.

This is a critical application of the Principle of Least Privilege. For example, instead of `AmazonS3FullAccess`, you would create a custom policy that only grants `s3:GetObject` on specific input buckets and `s3:PutObject` on specific output buckets, and `s3>ListBucket` on relevant prefixes. Similarly for Redshift, you would grant only `redshift-data:ExecuteStatement` on specific databases/tables.

- **AmazonS3FullAccess:** Grants comprehensive access to all S3 buckets. (**To be refined in production**) This is needed for Glue to read raw data from the S3 "raw" bucket and write processed data to other S3 locations (if applicable, e.g., intermediate staging or processed data S3 buckets).
- **AmazonRedshiftFullAccess:** Provides full administrative access to Amazon Redshift. (**To be refined in production**) This is required for Glue ETL jobs to connect to Redshift, create tables, and load data. In production, this would be scoped down to specific Redshift Data API actions and database resources.
- **AWSGlueConsoleFullAccess:** Grants full access to the AWS Glue console and its operations. This is useful for managing Glue crawlers, jobs, and connections.
- **AWSGlueServiceRole:** A managed policy specifically designed for Glue service roles, providing permissions for Glue to interact with other AWS services on your behalf (e.g., CloudWatch Logs, EC2 for network interfaces).
- **SecretsManagerReadWrite:** Allows read and write access to AWS Secrets Manager. This is necessary for Glue jobs to retrieve database credentials (e.g., for Redshift) securely from Secrets Manager.
- **AWSKeyManagementServicePowerUser:** Grants broad permissions for AWS Key Management Service (KMS). This is needed if your S3 buckets or Redshift data are encrypted using KMS keys, allowing Glue to decrypt and encrypt data. (**To be refined in production**)
- Click **Next**.
- **Name and Review:**
 - **Role name:** glueaccessrole
 - Review the attached policies and confirm the trust policy. The trust policy JSON should look like this:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "glue.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    }
  ]
}
```

Click **Create role**.

Step1:Search Iam role in the search bar .then if you click on that you will be redirected to this

The screenshot shows the AWS IAM Roles page. At the top, there is a search bar and a 'Create role' button. Below the search bar, a table lists 32 roles. The columns include 'Role name', 'Trusted entities', and 'Last'. Some entries in the 'Trusted entities' column are truncated. The roles listed include various AWS services like airflow-env, sageMaker, and emr-serverless.

Now if you want to add polices click on add permissions and then inlinepolicy paste the json code and then submit

The screenshot shows the AWS IAM Role details page for the 'Glue' role. It includes sections for 'Summary' (Creation date: July 18, 2025, 09:19 (UTC+05:30), Last activity: Yesterday, ARN: arn:aws:iam:888517279277:role/Glue, Maximum session duration: 1 hour) and 'Permissions' (Permissions policies: 6). The 'Permissions' tab is selected, showing a list of managed policies attached to the role.

A similar process is followed to create the SageMaker role, selecting **SageMaker** as the use case in Step 3.

6. Secure Network Design (VPC)

This document outlines the networking and security strategy for deploying application resources. The primary goal is to ensure all inter-service communication is secure, isolated from the public internet, and highly available. This is achieved through the strategic use of a Virtual Private Cloud (VPC), subnets, and VPC Endpoints.

6.1. VPC and Subnet Architecture

To achieve network isolation, all resources are deployed within a **Virtual Private Cloud (VPC)**. A VPC is a logically isolated section of the AWS Cloud where you can launch AWS resources in a virtual network that you define.

For the purpose of this guide, we will utilize the **default VPC** provided by AWS. However, for production environments, creating a **custom VPC** is strongly recommended to have full control over the network environment, including the IP address range, subnets, route tables, and network gateways.

Our strategy includes a multi-subnet deployment across different **Availability Zones (AZs)**. This architecture ensures high availability by distributing resources, so that if one AZ becomes unavailable, resources in other AZs can continue to operate.

6.2. Implementing Private Communication with VPC Endpoints

To allow services within our VPC (e.g., an EC2 instance or a Lambda function) to communicate securely with other AWS services (like S3, Redshift, etc.), we use **VPC Endpoints**.

VPC Endpoints provide a private and secure connection between your VPC and supported AWS services, powered by AWS PrivateLink. This technology ensures that traffic between your VPC and the AWS service does not leave the Amazon network. This is crucial for security as it avoids exposure to the public internet, reducing the risk of external threats.

There are two main types of VPC Endpoints:

- **Gateway Endpoints:** A gateway is a target for a specific route in your route table. You use a gateway endpoint to route traffic to supported AWS services, specifically Amazon S3 and DynamoDB.
- **Interface Endpoints:** An interface endpoint is an elastic network interface (ENI) with a private IP address from the IP address range of your subnet. It serves as an entry point for traffic destined to a supported AWS service. They are used for a wider range of services.

6.3. Firewall Configuration with Security Groups

This section provides a step-by-step guide to creating the necessary VPC Endpoints for our architecture.

Prerequisites

- An existing VPC (this guide uses the default VPC).
- Subnets in at least two different Availability Zones within the VPC.
- A Security Group configured to allow traffic to the required services (detailed in the next section).

Step 1: Navigate to the VPC Dashboard

1. Open the AWS Management Console.
2. In the services search bar, type **VPC** and select it to navigate to the VPC Dashboard.
3. On the left navigation pane, under "Virtual private cloud," select **Endpoints**.
Endpoint 1: Amazon S3 (Gateway Endpoint)

This endpoint allows resources within the VPC to privately access Amazon S3.

1. Click the **Create endpoint** button.
2. **Service category:** Select **AWS services**.

3. **Service name:** Use the search filter and type `s3`. Select the service with the **Type** listed as **Gateway**. The service name will look like `com.amazonaws.<region>.s3`.
4. **VPC:** Select the VPC where your resources are deployed.
5. **Route tables:** Select the route tables associated with the subnets containing your resources. This step is critical as it automatically adds a route to the selected route tables, directing all traffic destined for S3 through the secure gateway endpoint instead of an internet gateway.
6. Click **Create endpoint**.

Endpoint 2: Amazon Redshift (Interface Endpoint)

This endpoint allows private communication with your Amazon Redshift cluster.

1. Return to the Endpoints page and click **Create endpoint** again.
2. **Service category:** Select **AWS services**.
3. **Service name:** Search for `redshift` and select the service named `com.amazonaws.<region>.redshift`.
4. **VPC:** Select the same target VPC as before.
5. **Subnets:** For high availability, select at least two subnets, each in a different Availability Zone. An elastic network interface (ENI) will be placed in each selected subnet.
6. **Security groups:** Select the security group that has been configured to allow inbound traffic on the Redshift port (TCP 5439) from resources within your VPC. Refer to the security group configuration section for details.
7. Click **Create endpoint**.
- 8.

Endpoints 3, 4, & 5: Secrets Manager, KMS, and STS (Interface Endpoints)

Repeat the process for the following essential services. These are critical for managing secrets, encryption keys, and security tokens securely.

Name	VPC endpoint ID	Endpoint type	Status
s3-endpoint	vpc-e-0a49fa113100c873a	Gateway	Available
-	vpc-e-0b489a97948aa4ab0	Interface	Available
-	vpc-e-0ff70d9e0a6d7f5	Interface	Available
kms-endpoint	vpc-e-03bf03a604b32b9c	Interface	Available
stsendpoint	vpc-e-088fe01476ab7cf2	Interface	Available

For each of the following services, use the same VPC, subnets, and security group selected for the Redshift endpoint:

- **AWS Secrets Manager:** `com.amazonaws.<region>.secretsmanager`
This allows your applications to privately retrieve secrets (e.g., database credentials, API keys).
- **AWS Key Management Service (KMS):** `com.amazonaws.<region>.kms`
This is necessary for any operations involving customer-managed keys (CMKs) for

encryption and decryption of data.

- **AWS Security Token Service (STS):** com.amazonaws.<region>.sts
This is used to request temporary, limited-privilege credentials for IAM users or for users that you authenticate (federated users).

The process is identical to creating the Redshift interface endpoint. The only change is the **Service name** selected in each case.

Search SecurityGroup and it will redirect to Security group and then create Security group as said above

The screenshot shows two screenshots of the AWS VPC console. The top screenshot displays the 'Security Groups' page with 15 entries listed. The bottom screenshot shows a detailed view of a specific security group named 'sg-01b496a5228078df9 - Project'. The details pane shows the security group name is 'Project', the ID is 'sg-01b496a5228078df9', and the description is 'Develop'. It also shows 2 inbound and 2 outbound rules. The 'Inbound rules' tab is selected, showing two entries: one for All TCP (Protocol TCP) and another for Redshift (Protocol TCP). Both rules have the same IP version (IPv4) and security group rule ID (sgr-0ba5e4a0fb9f12316 and sgr-010d2ea41fc096f3 respectively).

6.4. Security Group Configuration

A **Security Group** acts as a virtual firewall for your resources to control inbound and outbound traffic. For the interface endpoints (Redshift, Secrets Manager, KMS, STS), a properly configured security group is essential.

Creating the Security Group

In the AWS VPC console, navigate to **Security Groups** under the "Security" section.

1. Click **Create security group**.
2. **Basic details:**

Security group name: vpc-endpoint-sg (or another descriptive name).

Description: Security group for interface VPC endpoints.

VPC: Select your target VPC.

Configuring Inbound Rules

Click **Add rule** in the "Inbound rules" section.

- **Rule 1: HTTPS for AWS Services**

- **Type:** HTTPS
- **Protocol:** TCP
- **Port range:** 443
- **Source:** Select the security group of the resources (e.g., your EC2 instances or Lambda functions) that need to communicate with these services. This ensures only your application resources can use the endpoints.
