
Lead Conversion Classification

Complete MLOPs implementation in AWS Cloud

Subramanya Rithwik Jakka

Table of Contents

Part 1: Business & System Overview

1. Executive Summary & Vision
 - o 1.1. Problem Statement
 - o 1.2. Business Objectives
 - o 1.3. Solution Overview
2. Success Criteria
 - o 2.1. Business Success Criteria
 - o 2.2. Machine Learning Success Criteria
3. System Architecture: A Deep Dive
 - o 3.1. Architectural Blueprint Diagram
 - o 3.2. Key Component Breakdown
 - o 3.3. End-to-End Data and Model Flow
4. Data Dictionary
 - o 4.1. Feature Descriptions

Part 2: Foundational AWS Infrastructure Setup

5. Identity & Access Management (IAM)
 - o 5.1. The Principle of Least Privilege
 - o 5.2. Creating the Core Service Roles
6. Networking & Security (VPC)
 - o 6.1. VPC and Subnet Strategy
 - o 6.2. Creating VPC Endpoints for Secure Communication
 - o 6.3. Configuring the Security Group Firewall

Part 3: Data Engineering & Warehousing

7. Data Lake & Warehouse Setup
 - o 7.1. S3 Bucket Configuration
 - o 7.2. Amazon Redshift Serverless Provisioning
 - o 7.3. Securing Credentials with AWS Secrets Manager
8. The Ingestion Pipeline: S3 to Redshift
 - o 8.1. ETL with AWS Glue: Overview
 - o 8.2. Creating the Glue Connection to Redshift
 - o 8.3. The Glue Crawler: Automatic Schema Discovery

- 8.4. Building the Visual ETL Job (`s3ToRedshift`)
- 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

- **Experiment Tracking with MLflow**
 - 16.1. Setting Up the MLflow Tracking Server
 - 16.2. Integrating MLflow into the Training Pipeline
- **Model Registration and Lifecycle Management**
 - 17.1. Saving and Registering the Champion Model
 - 17.2. The MLflow Model Registry
- **Serving and Inference**
 - 18.1. Loading the Production Model from the Registry
 - 18.2. Real-Time Prediction with Flask and Ngrok
- **Post-Deployment: Drift Detection**
 - 19.1. Monitoring for Data Drift with Evidently AI
 - 19.2. Automating Drift Analysis and Logging
- **Automation with MWAA**
 - 20.1. Orchestrating the MLOps Lifecycle with Airflow
 - 20.2. The Automated Retraining DAG
- **SageMaker Model Deployment**
 - 21.1 Folder Structure (Used in SageMaker)

- 21.2 Model Packaging and Upload
- 21.3 Upload the model bundle to S3
- 21.4 SageMaker Endpoint Creation via Notebook
- 21.5. Deployment Result
- 21.6 Root Cause Analysis of Failure

Part 6: Appendix

21. Technology & Libraries Stack

22. Resources

23. GitHub Repository

Part 1: Business & System Overview

1. Executive Summary & Vision

1.1. Problem Statement

In the highly competitive B2B sales environment, marketing and sales teams invest significant resources in acquiring leads. However, a large proportion of these leads fail to convert into paying customers, resulting in wasted efforts and suboptimal ROI.

Currently, there is a lack of intelligent mechanisms to prioritize high-potential leads based on historical data and behavioral attributes.

1.2. Business Objectives

The core business goal is to enable **data-driven sales optimization** by accurately identifying leads with the highest likelihood of conversion. By leveraging machine learning and modern MLOps practices, this system empowers sales representatives to focus their efforts on the most promising prospects.

This initiative directly supports the following strategic objectives:

- **Maximize Conversion Rate:** By predicting the probability of conversion for each lead, the system allows teams to allocate time and energy toward leads with higher success potential.
- **Maximize Lead Prioritization:** The model ranks leads based on predicted conversion likelihood, enabling targeted and prioritized outreach strategies.
- **Minimize Wasted Outreach:** Reduces time spent on low-quality or cold leads by eliminating guesswork through predictive analytics.
- **Maximize ROI:** Focused efforts on high-potential leads translate into improved sales efficiency, higher revenue, and better return on marketing and operational investments.

1.3. Solution Overview

To address the lead conversion prediction challenge, we have architected and implemented a **scalable, cloud-native machine learning pipeline** on **Amazon Web Services (AWS)**. This end-to-end system automates the entire ML lifecycle—from raw data ingestion to real-time inference—ensuring efficiency, reproducibility, and scalability.

The solution leverages a **modern AWS-native data and ML stack**:

- **AWS S3** for secure and scalable data storage
- **AWS Glue and Redshift** for data transformation and querying
- **Amazon SageMaker** for model development and training
- **MLflow** for experiment tracking and model lifecycle management

The resulting pipeline produces a **high-performing predictive model** that scores incoming leads and integrates seamlessly into downstream applications via a deployed **prediction API**. This empowers sales teams with **actionable intelligence** to prioritize high-value leads, optimize outreach, and ultimately increase conversion rates.

2. Success Criteria

The success of this project is measured against clearly defined business and technical metrics.

2.1. Business Success Criteria

1. Increase Conversion Rate by Prioritizing High-Potential Leads by at least 20%:

- Focus sales efforts on leads most likely to convert, as identified by predictive modeling.
- Elevate the overall conversion rate by at least 20% by systematically targeting high-value opportunities.

2. Improve Sales Team Efficiency by Reducing Time on Low-Value Leads

- Decrease the time and effort spent pursuing low-potential leads.
- Optimize resource allocation, enabling the sales team to concentrate on leads with greater conversion probability.

2.2. Machine Learning Success Criteria

1. Prediction Accuracy $\geq 90\%$

- The model must correctly classify at least 90% of all leads, indicating strong overall predictive performance.
- High accuracy ensures reliable identification and separation of high- and low-potential leads.

2. Precision for "Converted" Class $\geq 85\%$

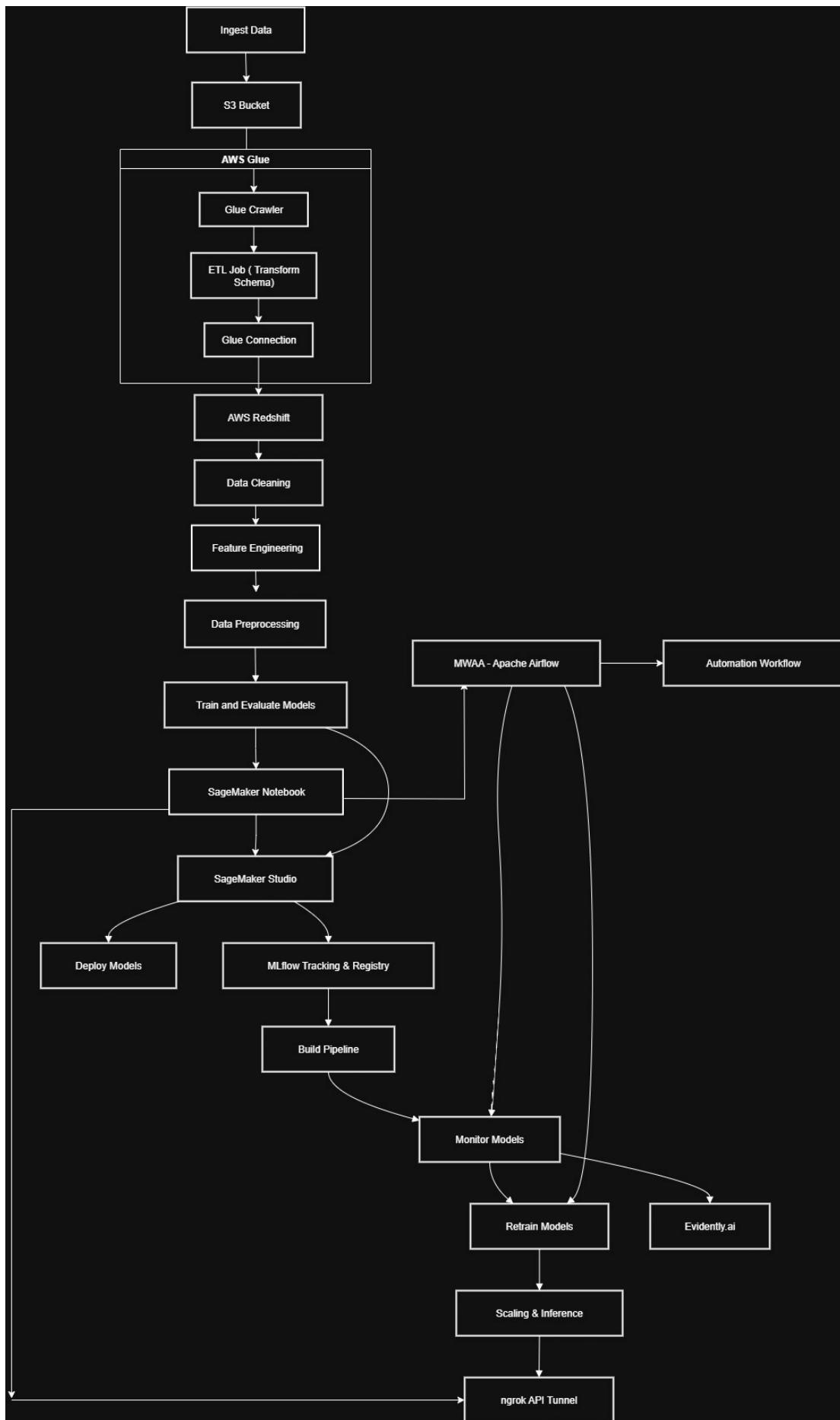
- For leads predicted as "converted," at least 85% must actually convert.
- This minimizes false positives, ensuring that sales efforts are focused on genuinely high-potential leads.

3. ROC AUC ≥ 0.90

- The model's ROC AUC (Receiver Operating Characteristic Area Under Curve) should be 0.90 or higher.
 - A high ROC AUC score demonstrates strong ability to distinguish between converted and non-converted leads, supporting confident prioritization and decision-making.
-

3. System Architecture: A Deep Dive

3.1. Architectural Blueprint Diagram : The system is designed as a modular, event-driven architecture that ensures scalability, security, and maintainability.



This diagram illustrates a cloud-based, automated end-to-end machine learning (ML) pipeline utilizing AWS services, workflow orchestration, and tools for monitoring, retraining, and deployment.

1. Data Ingestion and Storage

- **Ingest Data:** Data is collected from one or more raw sources.
- **S3 Bucket:** The ingested data is stored securely in AWS S3, a scalable object storage service.

2. Data Cataloging and Transformation (AWS Glue)

- **AWS Glue:** A serverless data integration service.
 - **Glue Crawler:** Scans the data in S3, infers schema, and creates metadata tables.
 - **ETL Job (Transform Schema):** Performs Extract, Transform, Load (ETL) operations to clean and reformat the data.
 - **Glue Connection:** Facilitates integration and migration of transformed data.

