# 1   Business Description

## 1.1   Business Overview

FitFlow Fitness is a multi-location gym and wellness chain that operates 30+ clubs across major cities and offers both in-person and digital fitness services. Members can access gyms, attend group classes, book personal training sessions, and join live or on-demand workouts through a mobile app. The business targets busy professionals, students, and fitness enthusiasts who want flexible access to facilities and coaching across locations.

## 1.2   Customers

FitFlow's primary customers are individual members who purchase monthly or annual memberships. They are segmented into tiers: 1. Basic: Single-location access; 2. Plus: Multi-location access; 3. Premium: Multi-location + classes + digital content.

Secondary customers include corporate clients who purchase bulk memberships for employees, and personal clients who pay for 1:1 or small-group coaching. Trainers and instructors are important "internal users" of the system, since they rely on scheduling, attendance, and performance.

## 1.3   Revenue Streams

The main revenue stream is recurring membership fees from subscription plans. Additional revenue comes from:

- Paid add-ons (class packs, personal training sessions, nutrition coaching)

- Drop-in day passes for non-members

- Corporate membership contracts

- Retail sales (merchandise, supplements, drinks at the front desk)

- Digital add-ons such as premium online programs and live-streamed classes

Because of the subscription model and add-ons, the business generates a high volume of transactional data related to membership billing, class bookings, and check-ins.

## 1.4   Data Needs

Operationally, FitFlow needs to track members, memberships, locations, staff, class schedules, bookings, and payments.

For the **OLTP system**, key data includes member profiles, active subscriptions, check-ins at each location, class sessions (who taught, where, when), member bookings, and payment records (amount, method, status). This data supports daily operations such as verifying access at the front desk, enforcing capacity limits in classes, allocating trainers, and handling billing and failed payments.

From an **analytical (OLAP) perspective**, FitFlow wants to analyze trends across time, locations, and member segments. Example questions include: Which locations have the highest utilization? What plans have the best retention? Which classes drive the most engagement and add-on revenue? This requires aggregating historical transactional data into fact tables (e.g., `fact_membership_billing`, `fact_class_attendance`) with dimensions such as time, location, member, plan, and class type. Overall, the business has high data throughput and clear separation between real-time transactional needs and historical analytics.
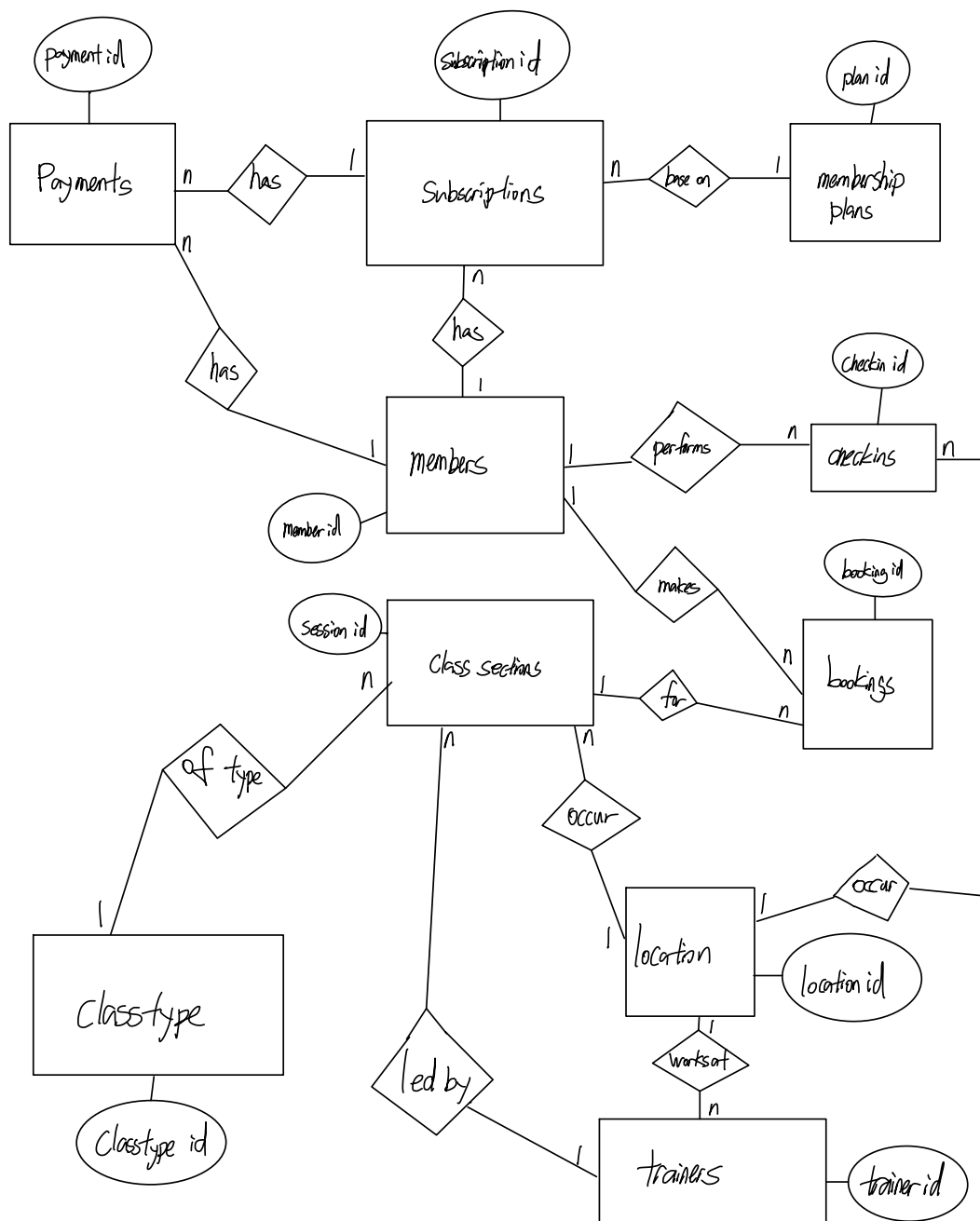
Figure 1: ER Diagram for OLTP Schema (Keys and Partial Attributes)

