Name: Harshit Hemant Gupta

IU ID: 2001096665

Subject: FA24-BL-ENGR-E534 BIG DATA APPLICATIONS

**Mini Project: Fall24**

1. ## Dataset Selection (0 Points)
   ● **Task:** Choose a dataset of your own choice that is large enough to derive meaningful insights.
   ○ **Examples:** Financial transactions, customer behavior, IoT sensor data.
   ○ Ensure the dataset contains diverse columns to allow for cleaning, transformation, and aggregation.

   **California Independent Medical Review**

   https://www.kaggle.com/datasets/prasad22/ca-independent-medical-review/data

Description:

This dataset, sourced from the California Department of Managed Health Care (DMHC), includes all decisions from Independent Medical Reviews (IMRs) conducted by the DMHC since January 1, 2001. IMRs are impartial evaluations of health care services that have been denied, delayed, or modified by a health plan, often on grounds of being deemed unnecessary, experimental, or non-urgent. When an IMR decision favors the enrollee, the health plan is required to approve the requested treatment or service.

File Information:

•      File Format: CSV

•      Number of Rows: 19246

•      Number of Columns: 11

Reference | Report Ye | Diagnosis | Diagnosis | Treatment | Treatment | Determina | Type | Age Range | Patient Ge | Findings
---|---|---|---|---|---|---|---|---|---|---
MN16-22€ | 2016 | Infectious | Hepatitis | Pharmacy, | Anti-virals | Overturne | Medical N | 41-50 | Male | Nature of Statutory Criteria/Case Summary: An enrollee has requested Harvoni for treatment of H
MN16-22€ | 2016 | Mental | Eating Dis | Mental He | Residentia | Upheld De | Medical N | 21-30 | Female | Nature of Statutory Criteria/Case Summary:  An enrollee has requested  residential treatment cen
MN16-22€ | 2016 | Autism Sp | Autism-PC | Autism Re | Speech Th | Upheld De | Medical N | 0-10 | Female | Nature of Statutory Criteria/Case Summary:  The parent of an enrollee has requested speech ther
EI16-2263 | 2016 | Prevention/Good He | Diagnostic | Mammog | Overturne | Experimer | 65+ | Female | Nature of Statutory Criteria/Case Summary: An enrollee has requested breast tomosynthesis for
EI06-5319 | 2006 | Cardiac/Circulatory | Cardio Vascular | | Upheld De | Experimer | 51-64 | Male | Physician 1: The patient is a 62-year-old male who is reported to have small vessel disease, not an
EI16-2263 | 2016 | Prevention/Good He | Diagnostic | Lab Work | Upheld De | Experimer | 21-30 | Male | Nature of Statutory Criteria/Case Summary: An enrollee has requested advanced lipoprotein test
EI16-2263 | 2016 | OB-Gyn/ F | Female Br | Diagnostic | Mammog | Overturne | Experimer | 65+ | Female | Nature of Statutory Criteria/Case Summary: An enrollee has requested breast tomosynthesis for
EI16-2263 | 2016 | OB-Gyn/ F | Female Br | Diagnostic | Mammog | Overturne | Experimer | 51-64 | Female | Nature of Statutory Criteria/Case Summary: An enrollee has requested breast tomosynthesis for
MN16-22€ | 2016 | Autism Sp | Autism-PC | Autism Re | Speech Th | Upheld De | Medical N | 0-10 | Female | Nature of Statutory Criteria/Case Summary:  The parent of an enrollee has requested occupation
EI16-2215 | 2016 | Digestive ! | Other | Diagnostic | Allergy Te | Upheld De | Experimer | 11_20 | Female | Nature of Statutory Criteria/Case Summary:  The parent of an enrollee has requested blood aller
EI16-2263 | 2016 | OB-Gyn/ F | Female Br | Diagnostic | Mammog | Overturne | Experimer | 51-64 | Female | Nature of Statutory Criteria/Case Summary: An enrollee has requested breast tomosynthesis for
EI16-2263 | 2016 | Prevention/Good He | Diagnostic | Mammog | Overturne | Experimer | 51-64 | Female | Nature of Statutory Criteria/Case Summary: An enrollee has requested breast tomosynthesis for
EI16-2262 | 2016 | Orthopedi | Fracture | Durable M | Other | Upheld De | Experimer | 41-50 | Male | Nature of Statutory Criteria/Case Summary: An enrollee has requested the ReWalk personal exos
MN16-22€ | 2016 | Mental | Depressio | Mental He | Acute Psy | Upheld De | Medical N | 51-64 | Female | Nature of Statutory Criteria/Case Summary:  An enrollee has requested inpatient mental health s
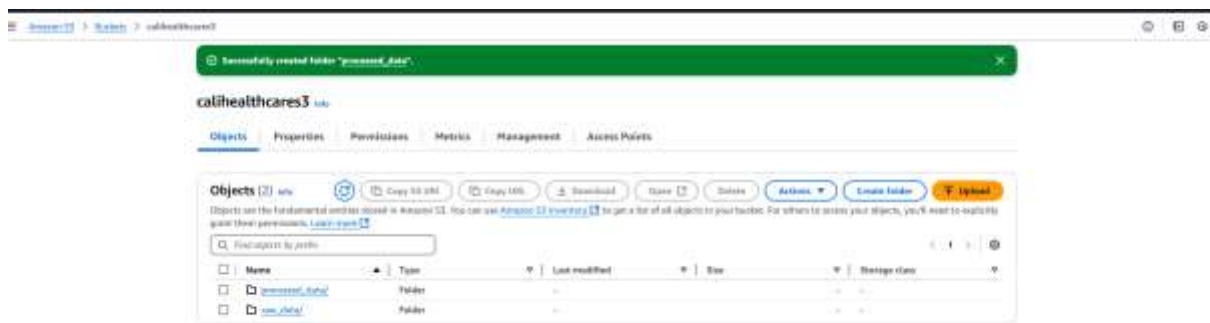
Independent Medical Reviews