3. Data Warehousing

- **AWS Redshift:** The transformed data is loaded into AWS Redshift, a scalable cloud data warehouse, for large-scale analytics and further processing.

4. Data Preparation

- **Data Cleaning:** Removal of noise, inconsistencies, or errors in the dataset.
- **Feature Engineering:** Creating new features or transforming variables to improve model effectiveness.
- **Data Preprocessing:** Further manipulation to ready data for modeling (e.g., normalization, encoding).

5. Model Development

- **Train and Evaluate Models:** ML models are built, trained, and evaluated.
- **SageMaker Notebook / Studio:** Amazon SageMaker is used for interactive development and model experimentation.

6. Model Deployment and Management

- **Deploy Models:** Trained models are deployed for inference (predictions).
- **MLflow Tracking & Registry:** Manages run metadata, model versions, and provenance.
- **Build Pipeline:** Automates deployment, testing, and integration steps.

7. Orchestration and Automation

- **MWAA - Apache Airflow:** Manages, schedules, and automates the pipeline's workflows.

- **Automation Workflow:** Enables repeatable, automated execution of steps from ingestion to deployment.

8. Monitoring and Feedback

- **Monitor Models:** Ongoing tracking of model performance and system health.
- **Evidently.ai:** Supports advanced model monitoring and drift detection.
- **Retrain Models:** Triggers retraining cycles if models deteriorate or new data arrives.
- **Scaling & Inference:** Supports model scaling and deploying to production for live inference.

9. API & External Integration

- **ngrok API Tunnel:** Creates secure tunnels to expose local web servers or APIs to the internet, enabling integrations, and remote access.

Key Features

- **Automated:** End-to-end orchestration using Apache Airflow.
- **Modular:** Clear separation of ingestion, transformation, modeling, deployment, and monitoring.
- **Scalable:** Built atop AWS's scalable infrastructure for storage, compute, and analytics.
- **Trackable:** MLflow and Evidently.ai ensure model versioning, traceability, and performance monitoring.
- **Continuous Feedback:** Automated retraining and monitoring loops enable adaptive and robust ML solutions.

3.2. Key Component Breakdown

- **Data Storage (Amazon S3):** Serves as the foundational storage layer. Different buckets/prefixes are used for raw-data, processed-data, and model-artifacts, ensuring a clean separation of concerns.
- **IAM (Identity and Access Management):** Provides the security backbone.
 - **AWS Glue Role:** Grants AWS Glue permissions to access S3 data and connect to the Redshift cluster.
 - **AWS SageMaker Role:** Allows SageMaker to read data from Redshift/S3 and write model artifacts back to S3.
- **AWS Glue (ETL Platform):** The serverless data integration engine.
 - **Connection to Redshift:** A pre-configured connector allowing Glue jobs to read from and write to Redshift.
 - **Crawler & Data Catalog:** Automatically scans and catalogs data schemas, making data discoverable and queryable.
 - **ETL Job (s3ToRedshift):** The core job that extracts raw S3 data, cleans it, and loads it into Redshift.
- **Networking (VPC):** Creates a private, isolated network for our resources.
 - **Security Group:** A virtual firewall controlling traffic to VPC endpoints.
 - **Endpoints:** Provide secure, private connections for S3, Redshift, Secrets Manager, KMS, and STS, preventing data exposure to the public internet.

- **Data Warehouse (Amazon Redshift):** The central repository for structured, analytics-ready data. It serves as the single source of truth for the machine learning model.
- **SageMaker (ML Platform):** The integrated development environment for all ML tasks.
 - **SageMaker Studio Lab:** Used for data preparation, interactive development, and model building.
 - **Artifact Management:** All trained models are saved to S3 for persistence and versioning.
- **Serving and Monitoring Stack:**
 - **MLflow:** Integrated for comprehensive experiment tracking, model logging, and lifecycle management via the Model Registry.
 - **Flask:** A lightweight web framework used to create a simple API to serve the model's prediction endpoint.
 - **Ngrok:** A utility to create a secure tunnel, exposing the local Flask API for testing and integration.
 - **MLflowUI:** The web dashboard for tracking, comparing, and visualizing ML experiments.

3.3. End-to-End Data and Model Flow

1. **Raw Data Ingestion:** Raw lead data (.csv) is uploaded to the designated Amazon S3 bucket.
2. **ETL with AWS Glue:** The Glue Crawler scans the S3 bucket and updates the Data Catalog. A scheduled Glue ETL job then extracts the data, transforms it, and loads it into Amazon Redshift.
3. **Redshift Data Access:** Redshift now holds clean, structured tables, ready for machine learning.
4. **Data Access for ML:** Amazon SageMaker securely accesses the Redshift data via VPC endpoints and its IAM role. Training, validation, and test sets are prepared.
5. **ML Development in SageMaker:** Models are built, trained, and evaluated in SageMaker Studio Lab. MLflow tracks all experiments and metrics. The final model artifacts are saved to S3.
6. **Model Logging & Serving:** The best model is registered in the MLflow Model Registry and promoted to "Production." It is then served as an API via a Flask application, exposed for testing with Ngrok.
7. **Feedback Loop & Automation:** Glue schedules ensure continuous data ingestion. The MLOps framework supports systematic retraining and versioning.

4. Data Dictionary

4.1. Feature Descriptions

The following table describes the features available in the dataset.

Feature Name	Explanation
id	Unique identifier for each lead.

Feature Name	Explanation
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.

Feature Name	Explanation
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...	Whether the lead requested a free resource (e.g., eBook) (Yes/No).
Last Notable Activity	The most significant recent action taken by the lead.
Export to Sheets	

Part 2: Foundational AWS Infrastructure Setup

5. Identity & Access Management (IAM)

5.1. The Principle of Least Privilege

Security is paramount. Our architecture adheres to the principle of least privilege, meaning every component is granted only the permissions essential for its designated function. This is achieved by creating highly specific IAM roles that services like Glue and SageMaker assume to perform their tasks, minimizing potential security risks.

5.2. Creating the Core Service Roles

We create two primary roles: one for Glue and one for SageMaker. The process below details the creation of the Glue role, which is the more complex of the two.

Console Walkthrough: Creating `glueaccessrole`

1. **Navigate to IAM:** In the AWS Management Console, go to the IAM service.
2. **Create Role:** Click on **Roles** in the left navigation pane, then click **Create role**.
3. **Select Trusted Entity:**
 - **Trusted entity type:** Select **AWS service**.
 - **Use case:** Choose **Glue** from the dropdown menu. This automatically creates a trust policy allowing the AWS Glue service (`glue.amazonaws.com`) to assume this role.
 - Click **Next**.
4. **Add Permissions:** On the "Add permissions" page, search for and attach the following AWS managed policies. For a production environment, these should be replaced with more fine-grained, custom policies.
 - `AmazonS3FullAccess`
 - `AmazonRedshiftFullAccess`
 - `AWSGlueConsoleFullAccess`
 - `AWSGlueServiceRole`
 - `SecretsManagerReadWrite`
 - `AWSKeyManagementServicePowerUser`
 - Click **Next**.
5. **Name and Review:**

- **Role name:** glue-lead-role
- Review the attached policies and the trust policy. The trust policy JSON should look like this:

JSON

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

- Click **Create role**.

Policy name	Type	Attached entities
AmazonRedshiftFullAccess	AWS managed	1
AmazonS3FullAccess	AWS managed	1
AWSGlueConsoleFullAccess	AWS managed	1
AWSGlueServiceRole	AWS managed	1
AWSKeyManagementService	AWS managed	1
SecretsManagerReadWrite	AWS managed	1

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

6. Networking & Security (VPC)

6.1. VPC and Subnet Strategy

To ensure all inter-service communication is secure and isolated from the public internet, all resources are deployed within a Virtual Private Cloud (VPC). We utilize the default VPC for this guide, but a custom VPC is recommended for production. The VPC contains multiple subnets across different Availability Zones for high availability.

6.2. Creating VPC Endpoints for Secure Communication

VPC Endpoints allow services within the VPC to communicate with other AWS services as if they were on the same network, without traffic ever leaving the AWS backbone.

Console Walkthrough: Creating Endpoints

1. Navigate to the VPC service in the AWS Console.
2. Select **Endpoints** and click **Create endpoint**.
3. **Endpoint 1: S3 (Gateway)**
 - o **Service category:** AWS services
 - o **Service name:** Search for `s3` and select the service with **Type: Gateway**.
 - o **VPC:** Select your target VPC.
 - o **Route tables:** Select the route tables for the subnets where your resources reside. This adds a route to the table, directing S3 traffic through the gateway.
 - o Click **Create endpoint**.
4. **Endpoint 2: Redshift (Interface)**
 - o Click **Create endpoint** again.
 - o **Service category:** AWS services
 - o **Service name:** Search for `redshift` and select the service `com.amazonaws.<region>.redshift`.
 - o **VPC:** Select your target VPC.
 - o **Subnets:** Select at least two subnets in different Availability Zones.
 - o **Security groups:** Select the security group created in the next section.
 - o Click **Create endpoint**.
5. **Endpoints 3, 4, 5: Secrets Manager, KMS, and STS (Interface)**
 - o Repeat the process for the following services, selecting the same VPC, subnets, and security group for each:
 - `com.amazonaws.<region>.secretsmanager`
 - `com.amazonaws.<region>.kms`
 - `com.amazonaws.<region>.sts`

Name	VPC endpoint ID	Endpoint type	Status
-	vpce-03aa60a0754af957	Interface	Available
s3-end-point	vpce-0d14956be8e372b1d	Gateway	Available
redshift-endpoint	vpce-0fd138d0366619c04	Interface	Available
secretmanager-endpoint	vpce-037e5224a58c24db	Interface	Available
kms-endpoint	vpce-07da1ea1628e42a4	Interface	Available
sts-end-point	vpce-06fda6c2e33970d7c	Interface	Available
-	vpce-06b25f18d55db2e8	Interface	Available

6.3. Configuring the Security Group Firewall

A security group acts as a stateful firewall for resources, controlling inbound and outbound traffic.

Console Walkthrough: Configuring the Security Group

1. Navigate to **VPC > Security Groups** and click **Create security group**.
2. **Name and VPC:** Name it `lead-security-group` and associate it with your VPC.
3. **Configure Inbound Rules:** Select the new security group and click the **Inbound rules** tab, then **Edit inbound rules**.
 - o **Rule 1 (For Redshift Client Access):**
 - **Type:** Redshift
 - **Protocol:** TCP
 - **Port range:** 5439
 - **Source:** My IP (to allow access from your local machine).
 - o **Rule 2 (For Internal Communication):**
 - **Type:** All TCP
 - **Protocol:** TCP
 - **Port range:** 0 - 65535
 - **Source:** Custom. Start typing the ID of the `lead-security-group` group itself (`sg-xxxxxxxx`). This self-referencing rule allows all resources within the security group to communicate with each other freely. This is critical for Glue, SageMaker, and Redshift to interact.
4. **Save Rules:** Click **Save rules**. The default outbound rule, which allows all traffic out, is sufficient for this project.