- **Description:** Allow inbound HTTPS for AWS service endpoints.

- **Rule 2: Redshift Access**

- **Type:** Redshift
- **Protocol:** TCP
- **Port range:** 5439
- **Source:** Select the security group of the resources that need to connect to Redshift.
- **Description:** Allow inbound traffic to Redshift endpoint.

After adding the rules, click **Create security group**. This new security group can now be attached to your interface endpoints as described in the previous section.

Part 3: Data Engineering & Warehousing

This document details the architecture and implementation of our data engineering and warehousing pipeline. The primary goal is to establish a robust and scalable system for ingesting raw data, transforming it, and loading it into a data warehouse for analytics and business intelligence.

7. Building the Data Lake and Warehouse

The foundation of our analytics platform consists of a data lake built on Amazon S3 and a data warehouse powered by Amazon Redshift Serverless.

7.1. S3 Data Lake Setup

Amazon S3 serves as our central data lake, providing a durable, scalable, and cost-effective storage solution for data at all stages of its lifecycle. We utilize a single S3 bucket, logically partitioned using prefixes to maintain organization and control access.

- **Bucket Name:** leadconversiondata

This bucket is structured with the following prefixes:

- s3://leadconversiondata/raw-data/: This prefix is the landing zone for all incoming, unaltered source data. Data here is considered immutable and serves as the single source of truth.
- s3://leadconversiondata/processed-data/: After data is cleaned, transformed, and enriched by our ETL process, the resulting datasets are stored here. This data is ready for analytics or machine learning.
- s3://leadconversiondata/model-artifacts/: This location is dedicated to storing machine learning assets, such as trained model binaries (.pkl, .h5), performance evaluation plots (e.g., SHAP value plots), and other related files.

Name	Type	Last modified	Size	Storage class
lead_scoring_100.csv	csv	July 18, 2025, 09:53:13 (UTC+05:30)	28.6 KB	Standard

7.2. Provisioning Redshift Serverless

For our data warehouse, we use **Amazon Redshift Serverless**. This option automatically provisions and scales the underlying data warehouse resources, allowing us to focus on data analysis rather than infrastructure management.

Provisioning Steps:

1. Navigate to the **Amazon Redshift** service in the AWS Console.
2. Click **Try Redshift Serverless** to begin the setup.
3. **Configuration:** For the initial setup, you can accept the default settings for the workgroup and namespace.
4. **Network and Security:** This is a critical step.
 - Associate the serverless workgroup with your **VPC**.
 - Attach the `gluesg` security group to allow inbound traffic from AWS Glue.
5. Review the configuration and click **Create**. The serverless endpoint will take several minutes to provision.

The top screenshot shows the 'Serverless dashboard' with the following metrics:

Total snapshots	Dashboards in my account	Dashboards requiring authorization	Dashboards from other accounts	Dashboards requiring association
0	0	0	0	0

The bottom screenshot shows the 'Namespace configuration' page for 'redshift' with the following details:

Namespace	Status	Workgroup	Status	Average query duration	Last hour
redshift	Available				

7.3. Managing Credentials with AWS Secrets Manager

To eliminate hardcoded credentials, a major security risk, we store the Redshift database user credentials in **AWS Secrets Manager**.

Configuration Steps:

1. In the Redshift console, select your newly created serverless instance and set a password for the default admin user, `awsuser`.
2. Navigate to the **AWS Secrets Manager** service.
3. Click **Store a new secret**.
4. **Secret type:** Choose **Credentials for Redshift cluster**.
5. **Credentials:** Enter the username (`awsuser`) and the password you just configured.
6. **Secret name:** Use a descriptive, hierarchical name, such as `prod/redshift/credentials`.
7. **Rotation:** For this guide, you can disable automatic rotation. In a production environment, enabling rotation is highly recommended.
8. Complete the process to store the secret. Take note of the secret's **ARN** for later use.

Secret name	Description	Last retrieved (UTC)
Redshiftsecret	glue owned service linked secret created for Glue ConnectionRedshift connection. Deleting the secret will fail the associated connection.	July 22, 2025
glue!888517279277-Redshiftconnection-1752812385302	Secret created for an AWS SQL Workbench connection	July 20, 2025
sqlworkbench!6d1f02b0-e38c-43e4-b11d-7078ba53ff0d	Secret created for an AWS SQL Workbench connection	July 22, 2025
sqlworkbench!03d57d4f-f430-4a8e-95e6-627fbe99340b	Secret created for an AWS SQL Workbench connection	July 22, 2025
sqlworkbench!3d7e8d9-45fa-4103-8b74-9173e9c5bc46	Secret created for an AWS SQL Workbench connection	July 22, 2025
sqlworkbench!ff0f94dbb-5d20-4033-9973-06d8949844cf	Secret created for an AWS SQL Workbench connection	July 22, 2025

Secret details

Encryption key aws/secretsmanager	Secret description
Secret name Redshiftsecret	
Secret ARN arn:aws:secretsmanager:ap-south-1:888517279277:secret:Redshiftsecret-oLjktj	

Overview **Rotation** **Versions** **Replication** **Tags**

Secret value Info Retrieve secret value

Retrieve and view the secret value.

8. The Ingestion Pipeline: S3 to Redshift

Our data ingestion pipeline uses **AWS Glue** to perform the Extract, Transform, and Load (ETL) process, moving data from our S3 data lake into the Redshift data warehouse.

8.1. AWS Glue for Serverless ETL

AWS Glue is a fully managed, serverless ETL service. The process involves three key components:

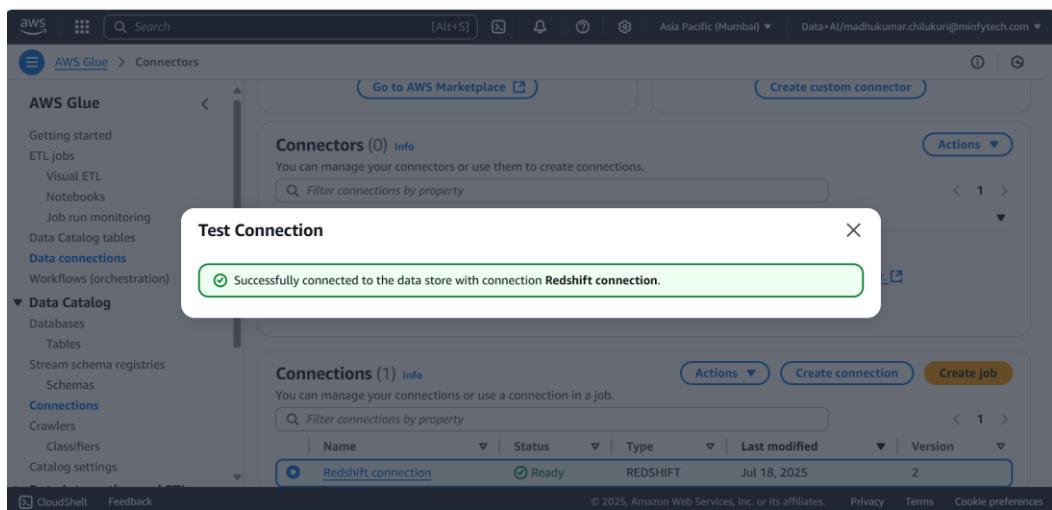
- A **Glue Connection** to securely access the Redshift data warehouse.
- A **Glue Crawler** to automatically discover the schema of our raw data.
- A **Glue Job** to execute the transformation and loading logic.

8.2. Establishing the Glue Connection

This connection enables Glue to communicate with your Redshift Serverless instance without exposing credentials in the ETL script.

Steps:

1. Navigate to the **AWS Glue** service in the console.
2. Select **Connections** from the left-hand menu and click **Create connection**.
3. **Connection details:**
 - **Name:** RedshiftConnection
 - **Connection type:** Amazon Redshift
4. Select **Connect to a Redshift Serverless workgroup** and choose your workgroup from the list.
5. **Authentication:** Choose the AWS Secret (`prod/redshift/credentials`) you created earlier.
6. **Test Connection:** After the connection is created, select it from the list and choose **Actions > Test connection**. Select the `glueaccessrole` IAM role. A "Success" message confirms that your networking (VPC, security groups) and permissions are correctly configured.



8.3. Automating Schema Discovery with Glue Crawlers

The Glue Crawler automatically scans the data in our S3 `raw-data` path, infers the schema, and creates a metadata table in the AWS Glue Data Catalog.

Steps:

1. In the Glue console, select **Crawlers** and click **Create crawler**.
2. **Name:** jobcrawler
3. **Data source:** Point the crawler to the S3 path containing the raw data:
`s3://leadconversiondata/raw-data/`.
4. **IAM Role:** Select the `glueaccessrole`, which grants Glue permissions to access S3 and create Data Catalog entries.
5. **Output:** Configure the crawler to add a new database in the Glue Data Catalog. Name it `lead_conversion_db`.
6. Run the crawler. Once it completes successfully, you will find a new table within the `lead_conversion_db` database, representing the schema of your raw data.

The screenshot shows the AWS Glue console with the 'Gluecrawler' crawler selected. The crawler properties are displayed, including its name ('Gluecrawler'), IAM role ('Glue'), database ('redshiftdatabase'), and state ('READY'). The 'Crawler runs' tab shows one run, with options to 'Stop run', 'View CloudWatch logs', and 'View run details'.

8.4. Designing the Visual ETL Job

We use **Glue Studio**, a visual interface for creating, running, and monitoring ETL jobs.

Steps:

1. In the Glue console, navigate to **Glue Studio** and create a **New Visual ETL job**.
2. **Source Node:**

- Select **AWS Glue Data Catalog** as the source.
- Choose the `lead_conversion_db` database and the table that was automatically created by the `jobcrawler`.

3. **Transform Node:**

- Click the + icon and add a transform node.
- For this example, you can use the **Change Schema** transform to ensure data types are correctly mapped for Redshift (e.g., changing a string to a specific `VARCHAR` length or a number to an `INTEGER`).

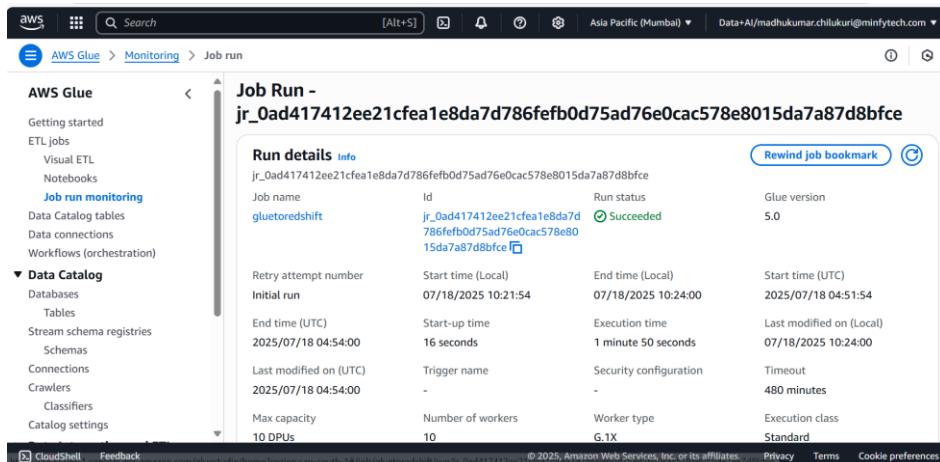
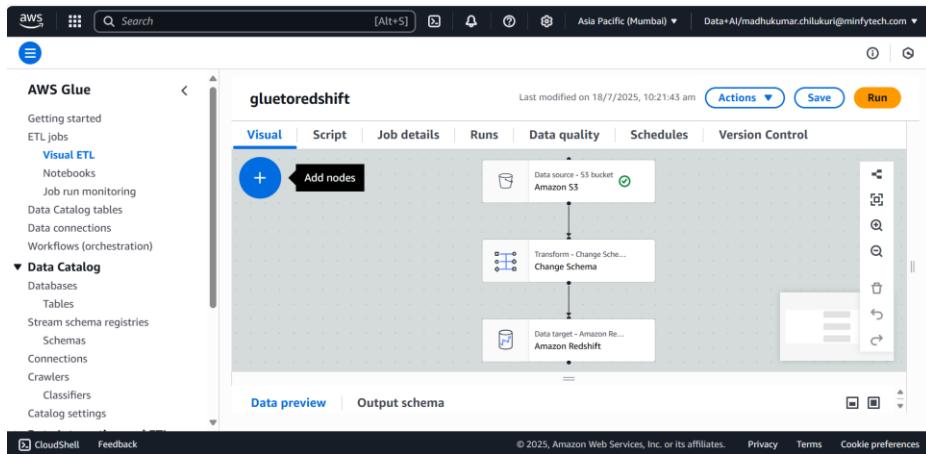
4. **Target Node:**

- Add a target node and select **Amazon Redshift**.
- **Connection:** Choose the `RedshiftConnection` created earlier.
- **Target options:** Specify the destination for the data, such as the `dev` database and a new table named `public.leadstable`.

5. **Job Details:**

- Navigate to the "Job details" tab.
- **Name:** `S3_to_Redshift_ETL`
- **IAM Role:** Assign the `glueaccessrole`.

6. **Save and Run** the job.



Once the job run succeeds, you can connect to your Redshift Serverless instance using a query editor and run `SELECT * FROM public.leadstable;` to verify that the data has been successfully ingested from S3.

9. Data Access from Amazon SageMaker

This section details the procedure for data scientists and machine learning engineers to securely access and load data from the Amazon Redshift data warehouse directly into an Amazon SageMaker notebook environment for analysis and model development.

9.1. Using the Redshift Data API with boto3

To bridge the gap between our data warehouse and our machine learning environment, we use the **Amazon Redshift Data API**. This modern, HTTP-based API provides a simplified and secure method to interact with a Redshift cluster.

Key advantages of this approach include:

- **No Driver Management:** It eliminates the need to install and manage cumbersome JDBC/ODBC drivers and their dependencies within the SageMaker environment.

- **Serverless-Friendly:** It is perfectly suited for interacting with our Redshift Serverless instance.
- **IAM-Based Authentication:** It leverages IAM roles and AWS Secrets Manager for secure, programmatic access without ever exposing database credentials in the code.
- **Asynchronous Execution:** It allows for running long-running queries without blocking the client application, making it efficient for large data-loading tasks.

The following Python script, designed to be run within a SageMaker notebook cell, automates the entire process of querying Redshift and loading the results into a Pandas DataFrame.

9.2. Complete Python Script for Data Loading

Here is the complete script for easy use. A detailed walkthrough of each section follows.

Python