**Data Schema:**

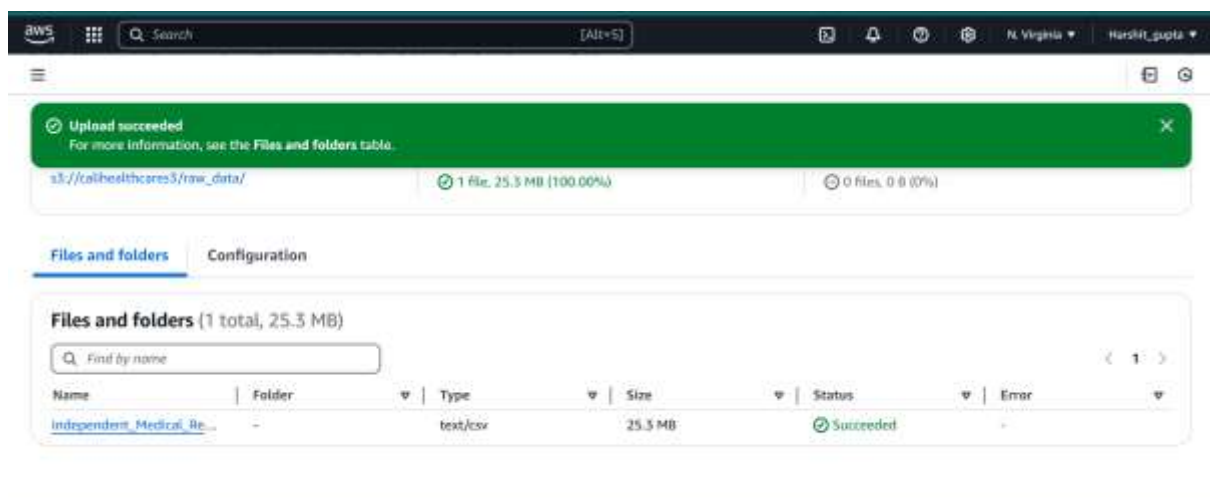| Field Name | Description | Data type |
|---|---|---|
| Reference ID | Unique identifier for the case | Plain Text |
| Report Year | Year the case was reported | Number |
| Diagnosis Category | The main diagnosis category | Plain Text |
| Diagnosis Sub Category | The secondary diagnosis category | Plain Text |
| Treatment Category | The main treatment category | Plain Text |
| Treatment Sub Category | The secondary treatment category | Plain Text |
| Determination | Indicates if the determination was upheld or overturned | Plain Text |
| Type | Indicates the type of case (Experimental/Investigational, Urgent Care, Medical Necessity) | Plain Text |
| Age Range | The age of the patient | Number |
| Patient Gender | The gender of the patient | Plain Text |
| Findings | A summary of the case findings | Plain Text |

## 2. Environment Setup (2.5 points)

☐ **AWS S3 for Data Storage**: Demonstrates the creation and configuration of an AWS S3 bucket, showcasing the process of uploading raw data into S3 to serve as a centralized storage for further data processing.

☐ **Linux Environment with PySpark**: Highlights the setup of a Linux-based environment (e.g., an AWS EC2 instance), installation of PySpark, and configuration of AWS CLI to enable seamless interaction with the S3 bucket for distributed data processing tasks.

### 1. AWS S3 for Data Storage (1 Point)

**a. Step 1:** Create an S3 bucket to store both raw and processed data.



**b. Step 2:** Upload the raw dataset to the S3 bucket.

## 2. Linux Environment with PySpark (1.5 Point)

**a. Step 1:** Set up a Linux-based environment, either locally or using an AWS EC2 instance. **Recommendation:** Use an AWS EC2 instance for better scalability and AWS integration.





**b. Step 2:** Install PySpark for distributed data processing.

**c. Step 3:** Configure AWS CLI to interact with S3 buckets.

# 3. Data Pipeline Tasks (6 points)

**Task 1: Data Ingestion from S3 (1 Point)**

● **Objective:** Pull raw data from S3 into the PySpark environment.

● **Steps:** 1. Use AWS CLI or PySpark's built-in S3 support to load the dataset directly.





2. Confirm successful ingestion by inspecting the dataset.

**Image**: Displays the raw data being pulled from an AWS S3 bucket into the PySpark environment.

**Explanation**: The image should show the use of either PySpark's spark.read functionality or the AWS CLI to load the dataset into the working environment. A preview of the successfully ingested data

confirms that the process was completed without errors, ensuring all data fields and rows are intact.

## Task 2: Data Processing with PySpark (2 Point)

All necessary techniques for cleaning data, such as transforming, aggregating, and removing outliers, should be applied.

1. **Data Transformation:** Create at least 2 new columns (e.g., Year , Month ) to aid in analysis.



2. **Data Aggregation:** Compute at least 5 key metrics, such as Total revenue by region, Monthly spending trends, Top 10 customers by transaction value (These metrics may vary depending on the specifics of your dataset, but the goal is to aggregate the data in ways that enable meaningful analysis and decision-making.)

### *Total cases by Report Year*

## *Distribution of Determination*



## *Cases by Diagnosis Category*



## *Gender-based distribution of cases*

## Age-range analysis



## Diagnosis and Treatment Summary

Explanation: Each image captures a different stage of processing—cleaning, transforming, or aggregating—depicting how raw data is converted into a structured and meaningful format suitable for analysis.

## Task 3: Store Processed Data Back to S3 (0.5 Point)

● **Objective:** Save the processed and aggregated data to a new S3 bucket or folder.

● **Steps:** 1. Export data in CSV or Parquet format.

The image should illustrate the use of PySpark's write.csv or write.parquet commands and subsequent confirmation of the file upload in the S3 bucket. This ensures the processed data is stored securely and is accessible for later analysis.

2. Upload the processed data to a designated S3 location for easy access. (Generating new CSV, downloading it and uploading to S3 bucket is fine).



**Task 4: Data Analysis Using Spark SQL (1 Point)**

• **Objective:** Use SQL to derive insights (Atleast 5 Queries).