The screenshot shows the AWS EC2 Security Groups console. The left sidebar is collapsed. The main area displays the details of the security group `sg-03265e475a479cd2f - lead-security-group`. The **Inbound rules** tab is selected, showing two rules:

IP version	Type	Protocol	Port range	Source
IPv4	Redshift	TCP	5439	14.194.147.50/32
-	All TCP	TCP	0 - 65535	<code>sg-03265e475a479cd2f..</code>

Part 3: Data Engineering & Warehousing

7. Data Lake & Warehouse Setup

7.1. S3 Bucket Configuration

Amazon S3 is the foundation of our data storage. We use a single bucket with a logical prefix structure to organize data by stage.

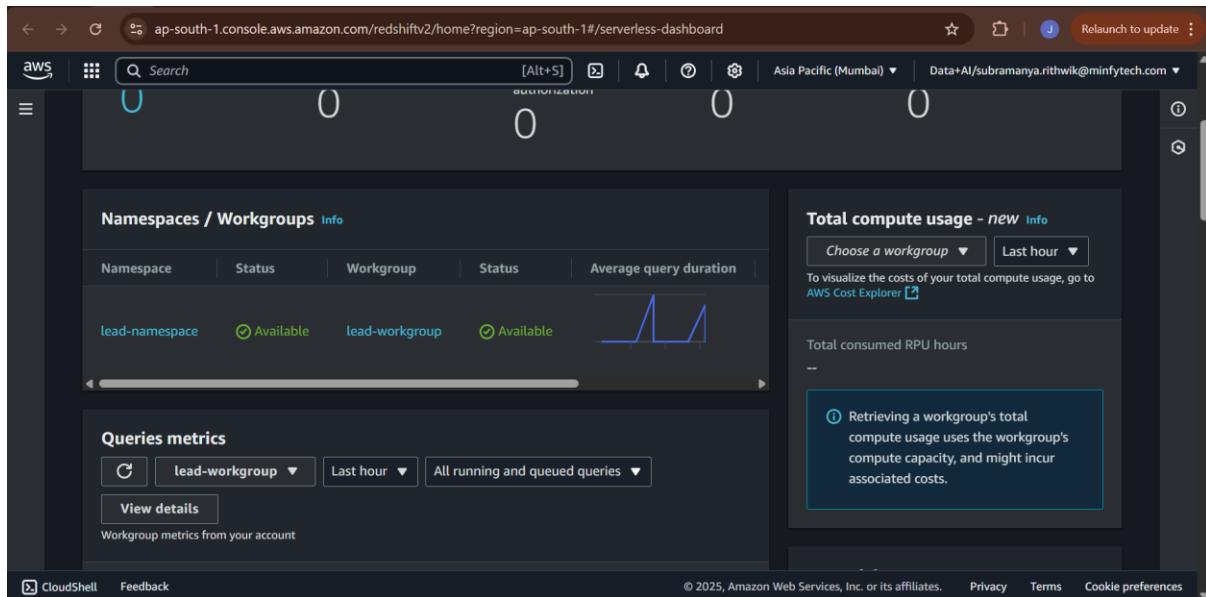
- **Bucket Name:**
- **Prefixes:**
 - s3://leadconversiondata/raw-data/: For incoming, unaltered lead data.
 - s3://leadconversiondata/processed-data/: For data cleaned and feature-engineered by Glue.
 - s3://leadconversiondata/model-artifacts/: For storing trained model binaries, SHAP plots, and other ML assets.

Name	Type	Last modified	Size	Storage class
lead_conversion_mod/el/	Folder	-	-	-
lead_scoring_100.csv	CSV	July 18, 2025, 09:32:53 (UTC+05:30)	28.6 KB	Standard
mlflow-artifacts/	Folder	-	-	-
model_bundle.tar.gz	gz	July 20, 2025, 13:12:57 (UTC+05:30)	8.4 KB	Standard

7.2. Amazon Redshift Serverless Provisioning

Redshift Serverless provides a simple, auto-scaling data warehouse.

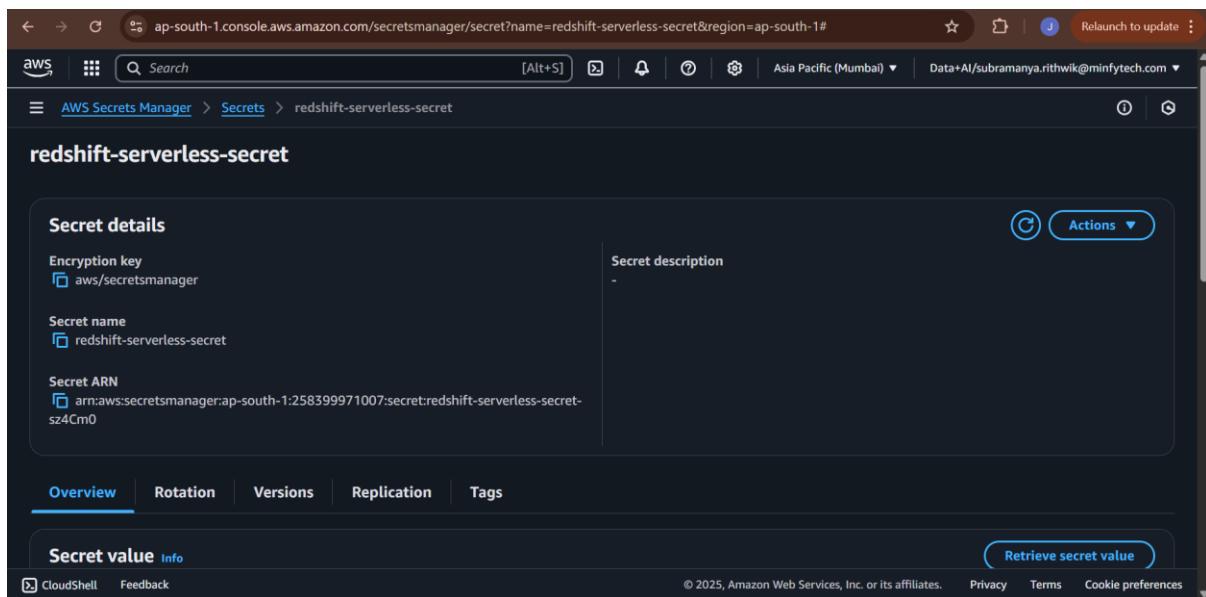
1. Navigate to the **Amazon Redshift** service and select **Try Redshift Serverless**.
2. Configure the workgroup and namespace, accepting the defaults.
3. In the **Network and security** section, associate the workgroup with your VPC and the `glue-lead-role` security group.
4. Review and create the serverless endpoint. This will take several minutes.



7.3. Securing Credentials with AWS Secrets Manager

We avoid hardcoding database credentials by storing them securely in Secrets Manager.

1. In the Redshift console, set a password for the default `awsuser`.
2. Navigate to **AWS Secrets Manager** and click **Store a new secret**.
3. **Secret type:** Select `Credentials` for Redshift cluster.
4. Enter the `awsuser` username and the password you set.
5. **Secret name:** Use a descriptive name like `prod/redshift/credentials`.
6. Complete the process, disabling automatic rotation for this guide. Note the ARN of the created secret.



8. The Ingestion Pipeline: S3 to Redshift

8.1. ETL with AWS Glue: Overview

AWS Glue provides a serverless framework to extract data from S3, transform it, and load it into Redshift. This process relies on a Glue Connection to securely access Redshift and a Glue Crawler to understand the data's structure.

8.2. Creating the Glue Connection to Redshift

1. Navigate to **AWS Glue** and select **Connections**.
2. **Create connection:**
 - o **Name:** Redshift Connection
 - o **Type:** Amazon Redshift
 - o Select Connect to a Redshift Serverless workgroup and choose your workgroup.
 - o For **Authentication**, select the AWS Secret (`prod/redshift/credentials`) you created.
3. **Test Connection:** After creating, select the connection and click **Actions > Test connection**. Choose the `glue-lead-role`. A "Success" message confirms that your networking and permissions are correct.

The screenshot shows the AWS Glue Studio interface. On the left, there is a sidebar with navigation links like 'Getting started', 'ETL jobs', 'Visual ETL', 'Notebooks', 'Job run monitoring', 'Data Catalog tables', 'Data connections', 'Workflows (orchestration)', 'Data Catalog', 'Databases', 'Tables', 'Stream schema registries', 'Schemas', 'Connections', 'Crawlers', 'Classifiers', and 'Catalog settings'. The 'Connections' link is highlighted. In the main content area, the title 'Redshift connection' is displayed above a 'Connection details' section. The details include:

- Connector type: REDSHIFT
- Subnet: subnet-0d2323ee771034b36
- Description: -
- Last modified: 2025-07-18 10:25:28.078000
- Database: dev
- Port: 5439
- Require SSL connection: -
- Security groups: sg-03265e475a479cd2f
- Created on: 2025-07-18 10:25:28.078000
- Version: 2
- Host name: lead-workgroup.258399971007.ap-south-1.redshift-serverless.amazonaws.com

At the bottom of the page, there are links for 'CloudShell', 'Feedback', and copyright information: '© 2025, Amazon Web Services, Inc. or its affiliates.' and 'Privacy Terms Cookie preferences'.

8.3. The Glue Crawler: Automatic Schema Discovery

The crawler scans our raw data and creates a metadata table in the Glue Data Catalog.

1. In Glue, select **Crawlers** and **Create crawler**.
2. **Name:** lead-crawler.
3. **Data source:** Point it to the S3 path `s3://leadconversiondata/raw-data/`.
4. **IAM Role:** Select `glue-lead-role`.
5. **Output:** Create a new database in the Glue Data Catalog named `glue-crawler-db`.
6. Run the crawler. Upon completion, a new table will be visible in the `glue-crawler-db` database.

The screenshot shows the AWS Glue Crawler properties page for 'lead-crawler'. The crawler is named 'lead-crawler' and has an IAM role 'glue-lead-role'. It is associated with a database 'glue-crawler-db' and is in a 'READY' state. The 'Crawler runs' tab shows one run listed.

8.4. Building the Visual ETL Job (`s3ToRedshift`)

1. In Glue, navigate to **Glue Studio** and create a new **Visual ETL** job.
2. **Source Node:** Select **AWS Glue Data Catalog** and choose the table created by the crawler.
3. **Transform Node:** Add any desired transformations. For a simple ingestion, you might use the **Change Schema** transform to ensure data types are correct.
4. **Target Node:** Select **Amazon Redshift**.
 - o **Connection:** Choose Redshift Connection.
 - o **Target options:** Specify the target database (`dev`) and table name (`public.lead_scoring`).
5. **Job Details:** Name the job `lead_job` and assign the `glue-lead-role`.
6. Save and **Run** the job. After it succeeds, you can query the `public. lead_scoring` table in Redshift.