```
import boto3
import pandas as pd
import time

# =====
# Step 1: Configuration
# =====
# Sets all necessary parameters for the connection and query.
# The secret_arn is retrieved from AWS Secrets Manager where Redshift credentials are stored.
region = 'ap-south-1'
workgroup_name = 'default-workgroup' # Or your custom workgroup name
database_name = 'dev'
# IMPORTANT: Replace the ARN with the one for your Redshift secret
secret_arn = 'arn:aws:secretsmanager:ap-south-1:123456789012:secret:prod/redshift/credentials-
xxxxxx'
sql = 'SELECT * FROM public.leadstable' # Using the table created by the Glue ETL job

# =====
# Step 2: Create Redshift Data API Client
# =====
# Initializes the low-level boto3 client to interact with the Redshift Data API.
print("Initializing Redshift Data API client...")
client = boto3.client('redshift-data', region_name=region)
print("Client initialized.")

# =====
# Step 3: Run the SQL Query
# =====
# Submits the SQL query for asynchronous execution on the Redshift Serverless workgroup.
# The API returns a unique statement_id to track the query's execution status.
print(f"Executing SQL query: {sql}")
response = client.execute_statement(
    WorkgroupName=workgroup_name,
    Database=database_name,
    SecretArn=secret_arn,
    Sql=sql
)
```

```

statement_id = response['Id']
print(f"Query submitted with Statement ID: {statement_id}")

# =====
# Step 4: Wait for Query Completion
# =====

# Polls the API to check the query's status. The loop continues until the status
# is 'FINISHED', 'FAILED', or 'ABORTED'.
print("Waiting for query to complete...")
desc = client.describe_statement(Id=statement_id)
status = desc['Status']

while status not in ['FINISHED', 'FAILED', 'ABORTED']:
    time.sleep(1) # Wait for 1 second before checking again
    desc = client.describe_statement(Id=statement_id)
    status = desc['Status']
    print(f"Current status: {status}")

if status == 'FAILED':
    raise Exception(f"Query failed: {desc['Error']}")
elif status == 'ABORTED':
    raise Exception("Query was aborted.")

print("Query finished successfully.")

# =====
# Steps 5 & 6: Retrieve and Parse Query Results
# =====

# Fetches the results of the completed query.
print("Retrieving query results...")
result = client.get_statement_result(Id=statement_id)

# The column names are extracted from the 'ColumnMetadata' field of the result.
columns = [col['name'] for col in result['ColumnMetadata']]
# The data records are in the 'Records' field.
rows = result['Records']

# Each record is a list of dictionaries. We extract the first value from each dictionary.
# A check for `None` is included to handle potential NULL values in the database.
data = []
for row in rows:
    processed_row = [list(col.values())[0] if col else None for col in row]
    data.append(processed_row)
print(f"{len(data)} rows retrieved.")

# =====
# Steps 7 & 8: Convert to Pandas DataFrame and Display
# =====

# The parsed data and column names are used to construct a Pandas DataFrame.
print("Converting to Pandas DataFrame...")
df = pd.DataFrame(data, columns=columns)

# Displays the first 5 rows of the DataFrame to verify successful data loading.
print("Data successfully loaded from Redshift into DataFrame:")

```

```
print(df.head())
```

9.3. Step-by-Step Code Walkthrough

Step 1: Configuration

This initial block sets up all the required parameters. Centralizing these variables makes the script easy to configure and maintain.

Python

```
import boto3
import pandas as pd
import time

# Sets AWS region, Redshift workgroup name, database name, secret ARN, and SQL query.
# The secret_arn is created in AWS Secrets Manager to store Redshift credentials securely.
region = 'ap-south-1'
workgroup_name = 'default-workgroup'
database_name = 'dev'
secret_arn = 'arn:aws:secretsmanager:ap-south-1:123456789012:secret:prod/redshift/credentials-xxxxxx' # Replace with your ARN
sql = 'SELECT * FROM public.leadstable'
```

Step 2: Create Redshift Data API Client

Here, we instantiate a `boto3` client. This object is our gateway to the Redshift Data API, providing methods that map directly to the API's endpoints.

Python

```
# Uses boto3 to create a Redshift Data API client.
# This is how you interact with Redshift without needing a persistent JDBC connection.
client = boto3.client('redshift-data', region_name=region)
```

Step 3: Run the SQL Query

We use the `execute_statement` function to send our SQL query to Redshift. The execution is **asynchronous**, meaning the API call returns immediately with a `statement_id` rather than waiting for the query to finish. This `Id` is the key to tracking our query's progress.

Python

```
# Executes the given SQL query on the Redshift Serverless workgroup.
# Returns a statement_id used to track the query's progress.
response = client.execute_statement(
    WorkgroupName=workgroup_name,
    Database=database_name,
    SecretArn=secret_arn,
    Sql=sql
)
statement_id = response['Id']
```

Step 4: Wait for Query Completion

Because the execution is asynchronous, we must programmatically check if the query is done. This loop calls `describe_statement` every second until the `Status` field changes from `SUBMITTED` or `PICKED` to a terminal state like `FINISHED`, `FAILED`, or `ABORTED`.

Python

```
# Repeatedly checks the query status using the statement_id.  
# Waits in a loop (sleeping for 1 second) until the query either FINISHES or FAILS.  
desc = client.describe_statement(Id=statement_id)  
while desc['Status'] not in ['FINISHED', 'FAILED', 'ABORTED']:  
    time.sleep(1)  
    desc = client.describe_statement(Id=statement_id)
```

Steps 5 & 6: Retrieve and Parse Query Results

Once the query is `FINISHED`, we can retrieve the data using `get_statement_result`. The returned object contains the data in a structured JSON format. We parse this structure to extract the column headers (from `ColumnMetadata`) and the data records (from `Records`).

Python

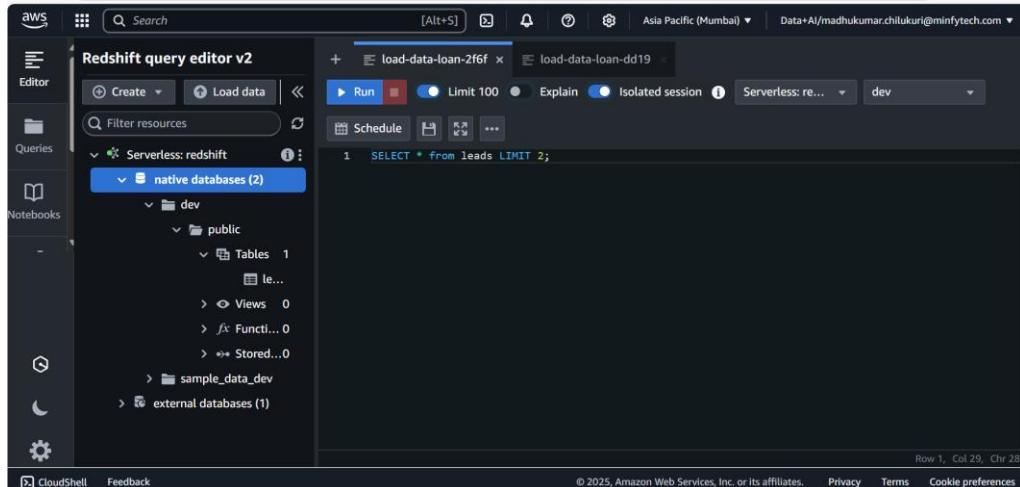
```
# Retrieves the actual data after the query is finished.  
result = client.get_statement_result(Id=statement_id)  
  
# Extracts column names from metadata.  
columns = [col['name'] for col in result['ColumnMetadata']]  
rows = result['Records']  
  
# Iterates over each record and extracts values.  
data = []  
for row in rows:  
    data.append([list(col.values())[0] if col else None for col in row])
```

Steps 7 & 8: Convert to Pandas DataFrame and Display

Finally, the extracted lists of columns and data are used to construct a **Pandas DataFrame**. The DataFrame is the de-facto standard data structure for data manipulation and analysis in Python's data science ecosystem. We print the `.head()` of the DataFrame to quickly verify that the data has been loaded correctly.

Python

```
# Converts the extracted data and columns into a standard Pandas DataFrame.  
df = pd.DataFrame(data, columns=columns)  
  