• **Example Queries:** 1. Identify top-performing regions. 2. Analyze month-over-month revenue growth. 3. Determine the most popular product categories.

```
24/12/07 23:53:32 INFO TaskSetManager: Finished task 0.0 in stage 4.0 (TID 3) in 249 ms on ip-172-31-81-146.ec2.internal (executor driver) (1/1)
24/12/07 23:53:32 INFO TaskSchedulerImpl: Removed TaskSet 4.0, whose tasks have all completed, from pool
24/12/07 23:53:32 INFO DAGScheduler: ResultStage 4 (showString at NativeMethodAccessorImpl.java:0) finished in 0.282 s
24/12/07 23:53:32 INFO DAGScheduler: Job 3 is finished. Cancelling potential speculative or zombie tasks for this job
24/12/07 23:53:32 INFO TaskSchedulerImpl: Killing all running tasks in stage 4: Stage finished
24/12/07 23:53:32 INFO DAGScheduler: Job 3 finished: showString at NativeMethodAccessorImpl.java:0, took 0.313829 s
24/12/07 23:53:32 INFO CodeGenerator: Code generated in 18.219701 ms
24/12/07 23:53:32 INFO CodeGenerator: Code generated in 19.891384 ms

+---------+----------+
|Age Range|TotalCases|
+---------+----------+
|    51-64|      5517|
|    41-50|      2824|
|    31-40|      1641|
|    11_20|      1634|
|     0-10|      1353|
+---------+----------+
```

i-0e80de008aa41d12d (CalihealthEC2)                                    ✕
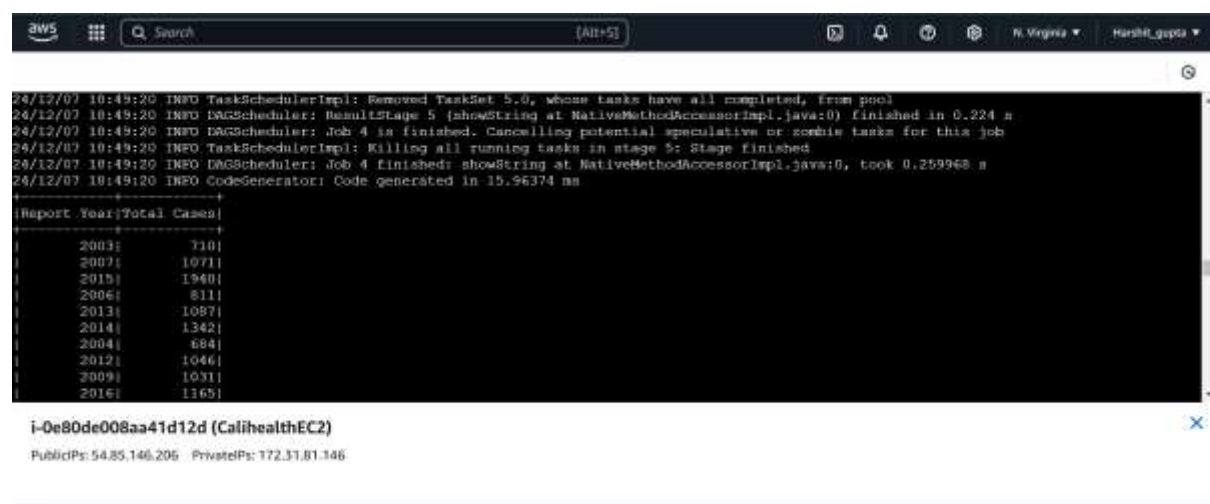
PublicIPs: 3.82.198.54   PrivateIPs: 172.31.81.146

```
24/12/07 23:53:34 INFO TaskSetManager: Finished task 0.0 in stage 7.0 (TID 5) in 94 ms on ip-172-31-81-146.ec2.internal (executor driver) (1/1)
24/12/07 23:53:34 INFO DAGScheduler: ResultStage 7 (showString at NativeMethodAccessorImpl.java:0) finished in 0.127 s
24/12/07 23:53:34 INFO TaskSchedulerImpl: Removed TaskSet 7.0, whose tasks have all completed, from pool
24/12/07 23:53:34 INFO DAGScheduler: Job 5 is finished. Cancelling potential speculative or zombie tasks for this job
24/12/07 23:53:34 INFO TaskSchedulerImpl: Killing all running tasks in stage 7: Stage finished
24/12/07 23:53:34 INFO DAGScheduler: Job 5 finished: showString at NativeMethodAccessorImpl.java:0, took 0.151954 s
24/12/07 23:53:34 INFO CodeGenerator: Code generated in 15.673697 ms

+----------------------------+---------+
|Diagnosis and Treatment Summary|CaseCount|
+----------------------------+---------+
|          Mental - Mental H...|     1514|
|          Infectious - Phar...|      924|
|           Cancer - Cancer T...|      681|
|         Orthopedic/ Muscu...|      611|
|         Orthopedic/ Muscu...|      523|
+----------------------------+---------+

Gender-based Distribution of Cases:
```

i-0e80de008aa41d12d (CalihealthEC2)                                    ✕

PublicIPs: 3.82.198.54   PrivateIPs: 172.31.81.146

```
1 and 0 (0.0 B) push-merged-local and 0 (0.0 B) remote blocks
24/12/07 23:53:35 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 3 ms
24/12/07 23:53:35 INFO Executor: Finished task 0.0 in stage 10.0 (TID 7). 5223 bytes result sent to driver
24/12/07 23:53:35 INFO TaskSetManager: Finished task 0.0 in stage 10.0 (TID 7) in 49 ms on ip-172-31-81-146.ec2.internal (executor driver) (1/1)
24/12/07 23:53:35 INFO TaskSchedulerImpl: Removed TaskSet 10.0, whose tasks have all completed, from pool
24/12/07 23:53:35 INFO DAGScheduler: ResultStage 10 (showString at NativeMethodAccessorImpl.java:0) finished in 0.072 s
24/12/07 23:53:35 INFO DAGScheduler: Job 7 is finished. Cancelling potential speculative or zombie tasks for this job
24/12/07 23:53:35 INFO TaskSchedulerImpl: Killing all running tasks in stage 10: Stage finished
24/12/07 23:53:35 INFO DAGScheduler: Job 7 finished: showString at NativeMethodAccessorImpl.java:0, took 0.086858 s

+--------------+----------+
|Patient Gender|TotalCases|
+--------------+----------+
|        Female|      8436|
|          Male|      6383|
+--------------+----------+

24/12/07 23:53:35 INFO BlockManagerInfo: Removed broadcast_12_piece0 on ip-172-31-81-146.ec2.internal:40879 in memory (size: 21.3 KiB, free: 413.8 MiB)
```

i-0e80de008aa41d12d (CalihealthEC2)                                    ✕

PublicIPs: 3.82.198.54   PrivateIPs: 172.31.81.146

Each image depicts a complete SQL workflow—from formulating queries to interpreting results. This showcases how Spark SQL facilitates powerful, SQL-based data exploration.