# 2 Transactional vs. Archival Data Needs

For FitFlow Fitness, the data architecture is bifurcated into two distinct categories based on operational requirements and analytical goals.

## 2.1 Transactional Data (OLTP)

The transactional data needs are driven by the daily operations of the 30+ gym locations. This data is characterized by high-frequency writes, real-time availability, and strict consistency (ACID properties). It's mainly used for:

- **Access Control & Check-ins:** The system must verify membership status in real-time when a member scans their badge.

- **Class Bookings:** The system tracks individual session bookings to enforce capacity limits. This requires immediate consistency to prevent overbooking.

- **Billing:** Payment processing generates critical transactional records that link specific members to specific subscription charges and add-on purchases.

## 2.2 Archival Data (OLAP)

The archival data needs focus on historical analysis to support strategic decision-making of the management. This data is read-heavy and often aggregated to identify trends across time, locations, and member segments.

- **Trend Analysis:** Management requires insights into utilization rates (e.g., "Which locations have the highest utilization?") and financial performance over long periods.

- **Aggregated Metrics:** Unlike the granular transactional view ("John Doe booked Spin Class at 5 PM", for example), the archival view aggregates data into metrics such as "Total Revenue per Location per Month" or "Retention Rate by Plan Tier".

**Notes:** Separating these systems is important in the context of the project as running heavy analytical queries (e.g., aggregating years of data) on the OLTP system would require extensive read locks. This, however, might inhibit current write operations, possibly causing latency for members trying to check in or book classes in real-time.

# 3 Relational Schema Design

## 3.1 OLTP Schema

The OLTP system is designed in 3rd Normal Form (3NF) to minimize redundancy and maintain data integrity. It strictly enforces referential integrity through foreign keys. The schema consists of 10 tables:

- **members**: Stores profile information (PK: member_id).

- **locations**: Stores gym branch details including capacity (PK: location_id).

- **membership_plans**: Defines membership tiers (Basic, Plus, Premium) and pricing (PK: plan_id).

- **class_types**: Defines the catalog of available classes (e.g., Yoga, HIIT) and their categories (PK: class_type_id).

- **trainers**: Stores trainer details and links them to their assigned "home" location (PK: trainer_id, FK: location_id).

- **subscriptions**: Links members to specific plans and tracks status (PK: subscription_id, FKs: member_id, plan_id).

- **checkins**: Records individual member visits to a location (PK: checkin_id, FKs: member_id, location_id).

- **class_sessions**: Represents specific scheduled class instances with an assigned trainer and location (PK: session_id, FKs: class_type_id, location_id, trainer_id).

- **bookings**: Records member reservations for specific class sessions (PK: booking_id, FKs: member_id, session_id).

- **payments**: Records transactional billing events linked to subscriptions (PK: payment_id, FKs: subscription_id, member_id).

## 3.2 OLAP Schema (Star Schema)

For the analytical system, we implemented a Star Schema to facilitate high-performance read queries. This involves denormalized dimension tables and centralized fact tables.

### 3.2.1 Dimension Tables

- **Dim_Location**: Denormalized location data (Attributes: City, State, Capacity).

- **Dim_Plan**: Plan details (Attributes: Plan Name, Tier, Price).

- **Dim_Class_Type**: Class categorization (Attributes: Category, Name).

- **Dim_Date**: A comprehensive date dimension for temporal analysis (Attributes: Year, Month, Day, Is_Holiday).

### 3.2.2 Fact Tables

- **Fact_Monthly_Revenue**: Aggregates financial performance.

  - *Granularity*: Per Location, Per Plan, Per Month.
  - *Measures*: Total Revenue, Transaction Count, Active Members.

- **Fact_Class_Attendance**: Tracks operational efficiency.