# Prints the first 5 rows to verify the data was loaded correctly.  
print("Data successfully loaded from Redshift:")  
print(df.head())
```



Part 4: Machine Learning Development

This document outlines the setup and configuration of the machine learning (ML) development environment. All ML activities, from data exploration to model deployment, are centralized within Amazon SageMaker.

10. Amazon SageMaker: The ML Development Environment

Amazon SageMaker is a fully managed AWS service that provides a comprehensive suite of tools to prepare data and to build, train, and deploy high-quality machine learning models efficiently. For this project, we leverage **SageMaker Studio**, a web-based integrated development environment (IDE) that consolidates all the necessary ML development tools.

The core of our setup is a **SageMaker Domain**, which provides a secure, collaborative space for data scientists.

10.1. Creating a SageMaker Domain

A SageMaker Domain is the foundational component for using SageMaker Studio. It organizes users, manages security, and provides shared resources, including an Amazon Elastic File System (EFS) volume for storing notebooks, data, and other artifacts.

Console Walkthrough: Creating the Domain

1. **Navigate to SageMaker:** In the AWS Management Console, use the search bar to find and navigate to the **Amazon SageMaker** service.
2. **Create a Domain:** On the SageMaker dashboard, select **Domain** from the left navigation pane and then click **Create Domain**.
3. **Setup:** Choose the **Standard setup** option for a comprehensive configuration experience.
4. **Authentication:** For this guide, select **AWS Identity and Access Management (IAM)** as the authentication method. This is the standard AWS method for managing access and permissions.
5. **Configure Network and Storage:** This is a critical step to ensure secure and efficient communication between SageMaker and other AWS resources.

- **VPC:** Select the **Virtual Private Cloud (VPC)** that houses your Redshift cluster and VPC endpoints. Placing SageMaker in the same VPC is essential for secure, low-latency communication that does not traverse the public internet.
 - **Subnet(s):** Choose the **private subnets** associated with your VPC. Selecting multiple subnets across different Availability Zones provides high availability for the SageMaker environment.
 - **Security Group:** Attach the `mlops-sg` security group. This firewall policy is configured to allow SageMaker to communicate securely with services like Redshift and an MLflow server through the private VPC endpoints established earlier.
6. **IAM Role:** The SageMaker execution role grants the necessary permissions to perform tasks on your behalf. Create a new IAM role or use an existing one with the following policies attached:
- `AmazonSageMakerFullAccess`: This AWS managed policy provides comprehensive permissions for SageMaker operations.
 - **S3 Bucket Policy:** A custom inline policy granting `s3:GetObject`, `s3:PutObject`, and `s3>ListBucket` permissions on the `leadconversiondata` bucket.
 - **Redshift Data API Policy:** A custom policy allowing the `redshift-data:ExecuteStatement` and `redshift-data:DescribeStatement` actions, enabling notebooks to query the data warehouse.
 - **Secrets Manager Policy:** A policy allowing access to the specific Redshift credentials secret stored in AWS Secrets Manager.
7. **Create User Profile:** Within the domain setup, add a default user profile. Assign it a descriptive name like `ml-developer`.
8. **Review and Submit:** Carefully review all configuration settings and click **Submit**. The domain creation process will take several minutes to complete.

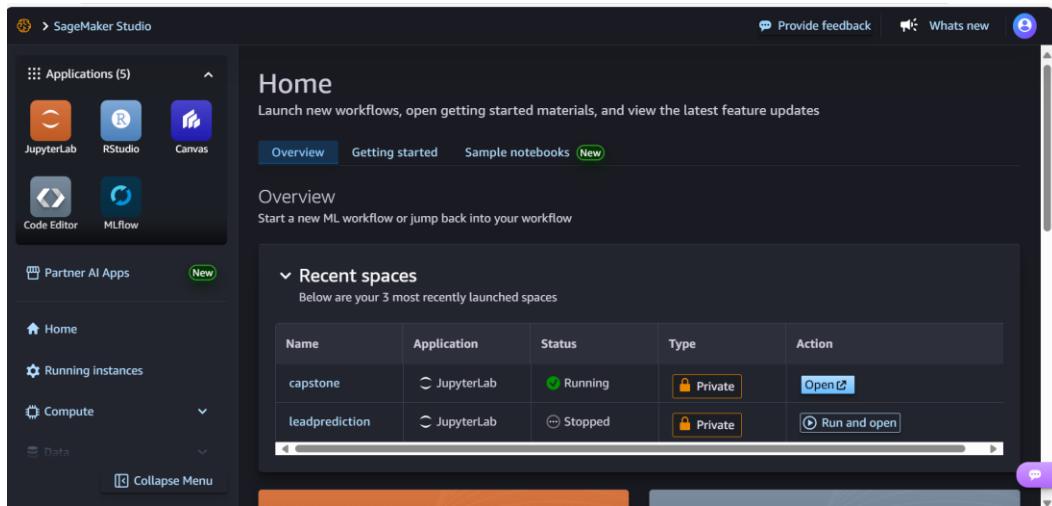
The screenshots show the 'User profiles' section of the 'Domain details' page for 'QuickSetupDomain-20250718T103832'. The top screenshot shows a single user profile named 'default-' with details: Jul 18, 2025 09:00 UTC and Jul 18, 2025 05:12 UTC. The bottom screenshot shows the same user profile with a context menu open over it, revealing options like 'Personal apps', 'Studio', 'Canvas', 'Collaborative', and 'Spaces'.

10.2. Launching SageMaker Studio

Once the Domain's status changes to "**Ready**," you can launch the SageMaker Studio IDE.

1. From the SageMaker **Domain** page, locate the `ml-developer` user profile in the user list.
2. On the right side of the user row, click the **Launch** dropdown menu and select **Studio**.

A new browser tab will open, and the SageMaker Studio environment will initialize. The first launch may take a few minutes as it provisions the underlying compute resources and attaches the EFS volume.



10.3. Opening the JupyterLab Notebook

SageMaker Studio provides a fully managed **JupyterLab** environment, which is the primary interface for all coding tasks.

1. **Open Launcher:** Once Studio has loaded, you will be presented with the main launcher screen. If this screen is not visible, you can open a new one by navigating to **File > New Launcher** in the top menu.
2. **Create Notebook:** Under the **Notebooks and compute resources** section of the launcher, click the **Notebook** button.
3. **Select Kernel:** A dialog box will appear, prompting you to select an image and kernel.
 - o **Image:** Choose **Data Science**.
 - o **Kernel:** Select **Python 3**.

```

# --- Configuration ---
region = 'ap-south-1' # change to your region
workgroup_name = 'redshift'
database_name = 'dev'
secret_arn = 'arn:aws:secretsmanager:ap-south-1:888517279277:secret:Redshiftsecret-oLjktj'
sql = 'SELECT * FROM leads LIMIT 100'

# --- Create boto3 Redshift Data API client ---
client = boto3.client('redshift-data', region_name=region)

# --- Execute query ---
response = client.execute_statement(
    WorkgroupName='redshift',
    Database='dev',
    SecretArn='arn:aws:secretsmanager:ap-south-1:888517279277:secret:Redshiftsecret-oLjktj',
    Sql='sql'
)

statement_id = response['Id']

# --- Wait for completion ---
desc = client.describe_statement(Id=statement_id)
while desc['Status'] not in ['FINISHED', 'FAILED', 'ABORTED']:
    time.sleep(1)
    desc = client.describe_statement(Id=statement_id)

```

This selection provisions a Python environment that comes pre-installed with essential data science libraries, including **pandas**, **NumPy**, **scikit-learn**, and **Matplotlib**, saving significant setup time.

A new notebook file (e.g., `Untitled.ipynb`) will open in a new tab. You are now in a fully configured development environment, ready to proceed with the machine learning workflow.

11. Exploratory Data Analysis (EDA) & Cleaning

Before any feature transformation, a thorough exploratory analysis was conducted to understand the dataset's characteristics and identify potential data quality issues.

11.1. Initial Data Inspection

The initial assessment involved several key steps to get a comprehensive overview of the data:

- `df.info()`: This revealed that the dataset contains 9240 entries and 37 columns, with a mix of numerical (`int64`, `float64`) and categorical (`object`) data types. It immediately highlighted columns with significant numbers of missing values, such as `Country`, `Specialization`, `Tags`, and `Lead Quality`.
- `df.describe()`: This provided statistical summaries for the numerical columns. It showed a wide range in values for features like `TotalVisits` and `Total Time Spent on Website`, indicating that scaling would be necessary.
- `df.isnull().sum()`: This gave a precise count of null values per column, confirming the findings from `df.info()` and providing a clear list of columns requiring imputation strategies.
- **Target Distribution**: An analysis of the `Converted` column showed a class distribution of approximately 61.5% not converted (0) and 38.5% converted (1), indicating a moderately imbalanced dataset.

11.2. Data Cleaning and Standardization

Based on the initial inspection, the following cleaning steps were performed:

- **Dropping Non-Predictive Columns:** Several columns were identified as irrelevant for predicting lead conversion and were programmatically dropped. This reduces noise and simplifies the model.
 - **Identifier Columns:** Prospect ID, Lead Number.
 - **Zero/Low Variance Columns:** Get updates on DM Content, Receive More Updates About Our Courses, I agree to pay the amount through cheque, Magazine, Update me on Supply Chain Content. These columns had the same value for nearly every row and thus offered no predictive power.
- **Standardizing Missing Value Representation:** The notebook handles placeholder strings like "Select" not by a global replacement, but within the custom mapping functions in the feature engineering step. This targeted approach ensures that such values are consistently mapped to a standardized "Unknown" category before being imputed.

12. Feature Engineering

Feature engineering is where raw data is transformed into features that better represent the underlying problem for the machine learning model. Our strategy focuses on custom mapping and consolidation to reduce complexity and embed domain knowledge.

12.1. Strategy: Custom Mapping and Transformation

We utilize Scikit-learn's `FunctionTransformer` to integrate custom Python functions directly into our main preprocessing pipeline. This allows for complex, reusable logic while maintaining a clean, end-to-end workflow.

The core of our feature engineering is the `map_categorical_columns` function, which consolidates high-cardinality categorical features into fewer, more meaningful groups.

12.2. Consolidation of Categorical Features

- **Geographic Grouping (`city, country`):** Raw location data is grouped into broader regions to capture market types rather than noisy individual locations. For example, Mumbai and Thane & Outskirts are mapped to Metro India, while various Middle Eastern countries are grouped into Middle East.
- **Occupation and Specialization Grouping:** What is your current occupation is simplified into categories like Student, Working, and Non-Working. The many Specialization values are consolidated into broader business fields like Marketing, Finance, and IT.
- **Behavioral & Intent Categorization (`tags`):** The highly variable `Tags` column, which contains free-text-like notes from sales agents, is mapped to a standardized set of intents. For instance, Will revert after reading the email becomes Pending Response, and switched off is grouped under Trying to Contact. This transforms noisy operational data into a clean feature representing the lead's status.

- **Lead Profile & Quality Simplification:** Lead Profile and Lead Quality values are mapped to standardized, more interpretable categories like High, Medium, Low.

12.3. Skewness Correction with Yeo-Johnson Transformation

Initial analysis revealed that several numerical features were highly skewed. To address this, the **Yeo-Johnson transformation** was applied. This is a power transformation technique that is highly effective at stabilizing variance and making the distribution of the data more normal.

The primary advantage of using Yeo-Johnson over a simpler method like a log transform is its **flexibility**. It can handle a wider range of data, including **zero and negative values**, which are present in our dataset (e.g., `TotalVisits` can be 0). This makes it a more robust and reliable choice for a production-grade preprocessing pipeline, ensuring the transformation works correctly without requiring special handling for non-positive numbers. This step is crucial for improving the performance of many machine learning algorithms that perform better with normally distributed data.

13. The Preprocessing Pipeline

To ensure all transformations are applied consistently and to prevent data leakage, we construct a unified preprocessing pipeline using Scikit-learn's `Pipeline` and `ColumnTransformer` objects.

13.1. Identifying Feature Types

The features are first programmatically segregated into three distinct groups for tailored processing:

- `numeric_cols`: All columns with `int64` or `float64` data types.
- `cat_ohe_cols`: Standard categorical features that will be one-hot encoded.
- `label_cols`: Ordinal features (`Lead Quality`, `Asymmetrique Activity Index`, etc.) that have an inherent order and will be label encoded.

13.2. Building the Scikit-learn Pipeline

The `build_full_pipeline` function assembles the complete preprocessing logic.

1. **Custom Transformations First:** The pipeline begins with the custom `FunctionTransformer` steps to ensure raw data is cleaned and engineered before any standard imputation or scaling:
 - `mapping_transformer`: Applies the categorical consolidation logic.
 - `replace_unknowns_transformer`: Converts any "Unknown" strings to `np.nan` so they can be imputed.
 - **Skewness Transformation:** The function applying the Yeo-Johnson transformation is included here.

2. **Column-Specific Sub-Pipelines:** A `ColumnTransformer` then applies specific pipelines to each feature group:

- o **Numeric Pipeline:**

1. `SimpleImputer(strategy='median')`: Fills missing numerical values with the column's median.
2. `MinMaxScaler()`: Scales features to a [0, 1] range.

- o **Categorical (Nominal) Pipeline:**

1. `SimpleImputer(strategy='most_frequent')`: Fills missing categorical values with the mode.
2. `OneHotEncoder(handle_unknown='ignore')`: Converts categories into binary columns, ignoring any new categories that might appear during inference.

- o **Categorical (Ordinal) Pipeline:**

1. `SimpleImputer(strategy='most_frequent')`: Fills missing ordinal values.
2. `OrdinalEncoder()`: Maps ordered categories to integer values (e.g., 'Low' -> 0, 'Medium' -> 1, 'High' -> 2).

This unified pipeline object encapsulates the entire data preparation workflow.

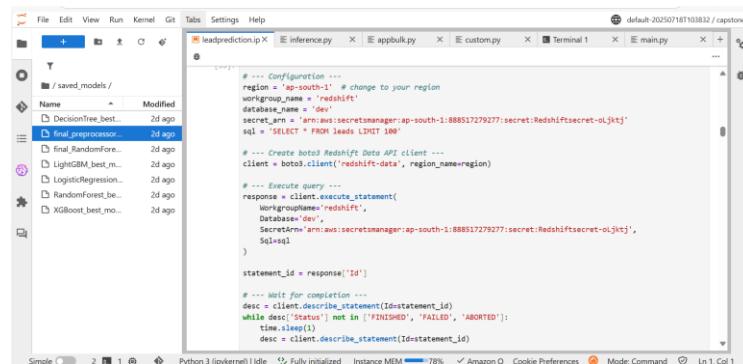
13.3. Saving the Preprocessing Pipeline

As part of the final model registration process in the `save_and_register_best_model_pipeline` function, the final, fitted preprocessor object is saved to disk using `joblib`.

Python