**Task 5: Machine Learning with AWS SageMaker Autopilot (1.5 Point)**

● **Objective:** Use **AWS SageMaker Autopilot** to train and evaluate machine learning models on the **processed dataset** stored in S3 without writing any code.

● Here is a tutorial on how to use AWS SageMaker AutoPilot

● **Steps:** 1. **Import Processed Data** : Load the processed dataset from S3 into SageMaker Autopilot.

2. **Run Autopilot Experiment** : ■ Select the target column. ■ Run the AutoML process to train and evaluate multiple models.

3. Consider taking screenshots of your working and query results.

4. **Review Results** : Analyze the model leaderboard and performance metrics. Address ethical issues like bias in training data and privacy concerns.

## Import tabular data                                                    ✕

ⓘ If your data has special character delimiters, use the advanced import settings to specify a custom delimiter. Learn More

part-00000-64de88d7-7ba6-494b-9c33-1dde6adf3583-c000.csv  ▾        🗑 Delete

| Report Year | Diagnosis Sub ... | Treatment Sub... | Determina |
|---|---|---|---|
| 2016 | Hepatitis | Anti-virals | Overtu |
| 2016 | Eating Disorder | Residential Treatment Center - Ad... | Upheld |
| 2016 | Autism-PDD-NOS | Speech Therapy | Upheld |

**Import settings**

Settings apply to all imported files. Learn more ⎘

Dataset name *

part-00000-64de88d7-7ba6-494b-9c33-1

**Sampling**

Sample your dataset for faster exploration. Your full dataset will be used for data export or model build. Learn more ⎘

Sampling method * ⓘ

Random                    ▾

Cancel    Back    **Import**

---

**My models** > Model_20241209_105458 > **Version 1**      + Create new version  ⟳  ⋮

| Select | **Build** | Analyze | Predict | Deploy |

**Select a column to predict**

Choose the target column. The model that you build predicts values for the column that you select.

Target column
◎  Age Range                    ▾

Value distribution
51-64
41-50
Other (5 Categories)

**Model type**

SageMaker Canvas automatically recommends the appropriate model type for your analysis.

🔆 3+ category prediction

Your model classifies Age Range into 3 or more categories.

**Configure model**

**Quick build**  ▾

**Preview model**

Dataset_20241209_105458  ▤ ▥ ⧩ 🔍 ≣ ◫ ≣ ⌁ View all          ⬙ Data visualizer  ≫

▥ Total columns: 7   ≣ Total rows: 14,819   ▥ Total cells: 103,733   ☑ Show dropped columns

---

**My models** > New model 2024-12-9 6:12:47 ... > **Version 1**     + Create new version  ⟳  ⋮

| Select | Build | **Analyze** | Predict | Deploy |

**Model overview**

Your model is being created. Quick build usually takes 2-15 minutes. You can now leave this view.

| | Expected build time | Build type | Detailed progress |
|---|---|---|---|
| | 2-15 minutes | Quick build | Generation column impact |

⊞ Dataset_20241209_105458   ▥ Total columns: 10   ≣ Total rows: 14,819   ▥ Total cells: 148,190   ◎ Patient Gender   🔆 2 category prediction

## Screenshot 1

My models > New model 2024-12-9 6:12:47 ... > Version 1

+ Create new version

Select | Build | Analyze | Predict | Deploy

**Model status** @ Quick build

Accuracy ⓘ    F1 ⓘ  Optimization metric

**53.524%   0.646**

Predict | Standard build | Deploy

The model predicts the correct Patient Gender 53.524% of the time. ⓘ

Overview | Scoring | Advanced metrics

Model leaderboard ≫

| Positive Class | F1 ⓘ  Optimization metric | Accuracy ⓘ | Precision ⓘ | Recall ⓘ | AUC-ROC ⓘ |
|---|---|---|---|---|---|
| Male  Female | 64.649% | 53.524% | 48.073% | 98.669% | 0.629 |

Predicted values

🖽 Dataset_20241209_105458   ▥ Total columns: 10   ▦ Total rows: 14,819   ▥ Total cells: 148,190   ◉ Patient Gender   ⬧ 2 category prediction   **Predict**

## Screenshot 2

My models > New model 2024-12-9 6:12:47 ... > Version 1

+ Create new version

Select | Build | Analyze | Predict | Deploy

| Column | Value |
|---|---|
| Report Year | 2015 |
| Diagnosis Sub Category | Other |
| Treatment Sub Category | Other |
| Determination | Upheld Decision ... ▾ |
| Type | Medical Necessity ▾ |

Patient Gender  ⧉ Copy

Prediction

# Male

■ New prediction
▦ Last prediction

Female    49.805% ⓘ

**☁ Download prediction ▾**

## Screenshot 3

My models > New model 2024-12-9 6:12:47 ... > Version 1

+ Create new version

Select | Build | Analyze | Predict

model drift.