The screenshot shows the AWS Glue Studio Visual ETL job editor for 'lead_job'. The job graph shows a flow from an Amazon S3 source, through a Change Schema transform, to an Amazon Redshift target.

The screenshot shows the AWS Glue Job Run details page. The job name is 'lead_job' and the run status is 'Succeeded'. The run was initiated at 07/18/2025 11:48:15 and completed at 07/18/2025 11:49:15, taking 13 seconds. The last modified time is 07/18/2025 06:19:50. The job has a max capacity of 100 workers and ran on a single worker.

Run details	Info
Job name	lead_job
Id	jr_13e9b69e81c644373113c95070760ad7198455457ddeba54a13fbf922e5c4e4f
Run status	Succeeded
Glue version	5.0
Retry attempt number	Initial run
Start time (Local)	07/18/2025 11:48:15
End time (Local)	07/18/2025 11:49:15
Start time (UTC)	2025/07/18 06:19:50
End time (UTC)	13 seconds
Execution time	1 minute 22 seconds
Last modified on (UTC)	2025/07/18 06:19:50
Trigger name	-
Security configuration	-
Timeout	480 minutes
Max capacity	100
Number of workers	1
Worker type	Standard

9. Data Access from SageMaker

9.1. Using the Redshift Data API with `boto3`

To load data into our SageMaker notebook for ML development, we use the Redshift Data API. This avoids managing traditional JDBC/ODBC drivers and is ideal for serverless interaction. The following Python script automates this process.

9.2. Step-by-Step Code Walkthrough

Step 1: Configuration This block sets up all the necessary parameters to connect to Redshift.

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 = 'lead-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.lead_scoring'
```

Step 2: Create Redshift Data API Client The `boto3` library is used to create a low-level client for the Redshift Data API service.

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 The `execute_statement` function sends the SQL query to Redshift asynchronously.

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 poll the API to check the query's status.

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)
```

Step 5 & 6: Retrieve and Parse Query Results Once finished, we fetch the results and parse them into a structure suitable for a DataFrame.

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])
```

Step 7 & 8: Convert to Pandas DataFrame and Display The final step is to create a Pandas DataFrame, the standard tool for data manipulation in Python.

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())
```

The screenshot shows the AWS Redshift Query Editor v2 interface. On the left, the sidebar displays the database schema with 'native databases (2)' under 'Serverless: lead-workgroup'. One of these databases is 'dev', which contains 'public' and 'sample_data_dev' schemas. The 'public' schema has 'Tables' (1), 'Views' (0), 'Functions' (0), and 'Stored Procedures' (0). The 'sample_data_dev' schema also contains one table. In the main pane, a query is running in 'Untitled 1': 'select * from lead_scoring;'. The results are displayed in a table titled 'Result 1 (100)'. The table has columns: city, do_not_call, converted, total_time_spent_on_lead, and magazine. The data shows various entries for Mumbai, with some rows having 'No' in the 'do_not_call' column and others having '1'. The 'converted' column has values 0 and 1. The 'total_time_spent_on_lead' column has values ranging from 0 to 1640. The 'magazine' column has values 'No' and 'Yes'. At the bottom of the results pane, it says 'Query ID 6137 Elapsed time: 5247 ms Total rows: 100'.

Part 4: Machine Learning Development

10. Amazon SageMaker: The ML Development Environment

All machine learning development for this project is conducted within **Amazon SageMaker**, a fully managed service that provides every developer and data scientist with the ability to prepare, build, train, and deploy high-quality ML models quickly. We specifically use **SageMaker Studio**, an integrated development environment (IDE) for machine learning.

10.1. Creating a SageMaker Domain

A SageMaker Domain is the primary entry point for using SageMaker Studio. It consists of a list of authorized users, a variety of security controls, and an Amazon Elastic File System (EFS) volume that contains the data, notebooks, and other artifacts for the users in the Domain.

Console Walkthrough: Creating the Domain

- Navigate to SageMaker:** In the AWS Management Console, go to the **Amazon SageMaker** service.
- Create a Domain:** On the SageMaker dashboard, click on **Domain** in the left navigation pane, then click **Create Domain**.
- Setup:** Choose the **Standard setup** option.
- Authentication:** For this guide, select **AWS Identity and Access Management (IAM)** as the authentication method.
- Configure Network and Storage:**
 - VPC:** Select the same Virtual Private Cloud (VPC) that your Redshift cluster and other resources are in. This is critical for secure communication.
 - Subnet(s):** Choose the private subnets associated with your VPC. Selecting multiple subnets across different Availability Zones ensures high availability.
 - Security Group:** Attach the `mlops-sg` security group that was configured in Part 2. This allows the SageMaker environment to communicate with other services like Redshift and the MLflow server through the private VPC endpoints.
- IAM Role:** Create a new IAM role or use an existing one that has the necessary permissions. This role should have, at a minimum:
 - `AmazonSageMakerFullAccess` policy.
 - Permissions to access the S3 bucket (`lead-classification-bucket`).
 - Permissions to interact with the Redshift Data API.
- Create User Profile:** Within the domain setup, add a default user profile. Give it a name like `ml-developer`.
- Review and Submit:** Review all configurations and click **Submit** to create the Domain. This process can take several minutes.

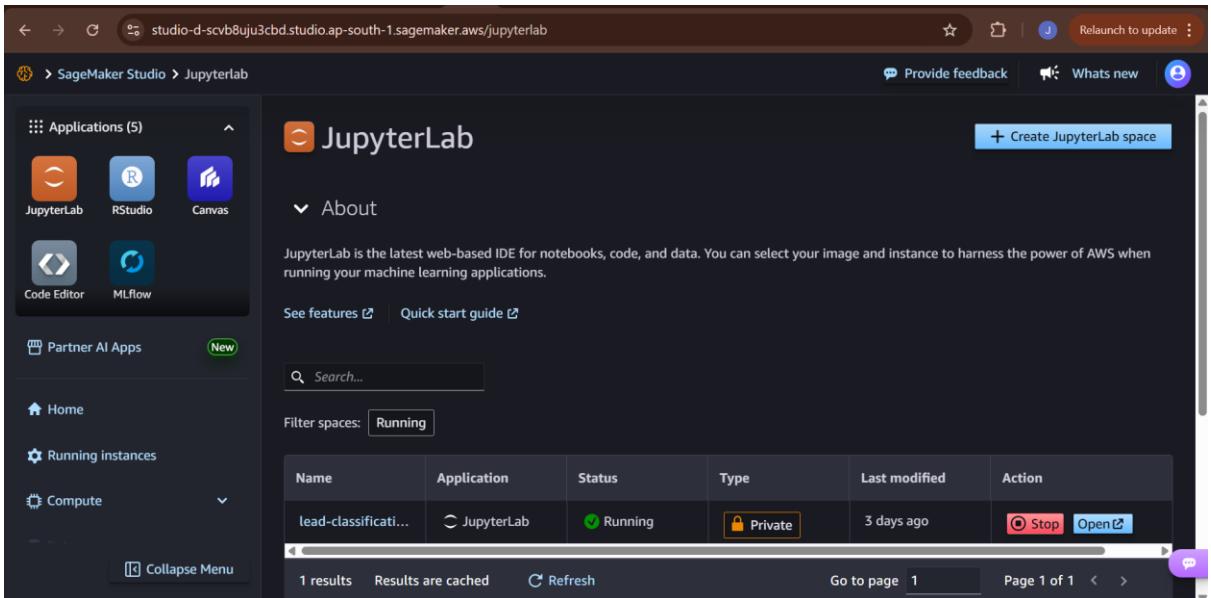
The screenshot shows the AWS Management Console for Amazon SageMaker AI. The URL is `ap-south-1.console.aws.amazon.com/sagemaker/home?region=ap-south-1#/studio/d-scvb8uju3cbd/user/default-20250718T120710`. The page title is "User Details: default-20250718T120710". The sidebar on the left shows "Amazon SageMaker AI" with sections for "Getting started" and "What's new". Under "Applications and IDEs", there are links for "Studio", "Canvas", "RStudio", "Notebooks", and "Partner AI Apps". Under "Admin configurations", there are links for "Domains" and "Role manager". At the bottom of the sidebar are "CloudShell" and "Feedback" buttons. The main content area shows the "User Details" page for the domain "default-20250718T120710". It includes a summary section with "General details about this user profile", a "Details" table with columns for Name, Execution role, and Status, and a "User profile rules" section with a "Manage rules" button. The status is shown as "InService". The execution role is `arn:aws:iam::258399971007:role/sagemaker-role`. The ID is `d-scvb8uju3cbd`. The page footer includes copyright information for 2025, links for "Privacy", "Terms", and "Cookie preferences", and a "Relaunch to update" button.

10.2. Launching SageMaker Studio

Once the Domain status is "Ready," you can launch the Studio IDE for your user.

- From the SageMaker Domain page, find your user profile (e.g., `ml-developer`) in the user list.
- On the right side of the user row, click **Launch** and select **Studio**.

3. A new browser tab will open, and the SageMaker Studio environment will begin to load. This can take a minute or two on the first launch as it provisions the necessary resources.



The screenshot shows the SageMaker Studio JupyterLab interface. On the left, there's a sidebar with icons for JupyterLab, RStudio, Canvas, Code Editor, MLflow, and Partner AI Apps. Below the sidebar are sections for Home, Running instances, Compute, and a Collapse Menu button. The main area is titled "JupyterLab" and contains an "About" section with a brief description of JupyterLab. It also has links for "See features" and "Quick start guide". A search bar and a filter button labeled "Running" are present. A table lists the running JupyterLab space, showing Name (lead-classificati...), Application (JupyterLab), Status (Running), Type (Private), Last modified (3 days ago), and Action (Stop, Open). At the bottom, there are buttons for "1 results", "Results are cached", "Refresh", "Go to page 1", and "Page 1 of 1".

Name	Application	Status	Type	Last modified	Action
lead-classificati...	JupyterLab	Running	Private	3 days ago	Stop Open

10.3. Opening the JupyterLab Notebook

SageMaker Studio provides a fully managed JupyterLab environment. This is where all coding, from data exploration to model training, will take place.

1. **Open Launcher:** Once Studio has loaded, you will be presented with a launcher screen. If not, you can open it from the **File > New Launcher** menu.
2. **Create Notebook:** Under the "Notebooks and compute resources" section, click on **Notebook**.
3. **Select Kernel:** A dialog will pop up asking you to select an image and kernel. Choose the **Data Science** image and the **Python 3** kernel. This provides a Python environment pre-loaded with common data science libraries like pandas, NumPy, and scikit-learn.
4. The new notebook file (e.g., `Untitled.ipynb`) will open, and you are now ready to begin the ML development process.
5. You can directly use your local .py files and make changes according to cloud.