```
# The preprocessor is saved as part of the main pipeline and also separately
preprocessor_path = os.path.join(save_dir, "final_preprocessor.pkl")
joblib.dump(preprocessor, preprocessor_path)
print(f"✅ Preprocessing pipeline saved at: {preprocessor_path}")
```

This serialized object is crucial for production, as it allows the exact same transformations to be applied to new, incoming data before making predictions, ensuring complete consistency between the training and inference environments.



```
# ... Configuration ...
region = 'ap-south-1' # change to your region
workgroup_name = 'redshift'
database_name = 'dev'
secret_arn = 'arn:aws:secretsmanager:ap-south-1:888517279277:secret:Redshiftsecret-oLjkjt'
sql = 'SELECT * FROM leads LIMIT 100'

# ... Create boto3 Redshift Data API client ...
client = boto3.client('redshift-data', region_name=region)

# ... Execute query ...
response = client.execute_statement(
    WorkgroupName='redshift',
    Database='dev',
    SecretArn=secret_arn,
    Sql=sql
)

statement_id = response['Id']

# ... Wait for completion ...
desc = client.describe_statement(Id=statement_id)
while desc['Status'] not in ['FINISHED', 'FAILED', 'ABORTED']:
    time.sleep(1)
    desc = client.describe_statement(Id=statement_id)
```

14. Model Training and Evaluation

This section outlines the complete methodology used for training, fine-tuning, and evaluating a suite of machine learning models to predict lead conversion. It details our specific strategy for handling class imbalance, the models tested, and the robust framework used for evaluation and explainability.

14.1. Strategy for Handling Class Imbalance

Our dataset exhibits a moderate class imbalance, where "non-converted" leads significantly outnumber "converted" ones. A naive model trained on such data would likely become biased towards the majority class, resulting in poor predictive performance on the minority class we care about most.

To counter this, we employ a sophisticated strategy that focuses on **model-level weighting** combined with a **robust, stratified validation structure**, rather than data-level resampling techniques like SMOTE.

Why Not SMOTE? The Case for Model-Level Weighting

While data resampling techniques like the **Synthetic Minority Over-sampling Technique (SMOTE)** are a common approach, they are not always the optimal choice, especially when working with powerful tree-based models. We deliberately chose *not* to use SMOTE for the following key reasons:

1. **Inherent Robustness of Tree-Based Models:** The core of our model suite consists of tree-based ensembles like Random Forest, XGBoost, and LightGBM. Unlike distance-based algorithms (like SVM) or some linear models, tree models make decisions by finding optimal split points based on criteria like Gini impurity or information gain. They can effectively learn to identify the minority class by creating specific decision rules, making them naturally more resilient to class imbalance.
2. **Risk of Introducing Noise:** SMOTE works by creating new, synthetic data points for the minority class. While this balances the dataset, it can also introduce noise and create artificial samples in areas of the feature space that are not representative of real-world data, potentially blurring the decision boundary and leading to less generalizable models.
3. **Superiority of Built-in Weighting Mechanisms:** Modern machine learning libraries provide a more direct and often more effective way to handle imbalance: **class weighting**. Instead of changing the data, we change the model's learning process.
 - For scikit-learn models like Random Forest and Decision Tree, we use the `class_weight='balanced'` parameter. This automatically adjusts the weights in the loss function to be inversely proportional to class frequencies, meaning the model is penalized more heavily for misclassifying a minority class sample.
 - For XGBoost and LightGBM, we use the `scale_pos_weight` parameter. This is a specific hyperparameter designed for binary classification that scales the gradient for the positive (minority) class, effectively telling the model to "pay more attention" to getting those predictions right.

This model-level approach forces the algorithm to focus on the minority class without artificially altering the underlying data distribution, leading to a more robust and reliable model.

Robust Evaluation with Stratified K-Fold Cross-Validation

To ensure our performance metrics are trustworthy and reflect real-world performance, we use **Stratified K-Fold Cross-Validation**. This technique splits the data into 'k' folds while preserving the original percentage of samples for each class in every fold. This guarantees that each training and validation split is a representative sample of the overall data distribution, preventing scenarios where a fold might accidentally contain very few minority class samples.

By combining model-level class weighting with stratified validation, we train our models to be sensitive to the crucial "converted" leads while evaluating them against a realistic, imbalanced benchmark.

14.2. Models Utilized

We evaluated a suite of powerful and well-established classification models to identify the best performer for our lead conversion problem:

- **Logistic Regression:** A strong, highly interpretable linear model that serves as a solid baseline.
- **Decision Tree:** Excellent for visualizing decision rules and capturing non-linear relationships.
- **Random Forest:** An ensemble of decision trees that reduces overfitting and improves robustness by averaging predictions.
- **XGBoost:** A high-performance, gradient-boosted decision tree implementation known for its state-of-the-art accuracy.
- **LightGBM:** A faster, more memory-efficient gradient boosting framework that excels on large, structured datasets.

14.3. The `train_log_and_shap_classification` Function: A Detailed Look

A single, comprehensive function, `train_log_and_shap_classification`, orchestrates the entire training, evaluation, and logging process, ensuring consistency and reproducibility.

Key Steps:

1. **Setup:** Initializes the MLflow tracking URI and experiment name, and creates local directories for saving model artifacts and SHAP plots.
2. **Model Loop:** Iterates through each of the models defined in the suite.

3. **Hyperparameter Tuning:** For each model, it performs hyperparameter tuning using `GridSearchCV` with 5-fold stratified cross-validation, optimizing for the `f1-score`.
4. **Evaluation:** The best estimator found by the grid search is used to make predictions on a held-out validation set. Our custom evaluation function calculates a full suite of performance metrics.
5. **MLflow Logging:** A new MLflow run is created for each model, and the following are meticulously logged:
 - The model's **best hyperparameters**.
 - All calculated **performance metrics** (F1, Precision, Recall, etc.).
 - The **best estimator object** itself, using `mlflow.sklearn.log_model` for easy reloading.
 - The **SHAP summary plot** as a `.png` artifact for model explainability.
6. **SHAP Explanations:** A SHAP explainer is generated to calculate feature importances. The results are visualized in a summary plot that is both saved locally and logged to MLflow.
7. **Results Aggregation:** All results are collected into a single Pandas DataFrame, allowing for easy, side-by-side comparison to identify the champion model.

14.4. Hyperparameter Tuning with GridSearchCV

`GridSearchCV` is the engine behind our model optimization. It automates the demanding process of finding the best model configuration. By systematically working through predefined combinations of model parameters (a "grid") and evaluating each combination using cross-validation, it identifies the set of hyperparameters that yields the best performance on our validation folds, ensuring we use a highly optimized version of each model for final comparison.

14.5. Model Evaluation Framework

Our `evaluate_classification_metrics` function provides a standardized and comprehensive report for each model, containing the most important classification metrics for our business problem:

- **Accuracy:** The overall percentage of correct predictions. While useful, it can be misleading on imbalanced datasets.
- **Precision:** Of all the leads the model predicted would convert, how many *actually* did? This is critical for minimizing the cost of engaging non-serious leads (false positives).
- **Recall:** Of all the leads that *actually converted*, how many did our model successfully identify? This is crucial for maximizing the number of captured opportunities (minimizing false negatives).
- **F1-Score:** The harmonic mean of Precision and Recall. It provides a single, balanced score that is especially useful for imbalanced classification, as it requires both high precision and high recall to be high. This was our primary metric for model selection.

- **ROC-AUC:** Measures the model's ability to distinguish between the positive and negative classes across all possible probability thresholds, giving a complete picture of its separability power.

15. Model Explainability with SHAP

To build trust in our model and move beyond just performance metrics, it is crucial to understand *why* it makes certain predictions. For this, we use the **SHAP (SHapley Additive exPlanations)** library, a state-of-the-art approach to explaining the output of any machine learning model.