  - *Granularity*: Per Location, Per Class Type, Per Day.
  - *Measures*: Total Bookings, Attendance Count, Utilization Rate (derived from Attendance / Capacity).

# 4 ETL Process and Implementation

## 4.1 Implementation Environment

We implemented both the OLTP and OLAP systems using Azure SQL Database. The ETL process was developed using Python, leveraging pandas and SQLAlchemy. This setup allows for flexible data manipulation while adhering to the logical separation of operational and analytical processing.

## 4.2 ETL Workflow Description

1. **Extract**

   To minimize the load on the live operational system and prevent potential table locks that could disrupt front-desk check-ins or member bookings, we adhered to two strict extraction rules:

   - **No Joins on OLTP:** We extracted raw data from the `Payments`, `Subscriptions`, and `Bookings` tables individually. We avoided performing any `JOIN` operations or complex aggregations on the OLTP side.

   - **Incremental Extraction:** Although we populated the initial database with synthetic data, our extraction logic is designed to be incremental. By checking the timestamp of the last ETL run, we only extract records that have been created or modified since that time, rather than reloading the entire dataset.

2. **Transform**

   All raw data was loaded into `pandas` DataFrames, which served as our temporary staging area. Processing data here ensures that transformations do not impact the OLTP system.

   - **Schema Mapping & Joins:** As the OLTP schema is normalized, we performed joins in the staging area to enrich the data before aggregation. For example, to generate the `Fact_Class_Attendance` table, we joined `bookings` with `class_sessions` to associate each booking with its specific class type and location. Although `trainers` data exists in OLTP for operational scheduling, it was not moved to OLAP as trainer performance analysis was outside the current scope of our Dimension design.

   - **Data Cleaning & Integrity:** We filtered out invalid records (e.g., failed payments) and standardized date formats. Transaction timestamps were converted into standard date keys (YYYYMMDD) to align with the `Dim_Date` table.

3. **Load**

   The final step involved loading the transformed data into the Azure OLAP system.

   - **"Dimensions First":** We strictly loaded the Dimension tables (`Dim Location`, `Dim Plan`, `Dim Class Type`) first. This ensures that when Fact tables are loaded, the foreign keys (Surrogate Keys) correctly reference existing dimension records.

   - **Fact Table Loading:** Finally, the aggregated transactional data was bulk-loaded into `Fact Monthly Revenue` and `Fact Class Attendance` using the `to_sql` method in append mode.

# 5    FitFlow Management Dashboard

We developed an interactive dashboard using Python (Dash/Plotly) connected to the Azure OLAP database. The dashboard provides management with real-time insights into revenue and operational efficiency, featuring interactive filtering capabilities.

## 5.1    Monthly Revenue Trends (Overall)

This visualization tracks the Total Revenue Over Time. When no location is selected, the chart aggregates revenue across all locations to provide an overall business view of the company's financial growth.
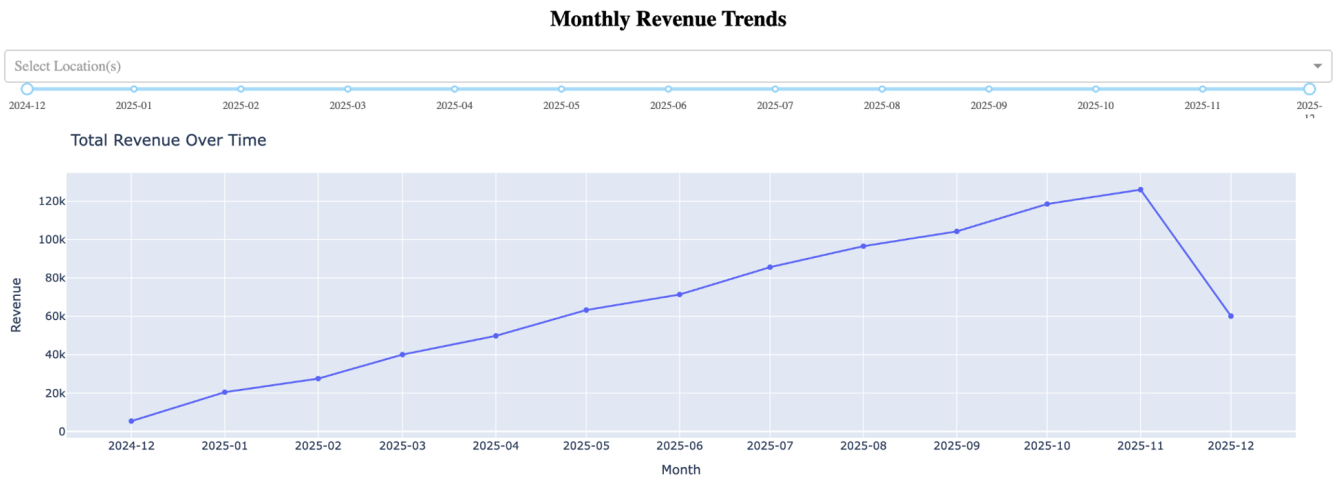