With the SageMaker Studio environment running, the first step in our ML workflow is to load the data from Redshift and begin the exploratory data analysis.

```

1 import boto3
2 import pandas as pd
3 import numpy as np
4 import time
5
6 def load_from_redshift(
7     region='ap-south-1',
8     workgroup_name='lead-workgroup',
9     database_name='dev',
10    secret_arn='arn:aws:secretsmanager:ap-south-1:258399971007:secret:redshift-serverless-secret-sz4Cm0',
11    sql='SELECT * FROM lead_scoring LIMIT 100'
12 ) -> pd.DataFrame:
13     """
14         Executes a SQL query on Redshift Serverless and returns the result as a pandas.DataFrame.
15     """
16     Parameters:
17         - region (str): AWS region
18         - workgroup_name (str): Redshift Serverless workgroup name
19         - database_name (str): Redshift database name
20         - secret_arn (str): ARN of AWS Secrets Manager secret with Redshift credentials
21         - sql (str): SQL query to execute
22
23     Returns:
24         - pd.DataFrame: Query result
25     """

```

```

1 from data_loader import load_from_redshift
2 from data_cleaner import clean_data
3 from feature_engineering import feature_engineering
4 from data_processor import preprocess_data
5 from model_trainer_evaluator_selector import train_evaluate_and_select_model_accuracy, train_evaluate_and_select_model_precision
6 from drift_detector import run_drift
7
8 def run_pipeline():
9     df = load_from_redshift()
10    df_cleaned = clean_data(df)
11    df_featured = feature_engineering(df_cleaned)
12    X_train, X_test, y_train, y_test = preprocess_data(df_featured)
13    best_model_accuracy, y_pred_test_accuracy = train_evaluate_and_select_model_accuracy(X_train, y_train, X_test, y_test)
14    best_model_precision, y_pred_test_precision = train_evaluate_and_select_model_precision(X_train, y_train, X_test, y_test)
15    run_drift()
16
17
18 if __name__ == "__main__":
19     run_pipeline()

```

11. Exploratory Data Analysis (EDA) and Cleaning

11.1 Initial Data Inspection

Before any modeling, the dataset loaded from Redshift is explored to understand its structure, completeness, and potential quality issues. The following steps are conducted:

- Data Preview (df.head())**
View the top few records to understand column structure and sample values.
- Data Types & Nulls (df.info())**
Check for missing values, data types (e.g., float64, object), and shape.

- **Descriptive Statistics (df.describe())**
Analyze numeric columns for mean, standard deviation, min, max, and percentiles.
 - **Missing Value Summary (df.isnull().sum())**
Detect columns with missing entries that need imputation or cleaning.
 - **Duplicate Records (df.duplicated().sum())**
Identify and remove duplicate rows that can bias learning.
 - **Column Name Review (df.columns)**
Validate column naming conventions before standardization.
-

11.2 Data Cleaning and Standardization

Implemented in `clean_data(df)` function. This step ensures that the dataset is free of inconsistencies and missing values are handled systematically:

- **Dropped Columns:**
 - **Identifiers:** 'Prospect ID', 'Lead Number' – unique IDs with no predictive power.
 - **Constant Value Columns:** 'Magazine', 'Receive More Updates About Our Courses', etc.
- **Handled Nulls and Placeholder Values:**
 - Replaced "Select" and "" with "Missing" or appropriate defaults.
 - Filled 'Country' with 'India', 'What is your current occupation' with 'Unemployed', etc.
- **Consolidated and Mapped Labels:**
 - Applied custom mappings to group sparse or noisy categorical values:
 - Lead Source (e.g., 'bing' → 'Other')
 - Specialization (e.g., 'e-business' → 'Business')
 - Tags (e.g., 'ringing', 'busy' → 'Not Reachable')
 - Lead Quality (e.g., 'low in relevance' → 'Low')
 - Lead Profile, How did you hear... are similarly grouped.
- **Ordinal Encoding:**

- Asymmetrique Profile Index, Asymmetrique Activity Index:
'01.High' → 3, '02.Medium' → 2, '03.Low' → 1
 - **Imputation:**
 - For numerical fields like 'TotalVisits', 'Asymmetrique Profile Score':
Used median to fill missing values after coercion to numeric.
 - **Text Normalization:**
 - Lowercased and stripped whitespace from all string cells using `applymap()`.
-

12. Feature Engineering

Executed through the `feature_engineering(df)` function. This step creates derived features using domain insights:

- **Engagement Score**
Weighted combination of:
 - Total Time Spent on Website (40%)
 - Page Views Per Visit (30%)
 - TotalVisits (30%)
- **Combined Asymmetrique Score**
 - Sum of Asymmetrique Activity Score and Asymmetrique Profile Score.
- **Is New Tag**
 - **Binary indicator:** 1 if 'Tags' contains 'student', else 0.
- **High Interaction**
 - Binary feature based on 'Last Activity' values such as 'SMS Sent', 'Email Opened'.
- **Potential Lead**
 - Derived from 'Lead Profile' — 1 if it contains 'potential'.

These features improve signal-to-noise ratio and help models capture behavioral patterns.

13. Preprocessing Pipeline

Implemented via preprocess_data(df). This stage ensures reproducibility and consistency in feature transformation.

13.1 Feature Grouping

- **Numerical:** Website behavior & profile scores
 - 'TotalVisits', 'Total Time Spent on Website', etc.
- **Ordinal:**
 - 'Asymmetrique Activity Index', 'Asymmetrique Profile Index', 'Lead Quality'
- **Categorical:** All remaining object-based columns

13.2 Data Splitting

- Performed an 80/20 stratified split into training and testing sets.
- Ensures balanced distribution of the target variable (Converted).

13.3 Outlier Handling

- Applied Winsorization (1% on both sides) on numerical features to cap extreme values and mitigate outliers' impact.

13.4 Pipeline Composition

- **Numerical Pipeline:**
 - SimpleImputer(strategy='mean')
 - PowerTransformer(method='yeo-johnson')
 - MinMaxScaler
- **Ordinal Pipeline:**
 - SimpleImputer(strategy='most_frequent')
 - OrdinalEncoder with predefined order for each feature
- **Categorical Pipeline:**
 - SimpleImputer(strategy='most_frequent')
 - OneHotEncoder(handle_unknown='ignore')

Combined via ColumnTransformer.

13.5 Saving Pipeline

- The final preprocessor is saved using joblib to preprocessor.pkl.
-

14. Model Training and Evaluation

Two separate pipelines target different business needs:

14.1 Training for Accuracy

Function: train_evaluate_and_select_model_accuracy()

- **Models Used:**
 - Logistic Regression, Random Forest, XGBoost, LightGBM
- **Cross-validation:**
 - 5-fold stratified CV using F1-score as the scoring metric
- **Evaluation Metrics:**
 - Accuracy, Precision, Recall, F1-score, ROC-AUC
- **Model Selection:**
 - Based on best test set accuracy
- **MLflow Integration:**
 - Logs hyperparameters, metrics, model artifact
 - Registers best model to SageMaker-hosted MLflow Registry
 - Promotes model to Staging and then Production

```
Best Model: xgboost
Test Accuracy: 0.9271
Test F1 Score: 0.9075
```

Classification Report:

	precision	recall	f1-score	support
0	0.93	0.95	0.94	955
1	0.92	0.89	0.91	637
accuracy			0.93	1592
macro avg	0.93	0.92	0.92	1592
weighted avg	0.93	0.93	0.93	1592

14.2 Training for Precision (Class 1)

Function: train_evaluate_and_select_model_precision()

- Same models and structure as above.
- Optimizes using precision_score for class 1 to reduce false positives.
- Selected model is also logged and registered in MLflow.

```

Best Model: xgboost
Accuracy: 0.8763
Precision (class 1): 0.9564
F1 Score: 0.8239

Classification Report:
      precision    recall  f1-score   support
0           0.84     0.98    0.90      955
1           0.96     0.72    0.82      637

accuracy                           0.88      1592
macro avg       0.90     0.85    0.86      1592
weighted avg    0.89     0.88    0.87      1592

```

15. Model Explainability with SHAP

SHAP analysis is used to build trust and interpretability in the deployed models:

- **Explainer Type:**
 - **TreeExplainer:** For tree-based models (RF, XGB, LGBM)
 - **LinearExplainer:** For Logistic Regression
- **SHAP Summary Plot:**
 - Shows top features pushing predictions toward 1 (Converted) or 0 (Not Converted)
 - Logged as .png in MLflow
- **Fail-Safe Handling:**
 - If SHAP fails due to model incompatibility or data mismatch, logs warning

16. Data Drift Detection

Performed via run_drift() and log_evidently_report() using Evidently AI.

- **Reference vs. Current Split:**

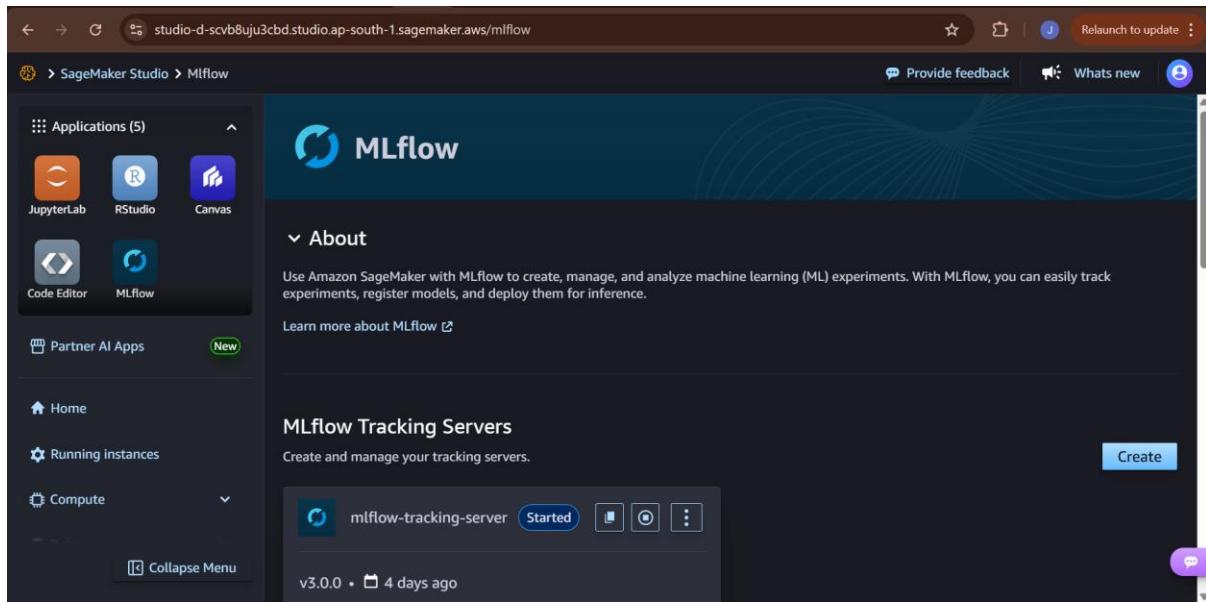
- Train vs. Test sets from initial stratified split
 - **Reports Generated:**
 - DataDriftPreset and DataSummaryPreset
 - Saved as .html and .json
 - Logged to MLflow experiment: *Lead Conversion Prediction Evidently*
 - **Metrics Logged:**
 - Drifted columns count and share
 - Row count and column count
 - Per-feature value drift
 - **Automation:**
 - Automatically aligns common columns and avoids drift failures due to schema mismatches
-

Part 5: MLOps, Deployment & Monitoring

16. Experiment Tracking with MLflow

16.1. Setting Up the MLflow Tracking Server

To centralize experiment tracking, we deploy an MLflow tracking server on a small EC2 instance within our VPC. This provides a persistent, shared location for all developers to log and view experiment results. The server is exposed via its private IP within the VPC, and the security group is configured to allow traffic on port 5000.