**Deployments**

Filter by status:  In service   Failed   Creating

**Selected model version**

New model 2024-12-9 6:12:47 AM

V1  ✅  Ready   Created: 12-09-2024-6:21 AM

| Deployment name | Status | Deployment URL |
|---|---|---|

**Deployment type**

Real-time ⓘ

No result found

**Deployment name**

Deployment name
new-deployment-12-09-2024-6-21-AM

Instance type ⓘ          Learn about pricing ☑

ml.m5.12xlarge  Recommended

← → C 🔒 d-rouutmo9zsym.studio.us-east-1.sagemaker.aws/canvas/default/models/New%20model%202024-12-9%206:12:47%20AM ☆ ⟲ 🔵 ⋮

⊞ 🔖 Under Emails, the L... 🔵 Google Analytics for... 🔗 linkedin.com/in/mic... 🔵 Models & Languages R More Python Extension Pa... 🔷 Issue navigator - Jira » 🗀 All Bookmarks

My models › New model 2024-12-9 6:12:47 ... › **Version 1**     + Create new version   ⟲   ⋮

| Select | Build | Analyze | Predict | **Deploy** |

## Deployments

Filter by status: ( In service ) ( Failed ) ( Creating )      🔍 Search deployments    **+ Create Deployment**

| Deployment name | Status | Deployment URL | Created |
|---|---|---|---|
| canvas-new-deployment-12-09-2024... | ✅ In service | https://runtime.sagemaker.us-east-1.amazonaws.com/endp... 🗑 | 12/09/24 06:21 AM   ⋮ |

---

Operations: Deployment **canvas-new-deployment-12-09-2024-6-21-AM** /      **Update configuration**   ⋮   ✕

| Details | **Test deployment** |

Modify values to predict **Patient Gender** in real time.

🔍 Filter columns

| Column | Value | | Patient Gender | 🗐 Copy |
|---|---|---|---|---|
| Report Year | 2015 | | Prediction | |
| | | | **M - l -** | |

### Invocation result

| Status | Execution length (ms) ⓘ Learn more 🗗 | Request time |
|---|---|---|
| ✅ Successful | 0 | 2024-12-09 07:22:06 AM |

---

Operations: Deployment **canvas-new-deployment-12-09-2024-6-21-AM** /      **Update configuration**   ⋮   ✕

| **Details** | Test deployment |

| Deployment name | Status | Deployment type | Model |
|---|---|---|---|
| canvas-new-deployment-12-09-2024-6-21-AM | ✅ In service | Real-time | New model 2024-12-9 6:12:47 AM  (v1) |

| Created | Average predictions per day | Last prediction |
|---|---|---|
| 12/09/24 06:21 AM | 1 | — |

| Instance type | Instance count | Inference response content | Input format |
|---|---|---|---|
| ml.t2.medium | 3 | predicted_label, probability, probabilities, labels | text/csv |

**Deployment URL** Learn how to invoke a real-time endpoint 🗗

https://runtime.sagemaker.us-east-1.amazonaws.com/endpoints/canvas-new-deployment-12-09-2024-6-21-AM/invocations 🗐

Image 1: Captures the data import process into SageMaker Autopilot, with the processed dataset loaded from S3.

Image 2: Shows the experiment configuration, including the selected target column and AutoML settings.

Image 3: Highlights the experiment results, such as the model leaderboard and performance metrics.

Explanation: The images provide a step-by-step overview of SageMaker Autopilot, demonstrating its ability to automate model training, evaluation, and comparison, while ensuring no coding is required. Insights into ethical considerations, such as addressing data bias, may also be visualized.
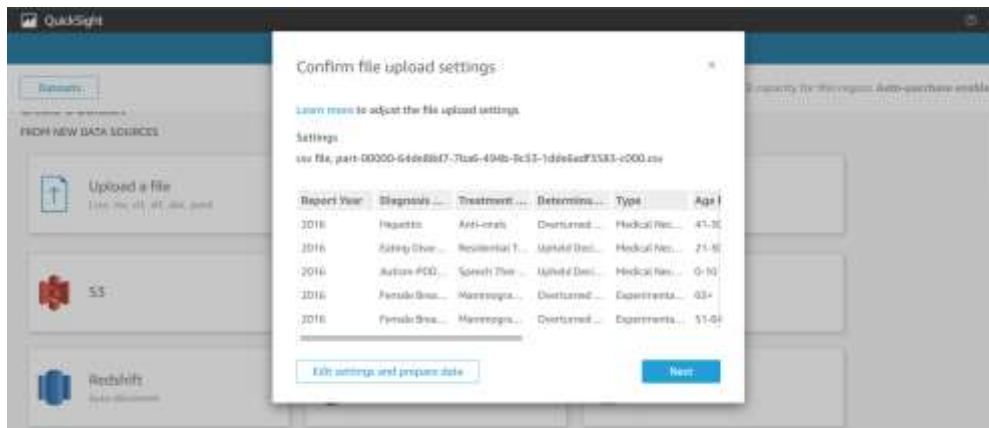
# 4. Visualization (1.5 Point)

**Tool Options:**

● **AWS QuickSight** for creating dynamic dashboards. (Note: AWS QuickSight offers a Free Tier for new users. The Free Tier includes the following: 1 GB of SPICE capacity. Here is a tutorial on how to use AWS QuickSight .)
**Task: Create Dashboards**

1. Connect QuickSight to the processed data in S3.

https://us-east-1.quicksight.aws.amazon.com/sn/accounts/914497062970/dashboards/d529018e-86d7-477c-8885-06e0a4a57962?directory_alias=hhgupta

2. Design a dashboard with at least 4 insightful visualizations. **Caution:** ∎ Be mindful of **QuickSight usage limits** as it may incur additional costs. ∎ Use Power BI for local visualization to avoid QuickSight charges.