Figure 2: Monthly Revenue Trends (Aggregated)

## 5.2    Monthly Revenue Trends (Filtered)

This version of the chart shows its interactivity. It is filtered by specific locations and a selected range of months, allowing for granular regional performance comparison.
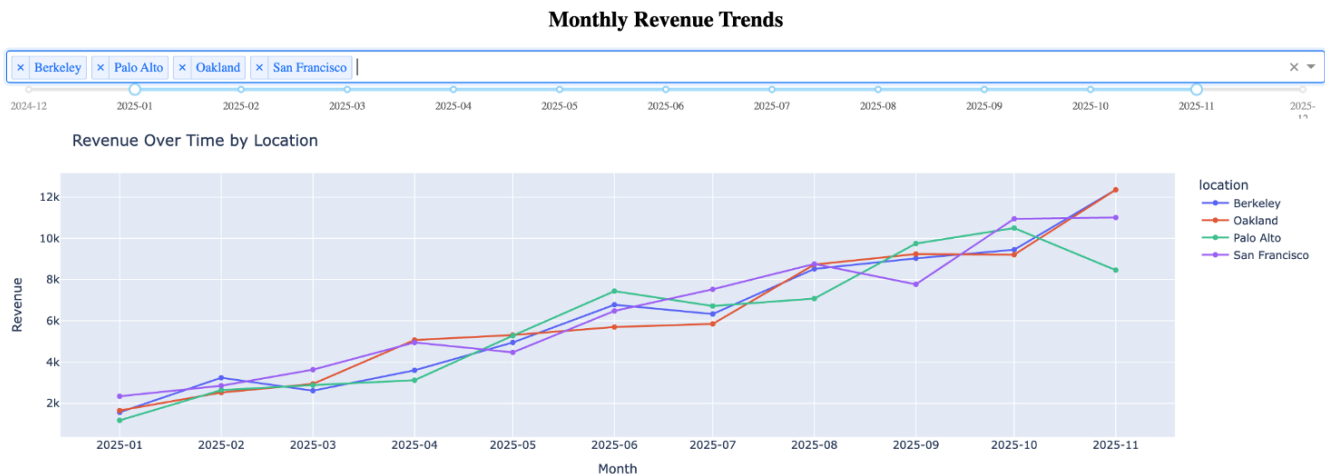


Figure 3: Monthly Revenue Trends (Filtered by Location)

## 5.3   Total Attendance by Class Type

This bar chart displays Class Popularity by Attendance volume. It helps identify which class types (e.g., HIIT Blast, Pilates Core) are driving the most foot traffic to the gyms.
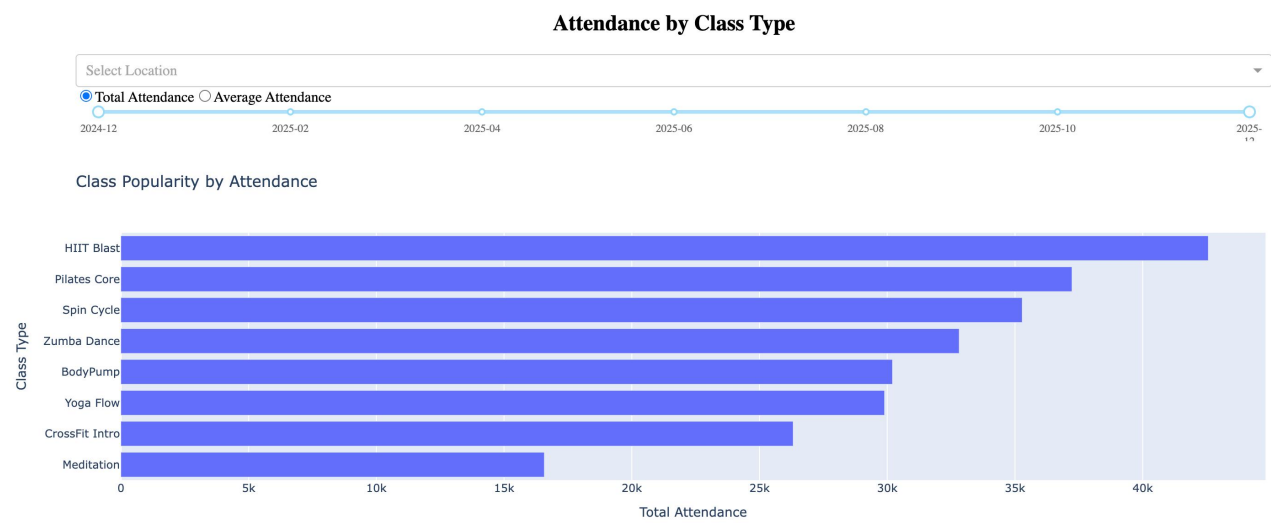