16.2. Integrating MLflow into the Training Pipeline

MLflow is seamlessly woven into our `train_pipeline` function.

- `mlflow.set_tracking_uri()`: Points the logging client to our EC2-hosted server.
- `mlflow.set_experiment()`: Organizes runs under a specific experiment name, "Lead Conversion Prediction".
- `with mlflow.start_run()::` Creates a new run context for each model, automatically capturing metadata.
- `mlflow.log_param()` & `mlflow.log_metric()`: Logs hyperparameters and evaluation scores.
- `mlflow.sklearn.log_model()`: Logs the trained model object, its dependencies, and a `conda.yaml` file, making it fully reproducible.
- `mlflow.log_artifact()`: Logs supplementary files like the SHAP plot.

17. Model Registration and Lifecycle Management

17.1. Saving and Registering the Champion Model

After evaluating all models, our pipeline identifies the "champion" model (typically the one with the highest ROC-AUC or F1-score). This model is then prepared for production.

1. **Retrain:** The champion model is retrained on the *entire* dataset (training + validation) to learn from all available data.

- Build Full Pipeline:** The retrained model is combined with the fitted preprocessor into a single, end-to-end `scikit-learn` pipeline. This ensures that raw data can be fed directly into the final artifact for prediction.
- Log and Register:** This full pipeline is logged to a final MLflow run and then registered in the MLflow Model Registry with a specific name (e.g., `Lead-Conversion-Classification-Precision`).

Run Name	Created	Dataset	Duration	Source	Mode
lightgbm_run	18 hours ago	-	8.1s	main.py	
xgboost_run	18 hours ago	-	7.6s	main.py	
random_forest_run	18 hours ago	-	9.6s	main.py	
logistic_regression_run	18 hours ago	-	11.7s	main.py	
lightgbm_run	2 days ago	-	7.9s	main.py	
xgboost_run	2 days ago	-	7.1s	main.py	
random_forest_run	2 days ago	-	9.1s	main.py	

17.2. The MLflow Model Registry

The Model Registry is a central hub for managing the lifecycle of our models. It allows us to:

- Version Models:** Every time we register a model with the same name, a new version is created.
- Manage Stages:** We can assign stages to model versions, such as **Staging**, **Production**, or **Archived**.
- Transition Models:** We can programmatically transition our newly registered model to the "Production" stage, signaling that it is the official version to be used for inference. This is a key step in controlled, automated deployments.

The screenshot shows the MLflow UI for a registered model named "Lead-Classification-Best-Model-Accuracy". The specific version shown is "Version 3". The "Stage" dropdown is set to "Production". The "Schema" section is expanded, showing a table with columns "Name" and "Type". There are two rows: one for "Inputs (0)" and one for "Outputs (0)".

18. Serving and Inference

18.1. Loading the Production Model from the Registry

To ensure our application always uses the correct, approved model, we have a function that automatically discovers and loads the latest "Production" version from the MLflow Model Registry.

- `get_latest_production_model_name()`: This function connects to the MLflow client, searches the registry, and returns the URI of the model currently in the "Production" stage.
- `mlflow.sklearn.load_model()`: This function takes the model URI and loads the complete, deployment-ready pipeline into memory.

The screenshot shows a JupyterLab interface with a terminal tab open. The terminal output shows the command `mlflow sklearn.load_model` being run, which then proceeds to download artifacts from the MLflow registry. The progress bar indicates the download is at 100% completion.

```

View run Lead-Conversion-Prediction-Drift-Detection at: https://ap-south-1.experiments.sagemaker.aws/#/experiments/69/runs/ba6c8df6d755
484a91c64e390c5b403
View experiment at: https://ap-south-1.experiments.sagemaker.aws/#/experiments/69
sagemaker-user@default:~$ python app.py
Traceback (most recent call last):
  File "/home/sagemaker-user/app.py", line 11, in <module>
    from pyngrok import ngrok
ModuleNotFoundError: No module named 'pyngrok'
sagemaker-user@default:~$ pip install pyngrok
Collecting pyngrok
  Using cached pyngrok-7.2.12-py3-none-any.whl.metadata (9.4 kB)
Requirement already satisfied: PyYAML<5.1 in /opt/conda/lib/python3.12/site-packages (from pyngrok) (6.0.2)
Using cached pyngrok-7.2.12-py3-none-any.whl (26 kB)
Installing collected packages: pyngrok
Successfully installed pyngrok-7.2.12
sagemaker-user@default:~$ python app.py
2025-07-22 04:22:52,616 - INFO - Loading model from: models:/Lead-Classification-Best-Model-Precision/Production
Downloading artifacts: 100% [██████████] / 5 / [00:00:00:00, 77.65it/s]
2025-07-22 04:22:54,836 - INFO - Updating auth token for default "config_path" of "ngrok_path": "/home/sagemaker-user/.config/ngrok/ngrok"
2025-07-22 04:22:55,044 - INFO - Opening tunnel named: http-8000-cb17378a-ed41-4ef3-9337-98d740f6b41a
2025-07-22 04:22:55,059 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="no configuration paths supplied"
2025-07-22 04:22:55,060 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="using configuration at default config path" path=/home/sagemake
r-user/.config/ngrok/ngrok.yml
2025-07-22 04:22:55,061 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="open config file" path=/home/sagemaker-user/.config/ngrok/ngrok
.yml err=nil
2025-07-22 04:22:55,206 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="starting web service" obj=web addr=127.0.0.1:4040 allow_hosts=[ ]
2025-07-22 04:22:55,904 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="client session established" obj=tunnels.session
2025-07-22 04:22:55,904 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="tunnel session started" obj=tunnels.session
2025-07-22 04:22:55,915 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="start pg/api/tunnels id=ae068fc445e0@73c

```

18.2. Real-Time Prediction with Flask and Ngrok

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

- **Flask App (`app.py`):**

1. **Initialization:** Sets up the Flask app and connects to the MLflow tracking server.
2. **Model Loading:** On startup, it calls the functions to load the latest production model from the registry.
3. **/predict Route:** Defines an API endpoint that accepts POST requests with lead data (e.g., in CSV or JSON format).
4. **Prediction Logic:** It takes the incoming data, passes it through the loaded pipeline's `.predict()` method, and returns the prediction (0 or 1) in a JSON response.

- **Ngrok Tunnel:**

1. The Flask app is run locally (`python app.py`).
2. The `ngrok http 8080` command is run in a separate terminal.
3. Ngrok creates a secure public URL that tunnels requests directly to the local Flask application, making it accessible from the internet for end-to-end testing.

The screenshot shows a Jupyter Notebook interface with a file tree on the left and a terminal window on the right. The terminal window displays the output of an `ngrok` command, which has started a tunnel to a local Flask application running on port 8080. The output includes logs about session establishment, tunnel creation, and the public URL generated by Ngrok.

```
2025-07-22 04:22:55,059 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="no configuration paths supplied"
2025-07-22 04:22:55,060 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="using configuration at default config path" path=/home/sagemaker-user/.config/ngrok/ngrok.yml
2025-07-22 04:22:55,061 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="open config file" path=/home/sagemaker-user/.config/ngrok/ngrok.yml err=nil
2025-07-22 04:22:55,206 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="starting web service" obj=web addr=127.0.0.1:4040 allow_hosts=[ ]
2025-07-22 04:22:55,904 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="client session established" obj=tunnels.session
2025-07-22 04:22:55,904 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg="tunnel session started" obj=tunnels.session
2025-07-22 04:22:55,915 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=start pgc/api/tunnels id=ae068fc4454e0b73c
2025-07-22 04:22:55,916 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=end pgc/api/tunnels id=ae068fc4454e0b73c status=200 dur=230.535μs
2025-07-22 04:22:55,916 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=start pgc/api/tunnels id=0d3ecc4be5e0b17f
2025-07-22 04:22:55,917 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=end pgc/api/tunnels id=0d3ecc4be5e0b17f status=200 dur=162.795μs
2025-07-22 04:22:55,917 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=start pgc/api/tunnels id=1822e9e178b13c8f
2025-07-22 04:22:55,918 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=end pgc/api/tunnels id=1822e9e178b13c8f status=200 dur=90.973μs
2025-07-22 04:22:55,919 - INFO - t=2025-07-22T04:22:55+0000 lvl=info msg=start pgc/api/tunnel id=8688a19385ff3a7e
2025-07-22 04:22:56,150 - INFO - t=2025-07-22T04:22:56+0000 lvl=info msg="started tunnel" obj=tunnels name=http-8000-cb17378a-ed41-4ef3-9337-98d740f6b41a addr=http://localhost:8000 url=https://c1e3c84258e1.ngrok-free.app
2025-07-22 04:22:56,151 - INFO - t=2025-07-22T04:22:56+0000 lvl=info msg=end pgc/api/tunnels id=8688a19385ff3a7e status=201 dur=231.890028ms
2025-07-22 04:22:56,151 - INFO - Public URL: https://c1e3c84258e1.ngrok-free.app
* Serving Flask app 'app'
* Debug mode: off
2025-07-22 04:22:56,153 - INFO - WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:8000
* Running on http://169.255.255.2:8000
2025-07-22 04:22:56,153 - INFO - Press CTRL+C to quit
```