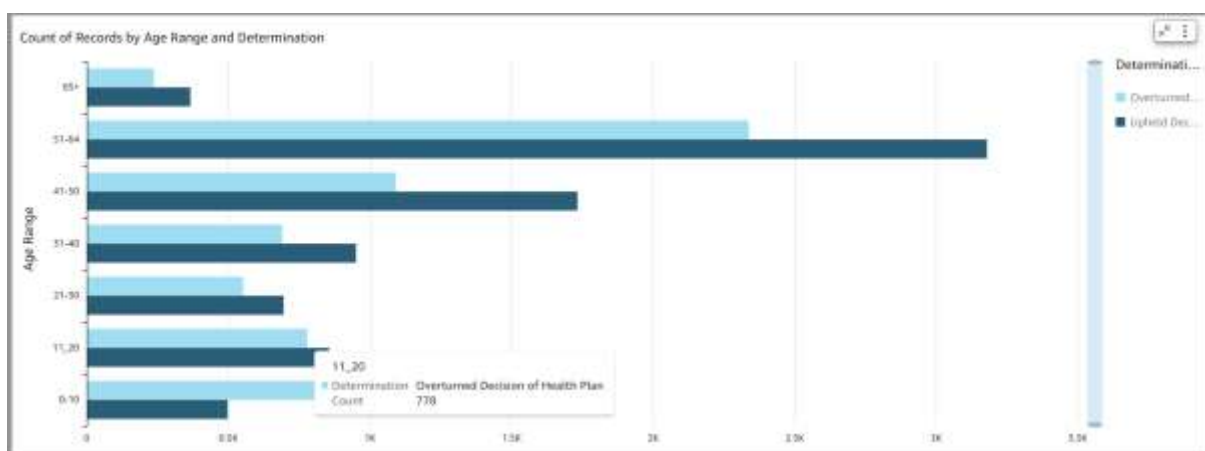# 1. Count of Records by Age Range and Determination

**Purpose**:
This bar chart visualizes the distribution of records across various age ranges, segmented by the determination outcome (e.g., "Overturned Decision of Health Plan" and "Upheld Decision of Health Plan").

**Insights Derived**:

- Higher concentrations of records are seen in the age range 51-64, indicating this group undergoes more health plan determinations.

- The "Overturned Decision" category is less common than "Upheld Decision," suggesting the majority of health plan decisions align with initial determinations.

**Filters or Parameters Applied**:

- Age ranges (0-10, 11-20, 21-30, etc.).

- Determination outcomes (e.g., overturned, upheld).

## 2. Count of Diagnosis and Treatment Summary by Report Year and Age Range
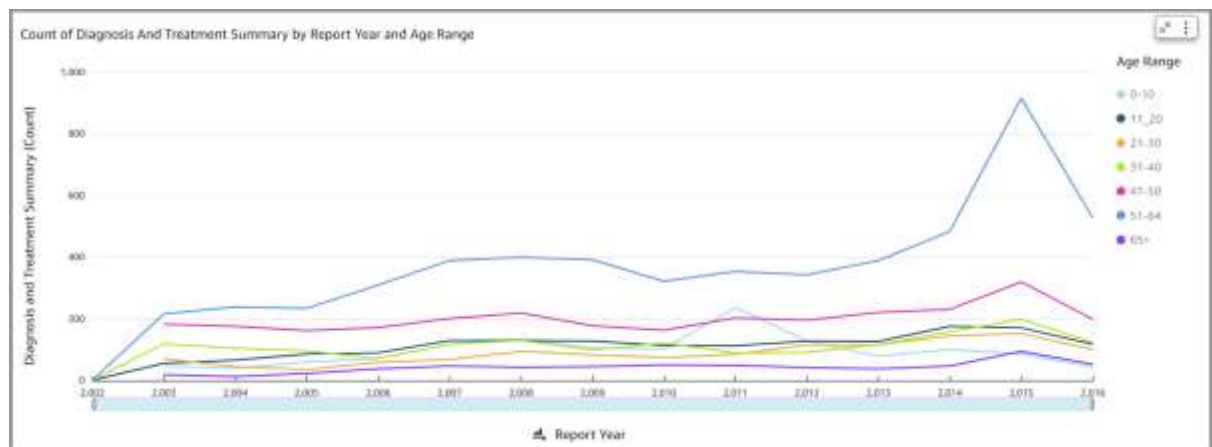
**Purpose**:
This line graph illustrates the number of diagnoses and treatments reported annually, categorized by age range.

**Insights Derived**:

- A steady increase in the number of diagnoses and treatments is observed over time, particularly in the 51-64 and 65+ age groups.

- Fluctuations in younger age ranges indicate inconsistent or lower levels of reporting.

**Filters or Parameters Applied**:

- Report years (2002-2016).

- Age groups (same as above).



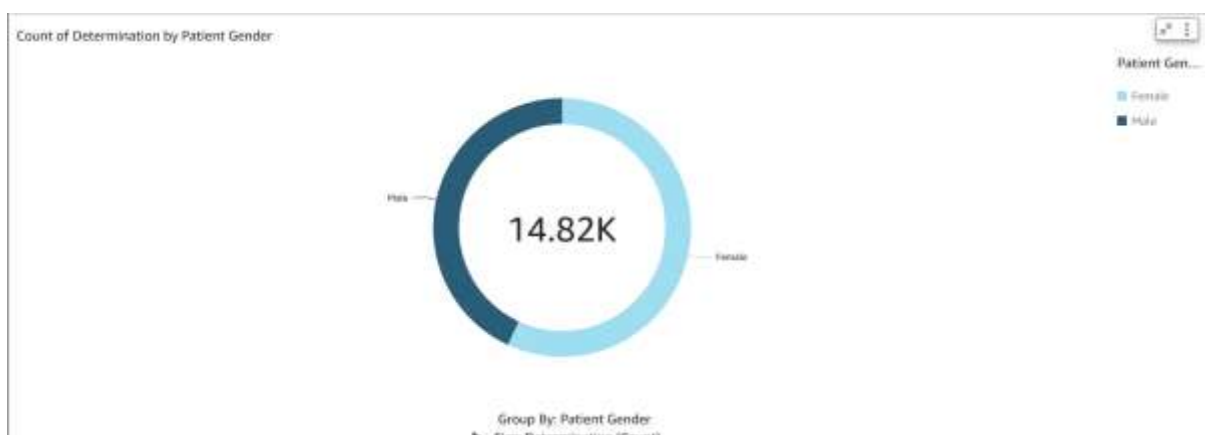## 3. Count of Determination by Patient Gender

**Purpose**:
This pie chart highlights the proportion of health plan determinations for males versus females.

**Insights Derived**:

- A balanced distribution between male and female patients, slightly skewed toward females receiving more determinations.

- Potential gender-related trends or disparities in health plan review processes.

**Filters or Parameters Applied**:

- Patient gender (Male, Female).

- Determination counts.

---



Count of Determination by Patient Gender

Patient Gen...