Figure 4: Total Attendance by Class Type

## 5.4   Average Attendance by Class Type

By toggling the metric to "Average Attendance", management can assess efficiency. This view normalizes the data to show attendance per individual session, highlighting classes that are consistently full even if they are offered less frequently.
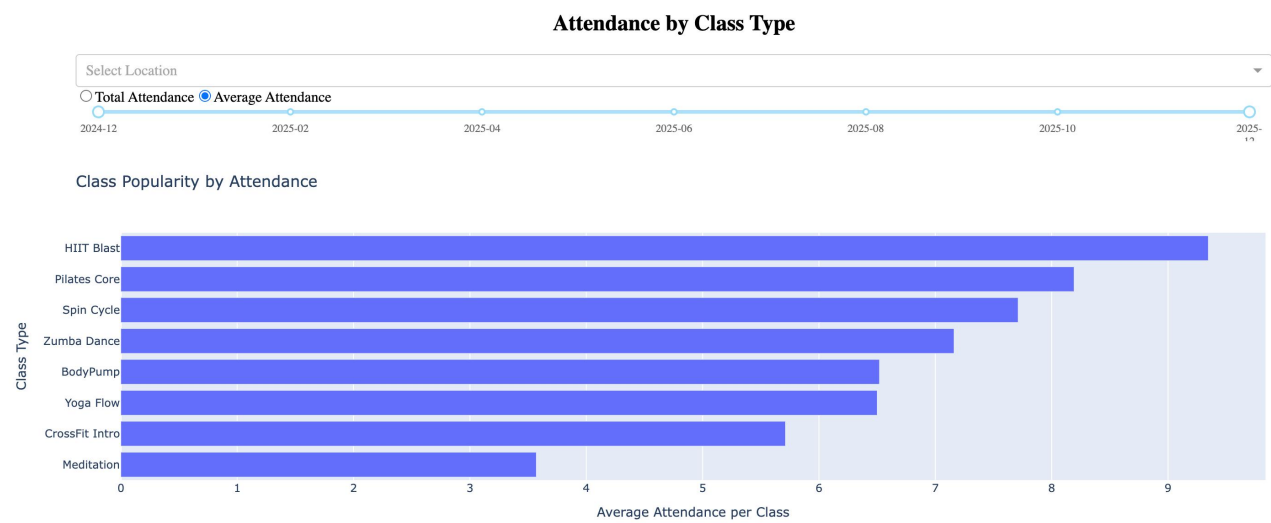


Figure 5: Average Attendance by Class Type

## 5.5   Attendance by Class Type (Filtered)

This plot shows attendance filtered by a specific location (Fremont) and a specific time range (Q3, 2025). This level of detail aids in local resource planning and schedule optimization for specific
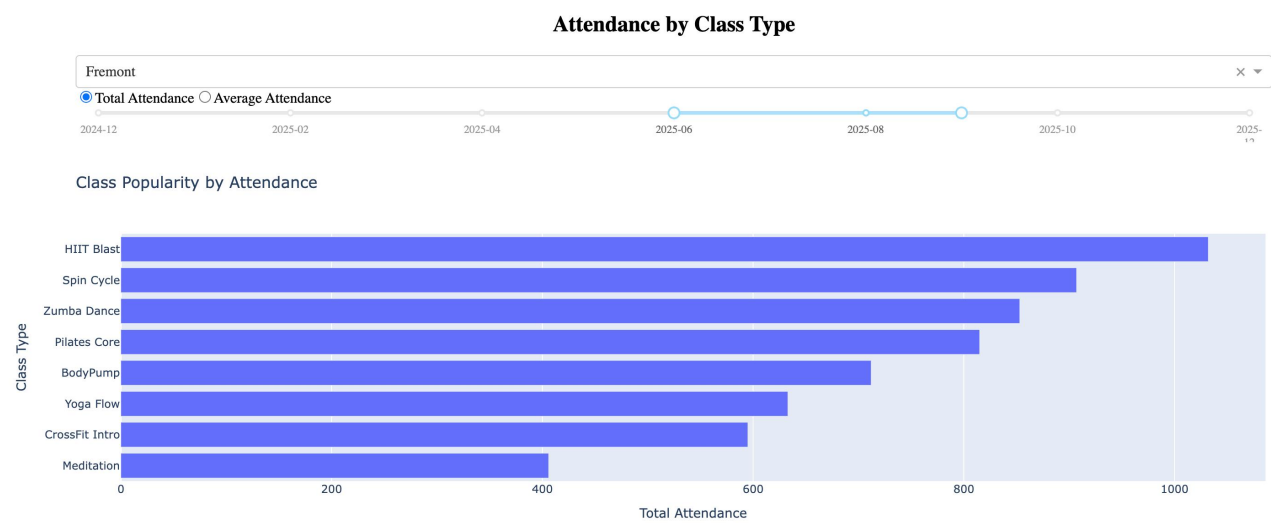
branches.