c1e3c84258e1.ngrok-free.app

Lead Conversion Prediction

Lead Origin

Lead Source

Do Not Email

Do Not Call

Last Activity

Country

c1e3c84258e1.ngrok-free.app

Lead Quality

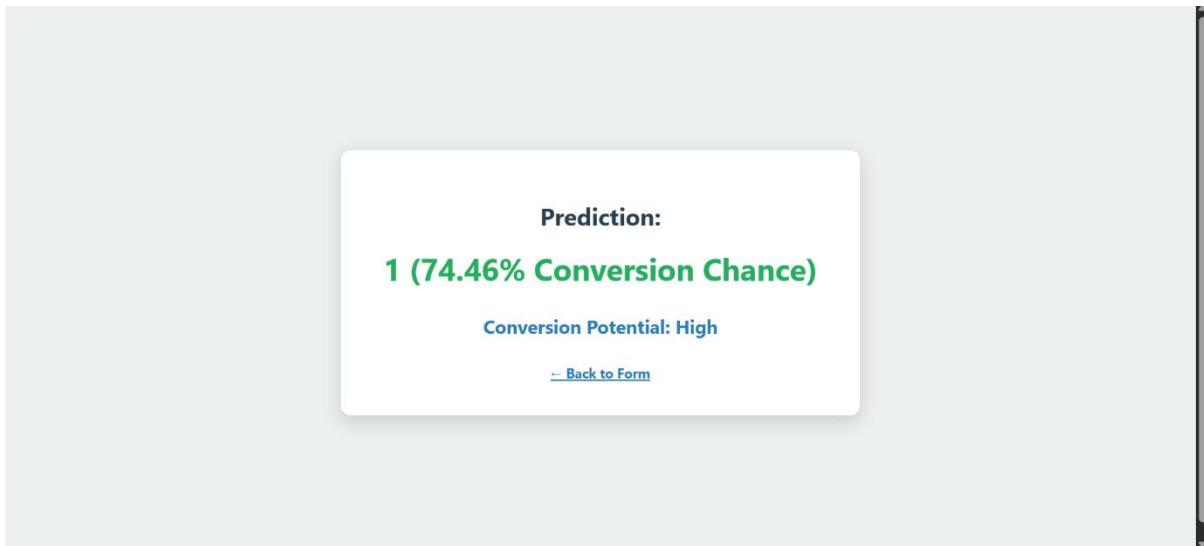
Predict

Upload CSV File

Choose file No file chosen

Upload and Predict

[Download Sample Template](#)



19. Post-Deployment: Drift Detection

19.1. Monitoring for Data Drift with Evidently AI

To ensure our model's performance doesn't degrade over time, we proactively monitor for **data drift**—a change in the statistical properties of the production data compared to the training data.

- **Tool:** We use the **Evidently AI** library, which specializes in generating comprehensive data drift and model performance reports.
- **Process:** We periodically run a script that takes our original training data as a *reference* dataset and a sample of recent production data as the *current* dataset.
- **Report Generation:** Evidently's `DataDriftPreset` generates a detailed report comparing the distributions of every feature between the two datasets.

19.2. Automating Drift Analysis and Logging

This process is automated and integrated with MLflow for tracking.

1. **Start MLflow Run:** A new run is initiated to store the drift analysis results.
2. **Generate Report:** The Evidently drift report is generated in two formats: a visual HTML file and a structured JSON object.
3. **Log HTML Report:** The interactive HTML report is logged as an artifact to MLflow, allowing for detailed visual inspection.
4. **Extract and Log Metrics:** The script parses the JSON report to extract key metrics, such as the `Share of Drifted Columns` and per-feature drift scores. These numerical values are logged as metrics in MLflow, making it easy to track drift trends over time and set up automated alerts.

20. Future State: Full Automation with MWAA

20.1. Orchestrating the MLOps Lifecycle with Airflow

To achieve full, hands-off automation, the entire MLOps workflow is orchestrated using **Amazon Managed Workflows for Apache Airflow (MWAA)**. MWAA is a managed service that makes it easy to run Apache Airflow on AWS. The entire automated process—from data drift detection to model retraining and deployment—is defined as a **Directed Acyclic Graph (DAG)**, which is a collection of all the tasks you want to run, organized in a way that reflects their relationships and dependencies.

20.2. Setting up the MWAA Environment

Before deploying the DAG, you must set up the MWAA environment and its required resources.

Step 1: Create an S3 Bucket for MWAA Assets MWAA uses an S3 bucket to store your DAGs, custom plugins, and Python dependencies (`requirements.txt`).

1. Create a new S3 bucket with a unique name (e.g., `mwaa-lead-conversion-dags`).
2. Enable **versioning** on this bucket to keep a history of your DAGs and configuration files.

Step 2: Configure the VPC and Networking MWAA requires a specific network setup to run securely and access other resources.

1. **Use the same VPC:** The MWAA environment must be launched in the same VPC as your Redshift, SageMaker, and MLflow resources to communicate with them privately via the VPC endpoints.

2. **Subnets:** You must provide at least two **private subnets** in different Availability Zones for high availability.
3. **Internet Access:** The private subnets must have a route to a **NAT Gateway**. This is required for the Airflow workers to install Python dependencies from public repositories like PyPI.
4. **Security Group:** The MWAA environment should use the same `mlops-sg` security group created in Part 2. This ensures it has the necessary inbound and outbound rules to communicate with the MLflow server, Redshift, S3, and other services.

Step 3: Create the MWAA Environment

1. Navigate to the **Managed Workflows for Apache Airflow** service in the AWS Console.
2. Click **Create environment**.
3. **Details:**
 - o **Name:** MyAirflowEnvironment
 - o **Airflow version:** Choose a recent, stable version (e.g., 2.8.1).
 - o **S3 bucket:** Browse and select the S3 bucket you created in Step 1.
 - o **DAGs folder path:** `dags`
 - o **Requirements file:** `requirements.txt`
4. **Networking:** Select your VPC, the private subnets, and the `mlops-sg` security group.
5. **Execution role:** Create a new IAM role that grants MWAA permissions to access its S3 bucket and CloudWatch Logs. Attach policies that also allow it to interact with SageMaker, Redshift, and any other services your DAG will orchestrate.
6. Click **Create environment**. Provisioning can take up to 30 minutes.

The screenshot shows the AWS MWAA Environment details page for 'MyAirflowEnvironment'. The page has a dark theme. At the top, there's a navigation bar with the AWS logo, search bar, and user information. Below the navigation, the breadcrumb trail shows 'Amazon MWAA > Environments > MyAirflowEnvironment'. The main content area has a title 'MyAirflowEnvironment' with 'Edit', 'Delete', and 'Open Airflow UI' buttons. On the left, there's a 'Details' section with 'Status' (Available), 'ARN' (arn:aws:airflow:ap-south-1:25839971007:environment/MyAirflowEnvironment), and a link to the 'Airflow UI'. On the right, there's a 'Weekly maintenance window start (UTC)' set for Saturday 19:30. At the bottom, there's a 'Last update' section showing 'Status' (SUCCESS) and 'Created at' (not specified). The footer includes links for CloudShell, Feedback, and various AWS terms like Privacy, Terms, and Cookie preferences.

20.3. Airflow Workflow Lead Conversion Retraining

Structure the Project S3 Folder: On your S3 bucket, create the following folder structure:

airflow-lead23/

```
|   └── dags/
|       └── lead_retrain_dag.py      ← Apache Airflow DAG
|
|   └── requirements/
|       └── requirements.txt      ← Python packages required by MWAA and our code
|
|   └── incoming/
|       └── new_data.csv      ← New uploaded dataset to trigger drift detection
|
|   └── scripts/
|       └── main.py      ← Contains entire ML pipeline (load → clean → train → register)
```

The screenshot shows the AWS S3 console interface. The left sidebar displays the 'Amazon S3' navigation pane with sections for General purpose buckets, Storage Lens, and CloudShell. The main content area shows the 'Objects' tab for the 'airflow-lead23' bucket. The 'Objects' section lists six items: 'dags/' (Folder), 'incoming/' (Folder), 'raw/' (Folder), 'requirements.txt' (txt file), 'requirements/' (Folder), and 'scripts/' (Folder). The 'requirements.txt' file is highlighted, showing its last modified date as July 21, 2025, at 21:55:59 (UTC+05:30), a size of 373.0 B, and a storage class of Standard.

Name	Type	Last modified	Size	Storage class
dags/_	Folder	-	-	-
incoming/_	Folder	-	-	-
raw/_	Folder	-	-	-
requirements.txt	txt	July 21, 2025, 21:55:59 (UTC+05:30)	373.0 B	Standard
requirements/_	Folder	-	-	-
scripts/_	Folder	-	-	-

Automate an end-to-end retraining pipeline that:

1. Detects data drift between incoming lead data (from S3) and reference data (from Redshift).

2. Retrains the lead conversion model (main.py in S3) only if drift > 0.3.
3. Skips retraining if drift is within acceptable threshold (≤ 0.3).

Architecture Overview:

Component	Technology Used
Workflow Engine	Apache Airflow (MWAA)
Cloud Provider	AWS
Data Storage	S3 (for incoming datasets and main.py)
Data Warehouse	Amazon Redshift Serverless
Drift Detection	Evidently AI
Model Training	Python (main.py script via BashOperator)

Data Sources

New Data

- Source: S3://airflow-lead23/incoming/new_data.csv
- Loaded by: boto3.download_file

Reference Data

- Source: Amazon Redshift table: lead_scoring
- Loaded by: boto3.client('redshift-data')

DAG: lead_conversion_retrain

A DAG scheduled daily to check for drift and conditionally retrain the model.

Workflow Steps

Step 1: Install Python Dependencies

- Operator: BashOperator
- Command: pip install -r /usr/local/airflow/requirements/requirements.txt

Step 2: Check for Data Drift

- Operator: PythonOperator
- Task: check_drift
- Performs the following:
 - Downloads new data from S3
 - Pulls reference data from Redshift using SQL
 - Renames and cleans columns
 - Runs Evidently's DataDriftPreset
 - Calculates drift score and pushes it to XCom

Step 3: Branch on Drift Score

- Operator: BranchPythonOperator
- Decides to retrain or skip based on drift threshold (0.3)

Step 4A: Retrain Model (If Drift > 0.3)

- Operator: BashOperator
- Command: aws s3 cp s3://airflow-lead23/scripts/main.py ./ && python main.py