15.1. Generating and Interpreting SHAP Values

For each model trained in our pipeline, a SHAP summary plot is automatically generated and logged.

Process: A SHAP `Explainer` is chosen based on the model type (e.g., `TreeExplainer` for Random Forest and XGBoost, which is highly optimized for tree-based models). This explainer computes SHAP values for every feature for each prediction made on our validation set. These values represent the contribution of each feature to pushing the model's output away from the base value.

The SHAP Summary Plot: This powerful visualization aggregates the SHAP values into a single plot that provides both global feature importance and local feature effects.

- **Feature Importance:** Features are ranked vertically from most important (top) to least important (bottom) based on their mean absolute SHAP value.
- **Impact on Prediction:** The horizontal axis represents the SHAP value.
 - **Positive SHAP values** indicate a feature's contribution towards a positive prediction (Converted = 1).
 - **Negative SHAP values** indicate a contribution towards a negative prediction (Not Converted = 0).
- **Original Feature Value:** The color of each point represents the original value of that feature for that observation (e.g., red for high values, blue for low values).

For example, if "Total Time Spent on Website" shows a wide spread of red dots on the right (positive SHAP value), it means that higher values for this feature strongly push the model to predict that a lead will convert.

Logging: This summary plot is saved as a `.png` file and logged as an artifact to the corresponding MLflow run. This creates a permanent, visual record of feature importance that can be reviewed for any model version at any time.

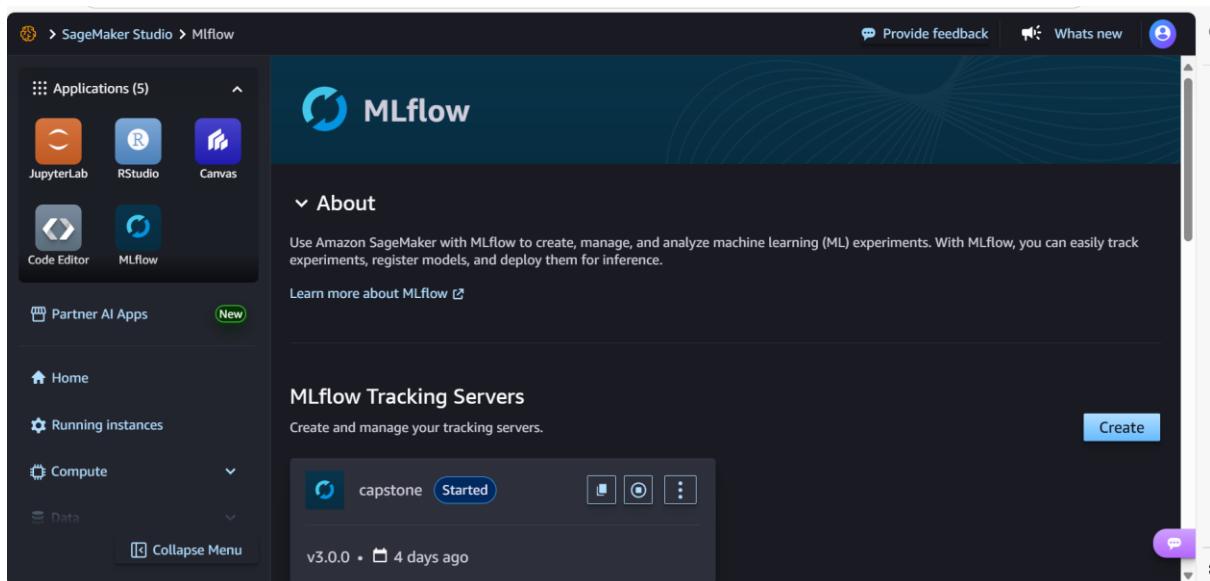
16. Experiment Tracking with MLflow

To manage the complexity of the machine learning experimentation process, we use **MLflow Tracking**. This component provides a centralized, organized repository for all training runs.

16.1. Setting Up the Centralized MLflow Tracking Server

To enable collaboration and ensure the persistence of our results, we deploy an **MLflow tracking server on a dedicated EC2 instance** within our VPC. This approach is superior to logging to local files because it:

- **Centralizes Data:** All team members can log to and view experiments from a single source of truth.
- **Ensures Durability:** Experiment data is safely stored on a remote server, independent of any developer's local machine.
- **Enhances Security:** The server is only accessible via its private IP within the VPC, and its security group is configured to allow inbound traffic only on port 5000 from trusted sources.



16.2. Integrating MLflow into the Training Pipeline

MLflow is seamlessly woven into our `train_log_and_shap_classification` function to automatically capture all relevant information for every model trained.

- `mlflow.set_tracking_uri()`: Points the MLflow client in our training script to the private IP address of our EC2-hosted tracking server.
- `mlflow.set_experiment()`: Organizes all runs under a specific experiment name (e.g., "LeadScoring_Simplified"), preventing clutter.
- `with mlflow.start_run()`: Creates a new run context for each model. MLflow automatically captures metadata like the source code version and run duration.
- `mlflow.log_param()` & `mlflow.log_metric()`: Explicitly logs the best hyperparameters found by GridSearchCV and the final evaluation scores (F1, ROC-AUC, etc.) on the validation set.
- `mlflow.sklearn.log_model()`: This is a critical step. It packages the trained model object, its library dependencies, and a `conda.yaml` file. This creates a fully reproducible model artifact that can be reliably loaded in different environments.
- `mlflow.log_artifact()`: Logs supplementary files, most importantly the SHAP summary plot, providing rich context for each run.

17. Model Registration and Lifecycle Management

After experimentation, the best model must be formally promoted for production use. This is managed by the **MLflow Model Registry**.

17.1. Saving and Registering the Champion Model

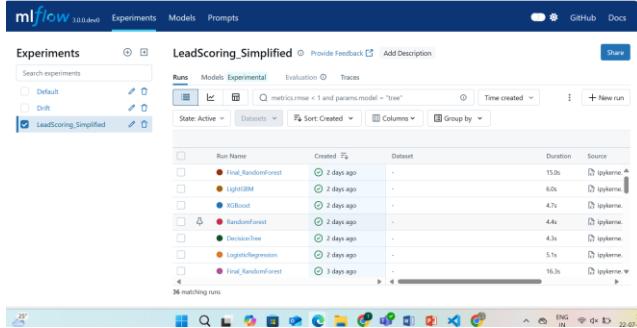
Our pipeline automatically identifies the "champion" model based on the highest F1-score. This model then undergoes a final preparation process:

1. **Retrain on All Data:** The champion model is retrained on the *entire* dataset (training + validation sets combined) to learn from all available information, maximizing its predictive power.
2. **Build Full Pipeline:** The retrained model is combined with the fitted preprocessing pipeline into a **single, end-to-end scikit-learn pipeline object**. This is a crucial best practice. This final artifact can take raw data as input and produce a prediction, encapsulating all necessary steps (imputation, scaling, encoding, and prediction) and preventing training-serving skew.
3. **Log and Register:** This complete pipeline is logged to a final MLflow run and then **registered** in the MLflow Model Registry under a consistent name (e.g., `lead-conversion-classifier`).

17.2. The MLflow Model Registry

The Model Registry acts as a central hub and single source of truth for managing the lifecycle of our production models. It allows us to:

- **Version Models:** Every time we register a model with the same name, a new, auto-incrementing version is created (v1, v2, etc.). This provides a complete lineage of all production candidates.
- **Manage Stages:** We use stages to govern the model's lifecycle. Key stages include:
 - **Staging:** For models that have passed initial tests and are ready for pre-production validation or integration testing.
 - **Production:** For the officially approved model version that is currently serving live traffic.
 - **Archived:** For retired model versions.
- **Transition Models:** Our pipeline programmatically transitions the newly registered model version to the "Staging" or "Production" stage. This controlled promotion is a key step in enabling automated CI/CD for machine learning.



18. Serving and Inference

The final step is to make our trained model available to accept new data and provide predictions.

18.1. Loading the Production Model from the Registry

To ensure our application always uses the correct, approved model, we do not hardcode model paths. Instead, we have a function that **dynamically discovers and loads the latest "Production" version** from the MLflow Model Registry.

- `get_latest_production_model()`: This helper function connects to the MLflow client, queries the registry for the model named `lead-conversion-classifier`, and retrieves the URI for the latest version in the "Production" stage.
- `mlflow.sklearn.load_model()`: This powerful function takes the model URI (e.g., `models:/lead-conversion-classifier/Production`) and loads the complete, deployment-ready pipeline object into memory.

This approach decouples our application from the model training cycle, allowing data scientists to update the production model without requiring any code changes in the application itself.

```

C2025-07-20 06:56:50,999 - INFO - t=2025-07-20T06:56:00+0000 lvl=info msg="received stop request" obj@app stopReq={err:<nil> res:<nil> tag:<nil>}
2025-07-20 06:56:51,000 - INFO - t=2025-07-20T06:56:51+0000 lvl=info msg="session closing" obj@tunnels.session err=nil
sagemaker-user@default: ~ python appbulk.py
2025-07-20 07:12:31,045 - INFO - [x] will load model: RandomForest (version 5)
2025-07-20 07:12:31,046 - INFO - [x] Loading model from URI: models:/RandomForest/Production
Downloading artifacts: 100% |██████████| 5 / [00:00<00:00, 42.58it/s]
2025-07-20 07:12:33,713 - INFO - [x] Model loaded successfully from MLflow
2025-07-20 07:12:33,715 - INFO - Updating auth token for default "config_path" of "ngrok_path": /home/sagemaker-user/.config/ngrok/ngrok
2025-07-20 07:12:33,727 - INFO - Opening tunnel named: http-8080-0e13714f-d7a9-4d8e-b300-1c5bde2e2f5a4
2025-07-20 07:12:33,730 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="no configuration paths supplied"
2025-07-20 07:12:33,740 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="using configuration at default config path" path=/home/sagemaker-user/.config/ngrok/ngrok.yml
2025-07-20 07:12:33,741 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="open config file" path=/home/sagemaker-user/.config/ngrok/ngrok.yml err=nil
2025-07-20 07:12:33,763 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="starting web service" obj=web addr=127.0.0.1:4040 allow_hosts=[]
2025-07-20 07:12:34,421 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="client session established" obj=tunnels.session
2025-07-20 07:12:34,421 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="tunnel session started" obj=tunnels.session
2025-07-20 07:12:34,433 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="start pg/api/tunnels id=9528d2ff61e9ad52"
2025-07-20 07:12:34,433 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="end pg/api/tunnels id=9528d2ff61e9ad52 status=200 dur=26 ms"
2025-07-20 07:12:34,434 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="start pg/api/tunnels id=5e11de4ee2142e47"
2025-07-20 07:12:34,434 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="end pg/api/tunnels id=5e11de4ee2142e47 status=200 dur=95 ms"

```