Figure 6: Attendance Filtered by Location (Fremont) and Month

## 5.6 Seasonality Analysis (Overall Heatmap)

This heatmap visualizes the seasonality in class attendance across all locations. Darker colors indicate higher attendance levels. Distinct patterns, such as "New Year" spikes in January, highlight participation behavior across the network.
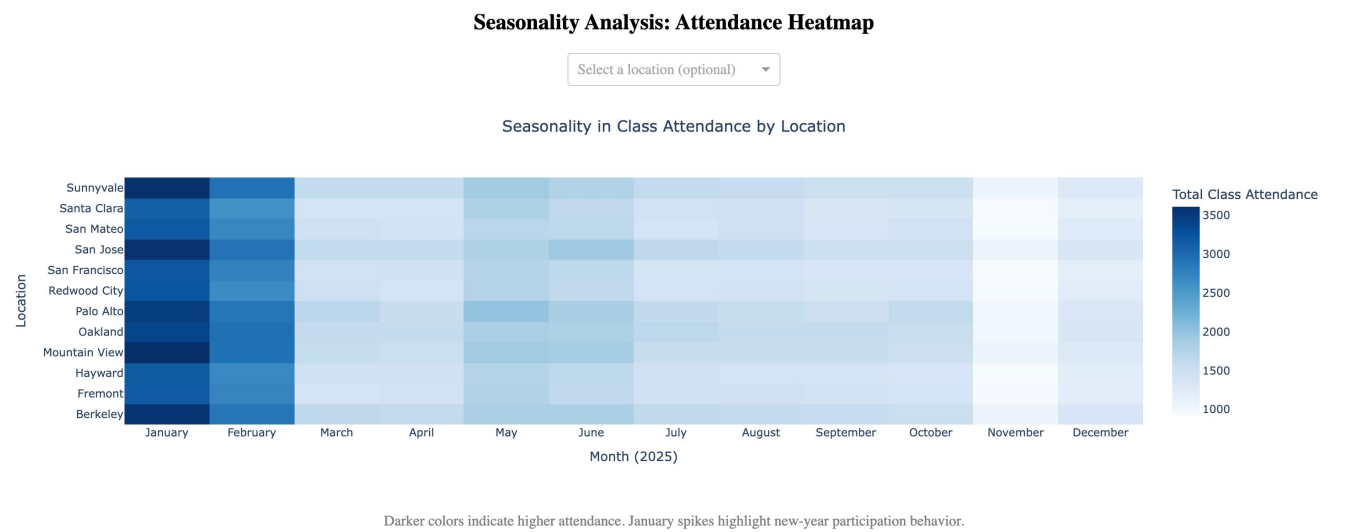


Figure 7: Seasonality Analysis: Attendance Heatmap

## 5.7 Seasonality Analysis (Filtered)

This heatmap is filtered to show data for a single location (Berkeley). It allows the specific branch manager to analyze their unique seasonal trends independent of the wider network.
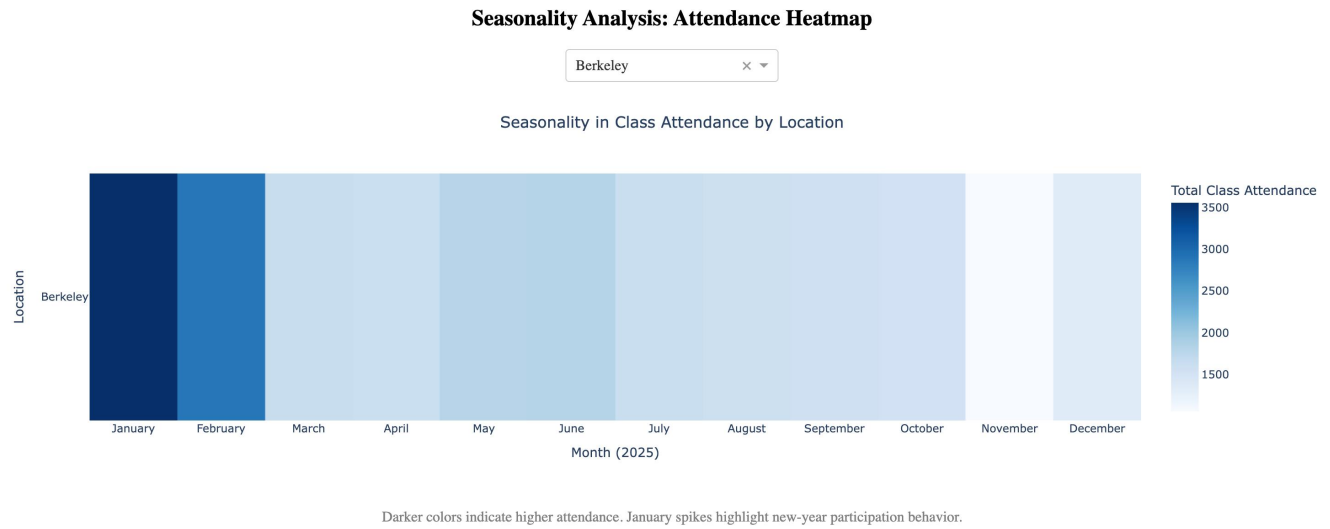
**Seasonality Analysis: Attendance Heatmap**



Figure 8: Seasonality Analysis Filtered (Berkeley)

# 6    Informative Analysis

## 6.1    Monthly Revenue Analysis by Subscription Plan

**Objective:** The objective of this analytical process is to understand how monthly revenue is distributed across different subscription plans over time. By analyzing revenue at the plan level, this process helps identify which plans contribute most to overall revenue growth and how customer preferences evolve across months.