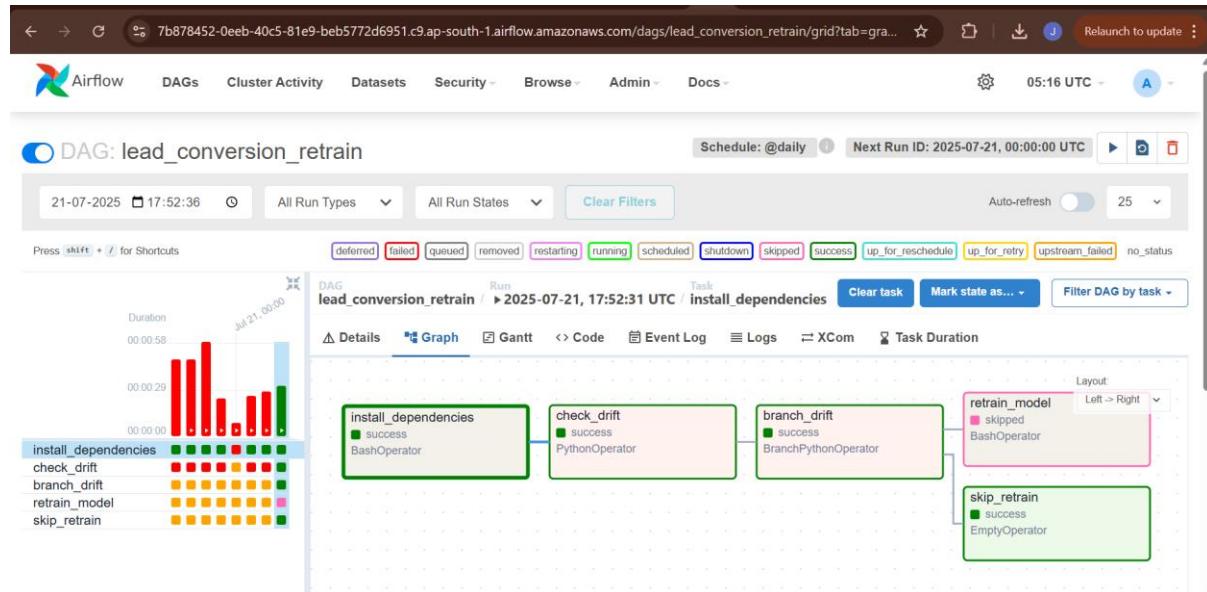
Step 4B: Skip Retraining (If Drift ≤ 0.3)

- Operator: DummyOperator
- Ends DAG branch with no retraining

Python Dependencies

- evidently >= 0.4.14
- boto3
- pandas
- numpy
- psycopg2

These are listed in /usr/local/airflow/requirements/requirements.txt



21. SageMaker Model Deployment

To deploy a trained lead conversion model as a SageMaker real-time inference endpoint using a custom script and required artifacts including model, preprocessor, and data transformation logic.

21.1 Folder Structure (Used in SageMaker)

All files are placed in the working directory: lead_model_deploy/.

```
lead_model_deploy/
├── inference.py          # Inference logic
├── data_cleaner.py        # Data preprocessing
└── feature_engineer.py    # Feature generation logic
```

```
└── best_model.pkl      # Trained model (pickle)
└── preprocessor.pkl   # Transformation pipeline (pickle)
└── requirements.txt    # Python dependency list
└── sagemaker_endpoint_trail.ipynb # Notebook to trigger deployment
└── model_bundle.tar.gz  # Archive uploaded to S3
```

21.2 Model Packaging and Upload

Compress the artifacts into a .tar.gz:

```
tar -czvf model_bundle.tar.gz inference.py data_cleaner.py feature_engineer.py  
best_model.pkl preprocessor.pkl
```

21.3 Upload the model bundle to S3:

```
aws s3 cp model_bundle.tar.gz s3://lead-classification-bucket/model_bundle.tar.gz
```

S3 Path: s3://lead-classification-bucket/model_bundle.tar.gz

21.4 SageMaker Endpoint Creation via Notebook

Using the sagemaker_endpoint_trail.ipynb, the model was deployed with the SKLearnModel class.

```
from sagemaker.sklearn.model import SKLearnModel  
from sagemaker import get_execution_role, Session  
  
role = get_execution_role()  
  
model = SKLearnModel(  
    model_data="s3://lead-classification-bucket/model_bundle.tar.gz",  
    role=role,  
    entry_point="inference.py",  
    framework_version="1.2-1",  
    sagemaker_session=Session()  
)  
  
predictor = model.deploy(  
    instance_type="ml.m5.large",  
    initial_instance_count=1,  
    endpoint_name="lead-conversion-endpoint"  
)
```

21.5. Deployment Result

Field	Value
Endpoint Name	lead-conversion-endpoint
Region	ap-south-1
Model Bundle	model_bundle.tar.gz
Status	Failed

The screenshot shows a Jupyter Notebook interface with the following details:

- Header:** The URL is `x7ulny7kzheffhz.studio.ap-south-1.sagemaker.aws/jupyterlab/default/lab/tree/lead_model_deploy/sagemaker_endpoint_trail... Relaunch to update`.
- File Menu:** File, Edit, View, Run, Kernel, Git, Tabs, Settings, Help.
- Toolbar:** Includes icons for file operations like +, -, &, %, and a search bar.
- Left Sidebar:** Shows a file tree with the path `/ lead_model_deploy /`. Files listed include `best_model.pkl`, `data_cleaner.py`, `feature_engineer.py`, `inference.py`, `model_bundle.tar...`, `model_bundle.zip`, `preprocessor.pkl`, `requirements.txt`, and `sagemaker_endpo...`.
- Code Editor:** Displays two code snippets from `sagemaker.py`:
 - Line 4864: `self.endpoint_arn = res["EndpointArn"]`
 - Line 4865: `if wait:`
 - Line 4866: `self.wait_for_endpoint(endpoint_name, live_logging=live_logging)`
 - Line 4867: `return endpoint_name`
 - Line 4868: `except Exception as e:`
 - Line 4869: `troubleshooting = (`
- Line 5673: `allowed_statuses=["InService"],`
- Line 5674: `actual_status=status,`
- Line 5675: `)`
- Line 5676: `raise exceptions.UnexpectedStatusException(`
- Line 5677: `message=message,`
- Line 5678: `allowed_statuses=["InService"],`
- Line 5679: `actual_status=status,`

- Output Panel:** Shows an error message:

`UnexpectedStatusException: Error hosting endpoint lead-conversion-endpoint: Failed. Reason: The primary production variant AllTraffic did not pass the ping health check. Please check Cloudwatch logs endpoint... Try changing the instance type or reference the troubleshooting page https://docs.aws.amazon.com/sagemaker/latest/dg/async-inference-troubleshooting.html`
- Bottom Status Bar:** Shows the kernel status as "Fully initialized", memory usage at 55%, and the current cell number as 1 in 1.

The screenshot shows a Jupyter Notebook interface with the following details:

- File Explorer:** On the left, it lists files and folders under the path `/lead_model_deploy/`. The files listed are: `best_model.pkl`, `data_cleaner.py`, `feature_engineer.py`, `inference.py`, `model_bundle.tar...`, `model_bundle.zip`, `preprocessor.pkl`, and `requirements.txt`.
- Code Editor:** The main area displays a terminal session with the following command history:

```
sagemaker-user@default:~/lead_model_deploy$ tar -czvf model_bundle.tar.gz inference.py data_cleaner.py feature_engineer.py best_model.pkl
preprocessor.pkl
inference.py
data_cleaner.py
feature_engineer.py
best_model.pkl
preprocessor.pkl
sagemaker-user@default:~/lead_model_deploy$ aws s3 cp model_bundle.tar.gz s3://lead-classification-bucket/
upload: ./model_bundle.tar.gz to s3://lead-classification-bucket/model_bundle.tar.gz
sagemaker-user@default:~$ python --version
Python 3.12.9
sagemaker-user@default:~$ cd lead_model_deploy
zip -r model_bundle.zip inference.py model/
utils/
bash: zip: command not found
sagemaker-user@default:~/lead_model_deploy$ sudo yum install -y zip
sudo: yum: command not found
sagemaker-user@default:~/lead_model_deploy$ sudo apt-get update && sudo apt-get install -y zip
Get:1 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]
Hit:2 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:3 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]
Get:4 http://security.ubuntu.com/ubuntu jammy-security/universe amd64 Packages [1267 kB]
Get:5 http://security.ubuntu.com/ubuntu jammy-security/restricted amd64 Packages [4932 kB]
Get:6 http://archive.ubuntu.com/ubuntu jammy-backports InRelease [127 kB]
Get:7 http://security.ubuntu.com/ubuntu jammy-security/main amd64 Packages [3148 kB]
Get:8 http://archive.ubuntu.com/ubuntu jammy-updates/restricted amd64 Packages [5139 kB]
Get:9 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages [1572 kB]
Get:10 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [3461 kB]
Get:11 http://archive.ubuntu.com/ubuntu jammy-updates/multiverse amd64 Packages [75.9 kB]
Get:12 http://archive.ubuntu.com/ubuntu jammy-backports/universe amd64 Packages [35.2 kB]
Get:13 http://archive.ubuntu.com/ubuntu jammy-backports/main amd64 Packages [83.2 kB]
```
- Toolbar:** At the top, there are standard browser-style navigation buttons (Back, Forward, Home) and a tab bar showing multiple open notebooks.
- Header:** The header includes the URL `x7ulny7kzheffhz.studio.ap-south-1.sagemaker.aws/jupyterlab/default/lab/tree/lead_model_deploy` and a "Relaunch to update" button.

The screenshot shows the Amazon SageMaker AI interface. On the left, there's a sidebar with links like 'Getting started', 'What's new', 'Applications and IDEs' (Studio, Canvas, RStudio, Notebooks, Partner AI Apps), and 'Admin configurations' (Domains, Role manager). The main area is titled 'Endpoints' and lists a single entry:

Name	ARN	Creation time	Status	Last updated
lead-conversion-endpoint	arn:aws:sagemaker:ap-south-1:25839971007:endpoint/lead-conversion-endpoint	7/20/2025, 1:29:31 PM	Failed	7/20/2025, 1:52:57 PM

At the bottom right of the main area, it says '© 2025, Amazon Web Services, Inc. or its affiliates.' and has links for 'Privacy', 'Terms', and 'Cookie preferences'.

21.6 Root Cause Analysis of Failure

The endpoint creation failed. The probable reasons is:

- Not passing health checks
-

Part 6: Appendix

22. Technology & Libraries Stack

- **Cloud Platform:** Amazon Web Services (AWS)
- **Data Storage:** Amazon S3, Amazon Redshift
- **Data Integration (ETL):** AWS Glue (Studio, Crawlers, Data Catalog)
- **Security & Networking:** AWS IAM, AWS VPC, AWS Secrets Manager, AWS KMS
- **ML Development:** Amazon SageMaker Studio Lab
- **MLOps & Tracking:** MLflow
- **Orchestration (Future):** Managed Workflows for Apache Airflow (MWAA)
- **Serving:** Flask, Ngrok
- **Core Python Libraries:**
 - **Data Handling:** pandas, numpy
 - **Visualization:** matplotlib, seaborn
 - **ML & Preprocessing:** scikit-learn
 - **Advanced Models:** xgboost, lightgbm
 - **Model Explainability:** shap
 - **Data Drift:** evidently

- **AWS Interaction:** boto3
- **Serialization:** joblib

23. Resources

- **Official Documentations:**
 - AWS Glue
 - Amazon SageMaker
 - MLflow
 - Evidently AI

24. GitHub Repository

- **GitHub link:** <https://github.com/Rithwik-Minfy/Lead-Conversion-Classification-Capstone>