```

2025-07-20 07:12:33,741 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="open config file" path=/home/sagemaker-user/.config/ngrok/ngrok.yml err=null
2025-07-20 07:12:33,763 - INFO - t=2025-07-20T07:12:33+0000 lvl=info msg="starting web service" obj=web addr=127.0.0.1:4040 allow_
hosts:[]
2025-07-20 07:12:34,421 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="client session established" obj=tunnel.session
2025-07-20 07:12:34,421 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="tunnel session started" obj=tunnel.session
2025-07-20 07:12:34,433 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="start pg/api/tunnels id=9528d2ff61e9ad52"
2025-07-20 07:12:34,433 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="end pg/api/tunnels id=9528d2ff61e9ad52 status=200 dur=26
6.164us
2025-07-20 07:12:34,435 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="start pg/api/tunnels id=5e11de4ee2142e47"
2025-07-20 07:12:34,435 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="end pg/api/tunnels id=5e11de4ee2142e47 status=200 dur=95
837us
2025-07-20 07:12:34,435 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="start pg/api/tunnels id=83e9baaf2d12e7c"
2025-07-20 07:12:34,435 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="end pg/api/tunnels id=83e9baaf2d12e7c status=200 dur=78
753us
2025-07-20 07:12:34,658 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="started tunnel" obj=tunnels name=http-8080-e13714f-d7a9
-4853-3534-3534-3534
2025-07-20 07:12:34,658 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="Public URL: https://85f926fbab9.ngrok-free.app"
2025-07-20 07:12:34,665 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="pg/api/tunnels id=f4388acb8b63f59e status=201 dur=22
1.853ms
2025-07-20 07:12:34,666 - INFO - t=2025-07-20T07:12:34+0000 lvl=info msg="WARNING: This is a development server. Do not use it in a production deployment. Use a production
NGINX server instead."
* Running on all addresses (0.0.0.0)

```

18.2. Real-Time Prediction with Flask and Ngrok

For demonstration and testing, we wrap our prediction logic in a lightweight **Flask** web application.

Flask App (`app.py`):

- Initialization:** On startup, the Flask app initializes and immediately calls the functions to load the latest production model from the MLflow Registry. This ensures the model is loaded into memory once, not on every request, for efficiency.
- /predict Route:** It defines a `/predict` API endpoint that accepts POST requests containing new lead data in a JSON format.
- Prediction Logic:** The endpoint takes the incoming JSON data, converts it into a DataFrame, passes it through the loaded pipeline's `.predict()` method, and returns the final prediction (0 or 1) (Convert or NotConvert)in a clean JSON response.

Ngrok Tunnel for Testing: To make the locally running Flask app accessible for end-to-end testing, we use **Ngrok**.

1. The Flask app is run locally: `python app.py`.
2. In a separate terminal, Ngrok is started: `ngrok http 5000`.

Ngrok creates a secure public URL that tunnels requests directly to our local Flask application. This allows us to test the entire inference pipeline, from sending a request over the internet to receiving a prediction from our production-ready model, without needing to deploy to a formal cloud environment.

The screenshot shows a web browser window at address 127.0.0.1:5001. The title bar reads "Lead Conversion Predictor". Below it is a form with a central area for file upload labeled "Click to upload or drag and drop CSV file with lead data". A blue button at the bottom of the form says "Predict Conversions". Above the form, a small note says "Upload a CSV file with lead data to predict which customers are likely to convert." The browser's status bar at the bottom right shows the date as 20-07-2025 and the time as 09:09.

Prediction Results

Analyzed 100 leads. Predicted 29 will convert.

PROSPECT ID	LEAD SOURCE	TOTAL TIME SPENT ON WEBSITE	PREDICTION
792702df-8bb0-4d29-09a2-bde0beafe620	Clark Chat	0	Not Converted
2a072236-5132-4136-86fe-dcc88c589482	Organic Search	674	Not Converted
8cc6c611-4219-4f35-ac23-fd7f2650bd8a	Direct Traffic	1532	Converted
0cc2df48-7c14-4e39-9de9-197979b38bcc	Direct Traffic	305	Not Converted
3256f528-e534-4826-9b53-4a8b8782852	Google	1428	Not Converted
2058ef08-2856-443e-a0f1-a9237b215ce	Clark Chat	0	Not Converted
9fae70f0-169d-48fb-af64-03d752542ed	Google	1640	Converted
20ef72a2-f3b3-45e0-924e-551c5fa59985	Clark Chat	0	Not Converted
cfa0128c-a0da-4658-9047-0aa4e7bf690	Direct Traffic	71	Not Converted
a4485d64-7204-4130-9e05-33231863c6b5	Google	58	Not Converted
2a393e35-0395-4ca9-9e4f-9d2775aa320	Organic Search	1351	Converted
9bc8cc93-0144-49b0-99b2-0a06c880fb3c	Direct Traffic	1343	Converted
8bf78a52-2478-476b-8818-1688e07874ad	Organic Search	1538	Converted
888870f7-3707-4753-a633-1c7d1d0535e6	Organic Search	170	Not Converted
a8537c22-fc11-48ff-a711-fb5a0f98a9f7	Direct Traffic	481	Not Converted
2594ac14-ff4b-4c22-9c61-b44e85e9198f	Organic Search	1012	Not Converted
3ab27c77-1634-4083-9a0f-861068220811	Clark Chat	0	Not Converted
e5c30eca-a066-4b3f-8c01-0919fbca3f2	Referral Sites	973	Not Converted
82cb16f0-2d97-4a39-a630-ab0fe2e7718c	Google	1688	Converted

19. Post-Deployment: Drift Detection

A machine learning model's performance can degrade silently over time if the live production data it encounters begins to differ from the data it was trained on. This phenomenon is known as **data drift**. We implement a proactive monitoring strategy to detect this drift early and ensure our model remains reliable.

19.1. Monitoring for Data Drift with Evidently AI

We use the **Evidently AI** library, a powerful open-source tool specializing in comprehensive data drift and model performance analysis, to monitor our model's health in production.

- Tool:** **Evidently AI** is chosen for its ability to generate detailed, easy-to-understand reports that compare datasets and highlight statistical differences.
- Process:** The core of our monitoring strategy involves periodically running an automated script. This script takes our original, stable training data as a **reference dataset** and compares it against a **current dataset**, which is a sample of recent, live production data that the model has been scoring.

- **Report Generation:** We use Evidently's `DataDriftPreset`. This powerful feature automatically generates a detailed report that compares the statistical distributions of every feature between the reference and current datasets. It checks for changes in mean, standard deviation, and value distributions for categorical features, providing a complete picture of any data drift.

19.2. Automating Drift Analysis and Logging

This entire drift analysis process is automated and tightly integrated with our centralized **MLflow Tracking Server** for versioning, review, and alerting.

- **Start MLflow Run:** A new MLflow run is initiated specifically to store the drift analysis results. This keeps them cleanly separated from our model training experiments, creating a clear audit trail for model monitoring.
- **Generate Report:** The Evidently drift report is generated in two essential formats:
 1. A visual, interactive **HTML file** for human analysis.
 2. A machine-readable **JSON object** for programmatic use and automation.
- **Log HTML Report:** The interactive HTML report is logged as an artifact to the MLflow run. This provides a rich, visual dashboard that data scientists can use for a deep-dive, manual inspection of any feature-level drift that is detected.
- **Extract and Log Metrics:** The script programmatically parses the JSON version of the report to extract crucial numerical metrics. These metrics, such as the overall **Share of Drifted Columns** and **per-feature drift scores** (e.g., p-values from statistical tests), are logged directly as metrics in MLflow. This is a critical step for automation, as it allows us to:
 - Build automated dashboards to track drift trends over time.
 - Set up automated alerts (e.g., via AWS Lambda and SNS) that can trigger if a drift score exceeds a predefined threshold, notifying the MLOps team of a potential issue before it significantly impacts business outcomes.

Metric	Value
train_vs_val_drift_ratio	0.3103
train_vs_val_asymmetrique_activity_sc...	0.9973
train_vs_val_asymmetrique_profile_sc...	0.9973
train_vs_val_total_time_spent_on_web...	0.3633
train_vs_val_a_free_copy_of_masterin...	0.0565
train_vs_val_asymmetrique_activity_in...	0.3151
train_vs_val_asymmetrique_profile_in...	0.6963
train_vs_val_city	0.4304

Part 9: Future State - Full MLOps Automation

This section outlines the architecture and implementation for transforming the manual, notebook-driven workflow into a fully automated, production-grade MLOps pipeline. The cornerstone of this future state is orchestration with **Amazon Managed Workflows for Apache Airflow (MWAA)**.

20. Future State: Full Automation with MWAA

The goal is to create a hands-off, scheduled system that proactively monitors for data drift and automatically retrains, validates, and deploys the model when necessary. This ensures the model remains accurate and relevant over time without requiring manual intervention.

20.1. Orchestrating the MLOps Lifecycle with Airflow

To achieve this automation, the entire MLOps workflow is orchestrated using **Amazon Managed Workflows for Apache Airflow (MWAA)**. MWAA is a managed service that runs Apache Airflow on AWS, removing the need to manage the underlying infrastructure.

The entire automated process—from data drift detection to model retraining and deployment—is defined as a **Directed Acyclic Graph (DAG)**. A DAG is a Python script that defines a collection of tasks, their relationships, and their dependencies, creating a clear, code-defined representation of our entire workflow.

20.2. Setting up the MWAA Environment

Before the automation DAG can be deployed, the MWAA environment and its required AWS resources must be provisioned and configured correctly.

Step 1: Create an S3 Bucket for MWAA Assets □

MWAA uses a dedicated S3 bucket to store all of its essential components.

- **Action:** Create a new S3 bucket with a unique, descriptive name (e.g., `mwaa-lead-conversion-dags`).
- **Best Practice:** Enable **versioning** on this bucket. This is crucial as it keeps a complete history of your DAGs and `requirements.txt` files, allowing for easy rollbacks and auditing.

Step 2: Configure the VPC and Networking

A secure and correctly configured network is essential for MWAA to communicate with other AWS services.

- **Use the Same VPC:** The MWAA environment **must be launched in the same VPC** as your Redshift, SageMaker, and MLflow resources. This ensures all communication happens over a private and secure network using VPC endpoints.
- **Subnets:** Provide at least **two private subnets**, each in a different Availability Zone. This is a requirement for MWAA to ensure high availability.
- **Internet Access:** The private subnets must have a route to a **NAT Gateway**. This allows the Airflow workers to securely access the public internet to install Python dependencies from PyPI, without being exposed to inbound internet traffic.