**Data Preparation:** This analysis uses data stored in an OLAP-style schema. The main data source is the `Fact_Monthly_Revenue` table, joined with `Dim_Plan` and `Dim_Date`. Using SQL, data is filtered for the period Dec 2024 - Dec 2025. Missing months are explicitly added to ensure a consistent time axis.

**Analytical Workflow:**

1. Query monthly revenue data by joining fact and dimension tables.

2. Aggregate total revenue by month and subscription plan.

3. Standardize the monthly timeline to ensure continuity.

4. Prepare the data for comparison across plans.



Figure 9: Monthly Revenue Contribution by Plan

## 6.2    Daily Utilization Rate Analysis and Forecasting

**Objective:** The objective is to analyze daily utilization rates over time and forecast short-term future utilization trends. Utilization rate reflects how efficiently resources (classes/facilities) are being used.

**Data Preparation:** Based on `Fact_Class_Attendance`. Missing dates are explicitly added to create a continuous daily time series, and missing utilization values are filled using time-based linear interpolation.

**Analytical Workflow:**

1. Retrieve daily utilization data using SQL joins.

2. Aggregate utilization rates at the daily level.

3. Enforce continuity in the time dimension and apply interpolation.

4. Train a Seasonal ARIMA (SARIMA) model on the cleaned daily time series.

5. Generate both in-sample fitted values and short-term forecasts.



Figure 10: Daily Utilization Rate Forecast (SARIMA)

## 6.3   Daily Revenue Trend Analysis and Forecasting

**Objective:** To analyze daily revenue patterns and forecast short-term future revenue trends. Daily revenue is a core business metric that reflects customer activity and overall financial performance.

**Data Preparation:** Based on `Fact_Monthly_Revenue`. Revenue values are aggregated by day. Missing dates are added, and missing revenue values are filled using time-based linear interpolation.

**Analytical Workflow:**

1. Retrieve daily revenue data using SQL joins.

2. Aggregate revenue at the daily level.

3. Enforce continuity and apply time-based interpolation.

4. Train a Seasonal ARIMA (SARIMA) model to capture trend and weekly seasonality.

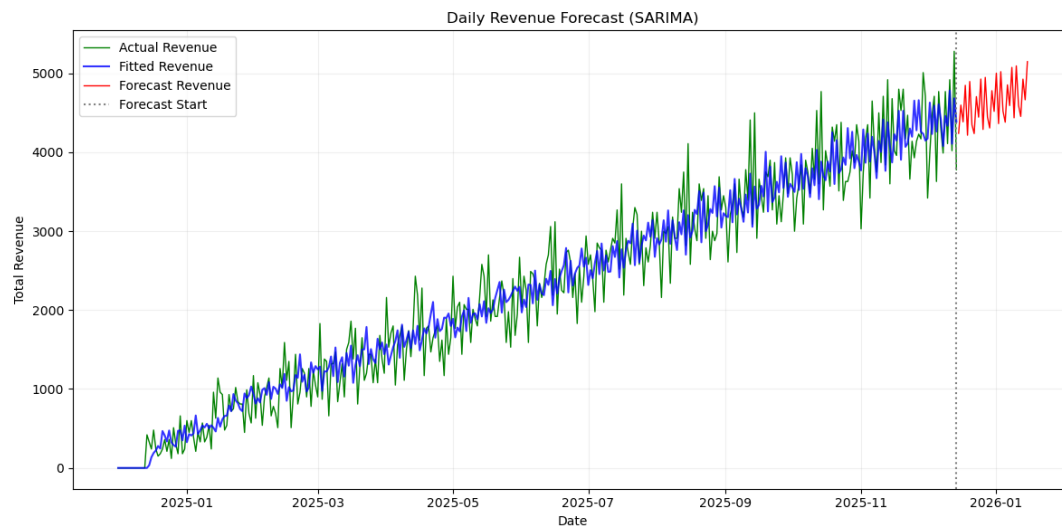5. Generate in-sample fitted values and short-term revenue forecasts.