■ Female
■ Male

Male

14.82K

Female

Group By: Patient Gender

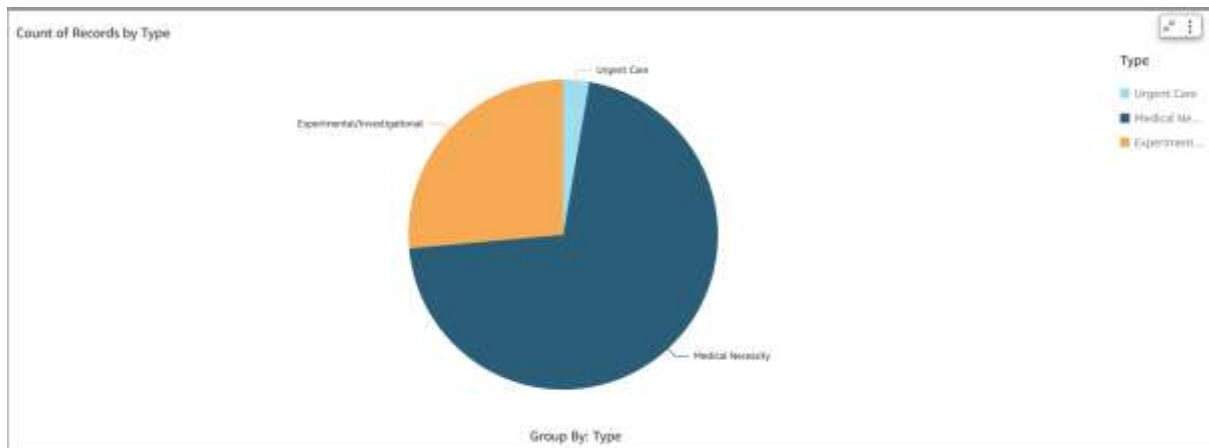## 4. Count of Records by Type

**Purpose**:
A horizontal bar chart showing the distribution of records based on the type of service (e.g., "Urgent Care," "Medical Necessity").

**Insights Derived**:

- "Medical Necessity" dominates the record count, reflecting its importance in health plan assessments.

- Urgent care records are significantly fewer, indicating lower frequencies of urgent cases under review.

**Filters or Parameters Applied**:

- Record types (e.g., "Urgent Care," "Medical Necessity")

Count of Records by Type

Group By: Type

## 5. Count of Determination by Patient Gender and Determination Outcome

**Purpose**:
A clustered bar chart comparing determination outcomes ("Overturned" vs. "Upheld") across genders.

**Insights Derived**:

- Females show a higher count of "Upheld Decision" cases compared to males.

- Disparity in "Overturned Decisions" between genders might point to gender-based discrepancies in health plan appeals.

**Filters or Parameters Applied**:

- Determination type and patient gender.

Count of Determination by Patient Gender and Determination

| Patient Gender | Determination | Determination |
|---|---|---|
| Female | Overturned Decision of Health Plan | 3,763 |
| Male | Overturned Decision of Health Plan | 2,779 |
| Female | Upheld Decision of Health Plan | 4,673 |
| Male | Upheld Decision of Health Plan | 3,604 |

## 6. Count of Records by Determination and Type

**Purpose**

This bar chart categorizes health records by the type of service (e.g., "Medical Necessity," "Urgent Care," "Experimental Treatments") and the determination outcome ("Overturned" or "Upheld").
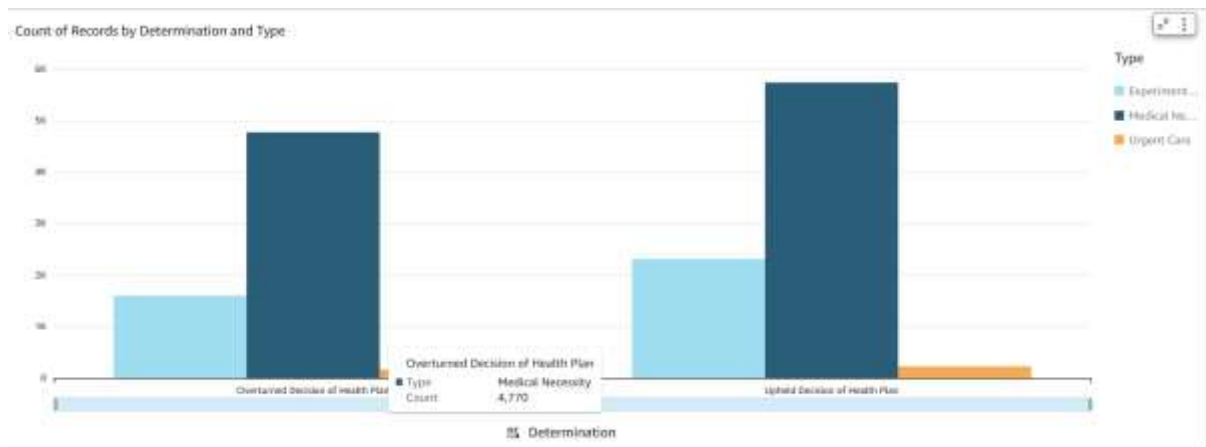
**Insights Derived**

- **Medical Necessity**: This category has the highest number of determinations, with most decisions being upheld. This suggests a higher degree of alignment with initial assessments for medical necessity cases.

- **Experimental Treatments**: A noticeable proportion of these cases are overturned, indicating potential challenges in meeting approval criteria for novel treatments.

- **Urgent Care**: Although a smaller category, urgent care records show a balanced distribution between overturned and upheld decisions.

**Filters or Parameters Applied**

- **Service Types**: Filtered by "Medical Necessity," "Urgent Care," and "Experimental Treatments."

- **Determination Outcomes**: Segmented into "Overturned" and "Upheld."