- **Security Group:** The MWAA environment should be associated with the `mlops-sg` security group. This pre-configured group ensures MWAA has the necessary inbound and outbound rules to communicate securely with the MLflow server, Redshift, S3, and other internal services.

Step 3: Create the MWAA Environment

- Navigate to the **Managed Workflows for Apache Airflow** service in the AWS Console and click **Create environment**.
- **Details:**
 - **Name:** LeadConversion-MLOps-Prod
 - **Airflow version:** Choose a recent, stable version (e.g., 2.8.1).
 - **S3 bucket:** Select the S3 bucket created in Step 1.
 - **DAGs folder path:** `dags`
 - **Requirements file path:** `requirements.txt`
 - **Networking:** Select your target VPC, the two private subnets, and the `mlops-sg` security group.
 - **Execution role:** Create a new IAM role that grants MWAA permissions to access its S3 bucket and write to CloudWatch Logs. **Crucially, attach additional policies** to this role that allow it to interact with SageMaker (for deployment), Redshift Data API (for ingestion), and MLflow (via SageMaker).
- Click **Create environment**. The provisioning process can take up to 30 minutes.

20.3. Preparing and Uploading the DAG File

- **Project Structure:** On your local machine, organize your project files as follows:
 - `mwaa_project/`
 - `└── dags/`
 - `└── data_drift_redshift_retrain.py`
 - `requirements.txt`
- **Create `requirements.txt`:** This file lists all Python packages that your DAG needs to run.
 - `apache-airflow-providers-amazon`
 - `mlflow-skinny`
 - `scikit-learn`
 - `pandas`
 - `numpy`
 - `boto3`
 - `evidently`
- **Upload to S3:** Upload the entire `dags` folder and the `requirements.txt` file to the root of your MWAA S3 bucket. MWAA will automatically synchronize the DAGs and install the specified Python requirements on its workers.

Amazon S3

General purpose buckets

- Directory buckets
- Table buckets
- Vector buckets [Preview](#)
- Access Grants
- Access Points for general purpose buckets
- Access Points for directory buckets
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- IAM Access Analyzer for S3

Block Public Access settings for

CloudShell Feedback

Objects (7)

[Create folder](#) [Upload](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Name	Type	Last modified	Size	Storage class
dags/	Folder	-	-	-
incoming/	Folder	-	-	-
raw/	Folder	-	-	-
requirements_5.txt	txt	July 21, 2025, 01:01:29 (UTC+05:30)	62.0 B	Standard
requirements.txt	txt	July 20, 2025, 13:12:41 (UTC+05:30)	320.0 B	Standard

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Amazon MWAA

Environments > MyAirflowEnvironment

MyAirflowEnvironment

[Edit](#) [Delete](#) [Open Airflow UI](#)

Details

Status Available

ARN [arn:aws:airflow:ap-south-1:888517279277:environment/MyAirflowEnvironment](#)

Airflow UI <http://cf107a0f-d0ae-47ef-b68e-725b3c2801e7.c0.ap-south-1.airflow.amazonaws.com>

Weekly maintenance window start (UTC)
Monday 23:30

Last update

Status SUCCESS

Created at 2025-07-21T19:47:17.000Z

Worker replacement strategy
Forced

CloudShell Feedback

21. The Automated Retraining DAG: A Detailed Look

The `data_drift_redshift_retrain.py` file defines the DAG that orchestrates our end-to-end MLOps workflow. Its primary purpose is to ensure the continued performance and relevance of our production model by proactively monitoring for data drift and automatically triggering a full retraining and deployment cycle when necessary. This process minimizes manual intervention and ensures the production model is always trained on the most relevant data.

21.1. Detailed Workflow Logic

The DAG is composed of a series of tasks executed in a predefined order with conditional logic.

Task 1: Scheduled Trigger

- Description:** The entire workflow is initiated on a recurring schedule, defined by a cron expression in the DAG definition (e.g., '`0 0 * * 0`' for a weekly run every Sunday at midnight).
- Purpose:** To regularly and proactively check the health of the model's data environment without waiting for performance metrics to drop significantly.

Task 2: Automated Drift Analysis

- **Description:** This task executes a Python script that performs data drift analysis using Evidently AI. It queries a **reference dataset** (the original training data) and a **current dataset** (e.g., the last 7 days of prediction data) from Amazon Redshift.
- **Output:** It generates a drift report and key metrics, such as the **Share of Drifted Columns**.
- **Integration:** The full HTML report and the summary metrics are logged to our MLflow Tracking Server for auditing and visualization.

Task 3: Conditional Retraining (Branching Logic)

- **Description:** This task acts as a decision gate, implemented using Airflow's `BranchPythonOperator`. It retrieves the drift analysis results from the previous task (via XComs) and evaluates them against a predefined, configurable threshold.
- **Condition:** For example, `if Share_of_Drifted_Columns > 0.2::`
- **If Drift Detected (TRUE):** The DAG proceeds down the retraining path.
- **If No Drift Detected (FALSE):** The DAG follows a short-circuit path, logs a message like "No significant data drift detected," and successfully concludes the run.

Task 4: Automated Model Retraining

- **Description:** If triggered, this task initiates the full model training pipeline. It uses the latest data to train a new "challenger" model, encapsulating all necessary steps: data ingestion, preprocessing, feature engineering, and training.
- **Purpose:** To create a new model candidate that is adapted to the latest data patterns.

Task 5: Automated Registration (Champion-Challenger Validation)

- **Description:** The performance of the newly trained "challenger" model is compared against the "champion" model (the version currently in production). The comparison uses a key business metric like the F1-score.
- **Condition:** If the challenger model's performance is better than the champion's, it is registered as a new version in the MLflow Model Registry.
- **If Not Better:** The new model is still logged for analysis, but it is not promoted, and the workflow may stop or send a specific "retraining failed to improve performance" alert.

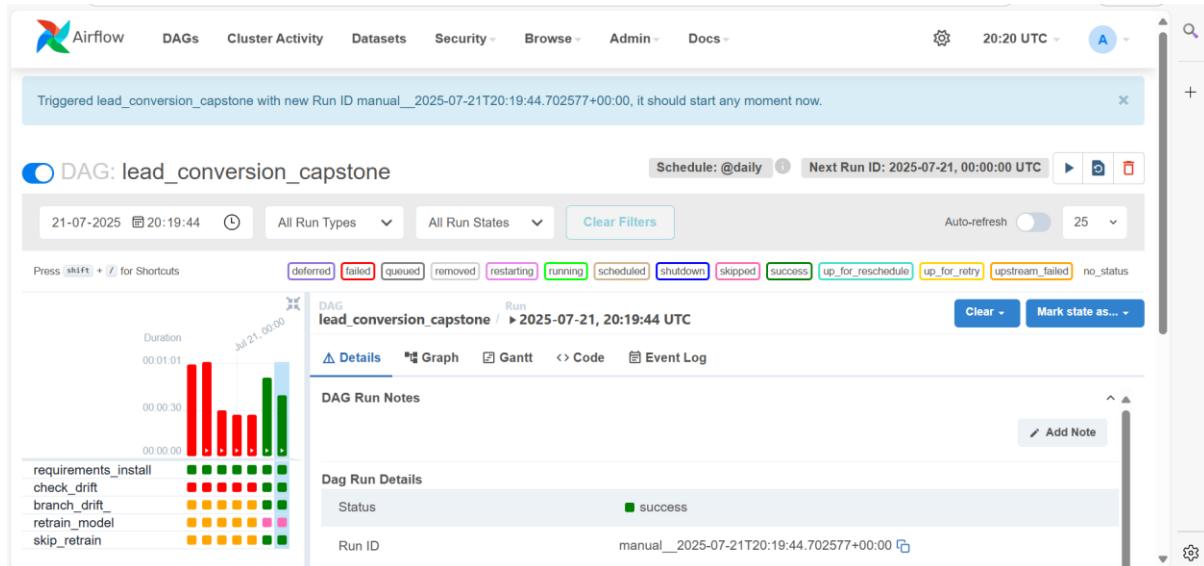
Task 6: Automated Deployment

- **Description:** This final task is triggered only if a new, superior model version was successfully registered. It interacts with the MLflow Model Registry to programmatically transition the stage of the new model version from "Staging" to "Production."
- **Impact:** This transition automatically makes the new model the active version to be served by the SageMaker Endpoint, completing the zero-touch deployment cycle.

Task 7: Alerting and Notification

- **Description:** This is a critical operational task integrated throughout the DAG. At key decision points, state changes, or failures, the DAG sends notifications to the MLOps team via channels like Slack or email.
- **Events Triggering Alerts:**
 - Significant data drift is detected.

- The retraining process begins.
- A new model is successfully deployed to production.
- Any task within the DAG fails.



Of course. Here is a refined and professionally formatted version of the Technology Stack appendix, with brief descriptions added for clarity.

Part 6: Appendix

20. Technology & Libraries Stack

This section provides a comprehensive overview of the cloud services, platforms, and Python libraries that constitute the technology stack for this project.

Cloud Platform & Services

- **Cloud Platform: Amazon Web Services (AWS)**
 - The primary cloud provider for all infrastructure, data, and machine learning services.
- **Data Storage: Amazon S3, Amazon Redshift**
 - S3 serves as the central data lake for raw data, processed data, and model artifacts.
 - Redshift Serverless acts as the high-performance data warehouse for structured, analytics-ready data.
- **Data Integration (ETL): AWS Glue (Studio, Crawlers, Data Catalog)**

- A managed, serverless ETL service used to extract data from S3, transform it, and load it into Redshift.
- **Security & Networking:** **AWS IAM, AWS VPC, AWS Secrets Manager, AWS KMS**
 - **IAM** manages secure access to AWS resources.
 - **VPC** provides a logically isolated network for all project resources.
 - **Secrets Manager** securely stores and manages database credentials.
 - **KMS** is used for managing encryption keys.
- **ML Development:** **Amazon SageMaker Studio**
 - The integrated development environment (IDE) for all machine learning activities, from data exploration to model training and deployment.
- **MLOps & Tracking:** **MLflow (on Amazon SageMaker)**
 - The platform for managing the end-to-end machine learning lifecycle, including experiment tracking and model registration.
- **Orchestration (Future):** **Managed Workflows for Apache Airflow (MWAA)**
 - The service planned for orchestrating the entire MLOps workflow as an automated, scheduled pipeline.
- **Serving (Prototyping):** **Flask, Ngrok**
 - **Flask** is used to create a lightweight web server to wrap the model in a REST API.
 - **Ngrok** exposes the local Flask server to the internet for testing and demonstration purposes.
- **Core Python Ecosystem**
- **Data Handling:** **pandas, numpy**
 - The foundational libraries for data manipulation, cleaning, and numerical operations.
- **Visualization:** **matplotlib, seaborn**
 - Used for creating static, animated, and interactive visualizations for exploratory data analysis.
- **ML & Preprocessing:** **scikit-learn**
 - The primary toolkit for building preprocessing pipelines, training baseline models, and evaluating performance.
- **Advanced Models:** **xgboost, lightgbm**
 - High-performance, gradient-boosting libraries used for training advanced and highly accurate models.
- **Model Explainability:** **shap**
 - Used to compute and visualize SHAP values to explain model predictions.
- **Data Drift:** **evidently**
 - The library used to generate data drift and performance monitoring reports.
- **AWS Interaction:** **boto3**
 - The AWS SDK for Python, used to programmatically interact with services like Redshift and SageMaker.
- **Serialization:** **joblib**
 - Used for saving and loading Python objects, specifically the trained model and preprocessing pipeline artifacts.

21. Code Repository & Resources

This section provides links to the project's source code repository and the official documentation for the key technologies utilized in this solution.

Code Repository

The complete source code for this project, including all scripts, notebooks, and configuration files, is hosted on GitHub.

- **GitHub Repository:** https://github.com/Madhu678-coder/lead_conversion_prediciton

Key Technology Documentation

For further information on the core services and libraries used, please refer to their official documentation.

- **AWS Glue:** [Official Documentation](#)
- **Amazon SageMaker:** [Official Documentation](#)
- **MLflow:** [Official Documentation](#)
- **Evidently AI:** [Official Documentation](#)