Count of Records by Determination and Type

# 5. Bonus Task: Automation of the Pipeline (5 Points)

**Objective: Automate the entire data pipeline from ingestion to visualization**:

1. **Develop an Automation Script (1 Point):** ○ Automate data retrieval, processing, and storage steps.



i-0e80de008aa41d12d (CalihealthEC2)

PublicIPs: 54.89.64.27   PrivateIPs: 172.31.81.146

2. **Set Up Scheduling Tools (1.5 Points):** ○ Use **cron jobs** or **AWS Lambda** to trigger the pipeline on a schedule or upon new data uploads to S3 bucket.

```
ubuntu@ip-172-31-81-146:~$ chmod +x /home/ubuntu/daily_task.sh
ubuntu@ip-172-31-81-146:~$ crontab -e
no crontab for ubuntu - using an empty one
crontab: installing new crontab
ubuntu@ip-172-31-81-146:~$ crontab -l
# Edit this file to introduce tasks to be run by cron.
#
# Each task to run has to be defined through a single line
# indicating with different fields when the task will be run
# and what command to run for the task
#
# To define the time you can provide concrete values for
# minute (m), hour (h), day of month (dom), month (mon),
# and day of week (dow) or use '*' in these fields (for 'any').
#
# Notice that tasks will be started based on the cron's system
# daemon's notion of time and timezones.
#
# Output of the crontab jobs (including errors) is sent through
```

## 3. Implement Logging and Notifications (1.5 Points): ○ Use AWS SNS or email to notify users of pipeline status (success or failure).



[External] Amazon S3 Notification

AN   AWS Notifications<no-reply@sns.amazonaws.com>    ☺ ↩ ↩ ↪ ···
     To: Gupta, Harshit Hemant                        Mon 12/9/2024 2:35 AM

[You don't often get email from no-reply@sns.amazonaws.com. Learn why this is important at https://aka.ms/LearnAboutSenderIdentification ]

This message was sent from a non-IU address. Please exercise caution when clicking links or opening attachments from external sources.

{"Service":"Amazon S3","Event":"s3:TestEvent","Time":"2024-12-09T07:35:06.080Z","Bucket":"calihealthcares3","RequestId":"7G07JWGSFFBF8315","HostId":"KpoR5EpQIBCdEOCPnl3BCD5G8Q6au4c2MDn/Hbgo9erjLjjLEEkvOJ2dYKKYxiqTVGM1jnYUPmU="}

--

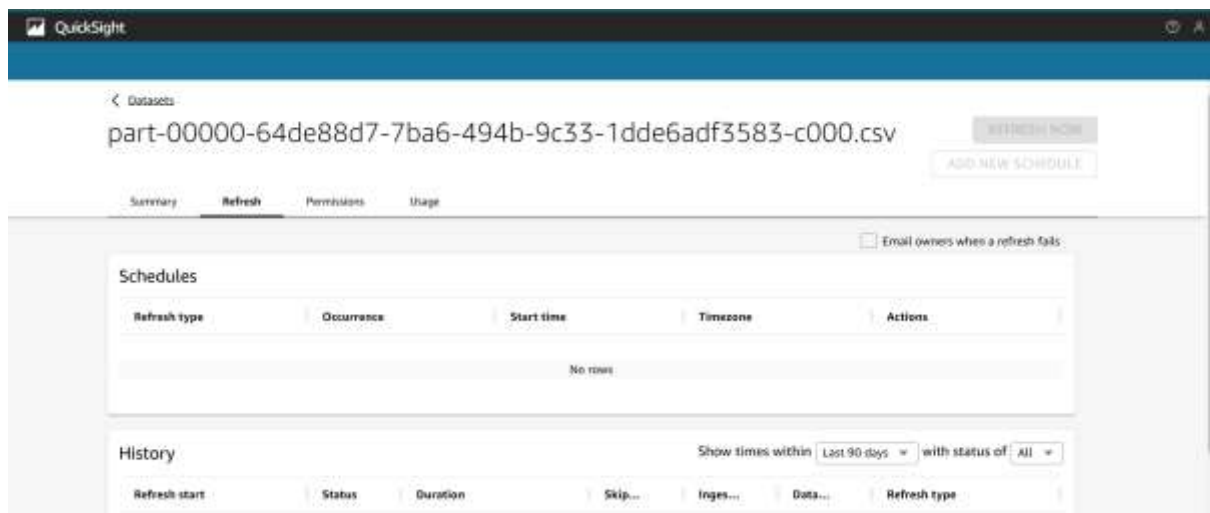If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:
https://nam12.safelinks.protection.outlook.com/?url=https%3A%2F%2Fsns.us-east-

This message was sent from a non-IU address. Please exercise caution when clicking links or opening attachments from external sources.

{"Service":"Amazon S3","Event":"s3:TestEvent","Time":"2024-12-09T07:35:06.080Z","Bucket":"calihealthcares3","RequestId":"7G07JWGSFFBF8315","HostId":"KpoR5EpQIBCdEOCPnl3BCD5G8Q6au4c2MDn/Hbgo9erjLjjLEEkvOJ2dYKKYxiqTVGM1jnYUPmU="}

--

If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:

https://nam12.safelinks.protection.outlook.com/?url=https%3A%2F%2Fsns.us-east-1.amazonaws.com%2Funsubscribe.html%3FSubscriptionArn%3Darn%3Aaws%3Asns%3Aus-east-1%3A914497062970%3AS3UploadNotification%3Aa173f175-ac50-45c2-9117-b334fe04d09e%26Endpoint%3Dhhgupta%40iu.edu&data=05%7C02%7Chhgupta%40iu.edu%7C321db7d33a9e434b230608dd182400b7%7C1113be34aed14d00ab4bcdd02510be91%7C0%7C0%7C638693265107304161%7CUnknown%7CTWFpbGZsb3d8eyJFbXB0eU1hcGkiOnRydWUsIlYiOiIwLjAuMDAwMCIsIlAiOiJXaW4zMiIsIkFOIjoiTWFpbCIsIldUIjoyfQ%3D%3D%7C0%7C%7C%7C&sdata=vWNXsAmlKQsN6%2B3GEwtjLu2nQtSIZCorw%2BDowyr%2BsK0%3D&reserved=0

Please do not reply directly to this email. If you have any questions or comments



**4. Automate Dashboard Updates (1 Point):** ○ Configure QuickSight or Power BI to refresh data periodically for real-time insights.

**Bonus:** Submit an architecture diagram showing the entire pipeline.