Figure 11: Daily Revenue Forecast (SARIMA)

# A  Appendix: SQL and Code Implementation

## A.1  SQL CREATE Statements (OLTP & OLAP)

```sql
-- 1. Independent Tables
CREATE TABLE locations (
    location_id VARCHAR(50) PRIMARY KEY,
    city VARCHAR(50), state VARCHAR(50), capacity INT
);
CREATE TABLE class_types (
    class_type_id VARCHAR(50) PRIMARY KEY,
    name VARCHAR(50), category VARCHAR(20)
);
CREATE TABLE membership_plans (
    plan_id VARCHAR(50) PRIMARY KEY,
    plan_name VARCHAR(50), price DECIMAL(10, 2), tier VARCHAR(20)
);
CREATE TABLE members (
    member_id VARCHAR(50) PRIMARY KEY,
    first_name VARCHAR(50), last_name VARCHAR(50), join_date DATE
);

-- 2. Dependent Tables (Level 1)
CREATE TABLE trainers (
    trainer_id VARCHAR(50) PRIMARY KEY,
    name VARCHAR(100),
    location_id VARCHAR(50) FOREIGN KEY REFERENCES locations(location_id)
);
CREATE TABLE subscriptions (
    subscription_id VARCHAR(50) PRIMARY KEY,
    member_id VARCHAR(50) FOREIGN KEY REFERENCES members(member_id),
    plan_id VARCHAR(50) FOREIGN KEY REFERENCES membership_plans(plan_id),
    start_date DATE, status VARCHAR(20)
);
CREATE TABLE checkins (
    checkin_id VARCHAR(50) PRIMARY KEY,
    member_id VARCHAR(50) FOREIGN KEY REFERENCES members(member_id),
    location_id VARCHAR(50) FOREIGN KEY REFERENCES locations(location_id),
    checkin_time DATETIME
);

-- 3. Dependent Tables (Level 2 & 3)
CREATE TABLE class_sessions (
    session_id VARCHAR(50) PRIMARY KEY,
    class_type_id VARCHAR(50) FOREIGN KEY REFERENCES class_types(class_type_id),
    location_id VARCHAR(50) FOREIGN KEY REFERENCES locations(location_id),
    trainer_id VARCHAR(50) FOREIGN KEY REFERENCES trainers(trainer_id),
    start_time DATETIME, duration_minutes INT
);
CREATE TABLE bookings (
    booking_id VARCHAR(50) PRIMARY KEY,
    member_id VARCHAR(50) FOREIGN KEY REFERENCES members(member_id),
    session_id VARCHAR(50) FOREIGN KEY REFERENCES class_sessions(session_id),
    booking_date DATETIME, status VARCHAR(20)
);
CREATE TABLE payments (
    payment_id VARCHAR(50) PRIMARY KEY,
```

```
54      subscription_id VARCHAR(50) FOREIGN KEY REFERENCES subscriptions(subscription_id),
55      member_id VARCHAR(50) FOREIGN KEY REFERENCES members(member_id),
56      amount DECIMAL(10, 2), payment_date DATETIME, status VARCHAR(20)
57  );
```

Listing 1: OLTP Schema Creation

```
1  -- Dimension Tables
2  CREATE TABLE Dim_Location (
3      location_key INT IDENTITY(1,1) PRIMARY KEY,
4      location_id VARCHAR(50), city VARCHAR(50),
5      state VARCHAR(50), capacity INT
6  );
7  CREATE TABLE Dim_Plan (
8      plan_key INT IDENTITY(1,1) PRIMARY KEY,
9      plan_id VARCHAR(50), plan_name VARCHAR(50),
10     tier VARCHAR(20), price DECIMAL(10, 2)
11 );
12 CREATE TABLE Dim_Class_Type (
13     class_type_key INT IDENTITY(1,1) PRIMARY KEY,
14     class_type_id VARCHAR(50), name VARCHAR(50), category VARCHAR(20)
15 );
16 CREATE TABLE Dim_Date (
17     date_key INT PRIMARY KEY, -- YYYYMMDD
18     full_date DATE, year INT, month INT,
19     day_name VARCHAR(20), is_holiday BIT
20 );
21
22 -- Fact Tables
23 CREATE TABLE Fact_Monthly_Revenue (
24     revenue_key INT IDENTITY(1,1) PRIMARY KEY,
25     location_key INT FOREIGN KEY REFERENCES Dim_Location(location_key),
26     plan_key INT FOREIGN KEY REFERENCES Dim_Plan(plan_key),
27     date_key INT FOREIGN KEY REFERENCES Dim_Date(date_key),
28     total_revenue DECIMAL(15, 2),
29     transaction_count INT,
30     active_members INT
31 );
32 CREATE TABLE Fact_Class_Attendance (
33     attendance_key INT IDENTITY(1,1) PRIMARY KEY,
34     location_key INT FOREIGN KEY REFERENCES Dim_Location(location_key),
35     class_type_key INT FOREIGN KEY REFERENCES Dim_Class_Type(class_type_key),
36     date_key INT FOREIGN KEY REFERENCES Dim_Date(date_key),
37     total_bookings INT,
38     attendance_count INT,
39     utilization_rate DECIMAL(5, 4)
40 );
```

Listing 2: OLAP Schema Creation

## A.2   Python Code for Synthetic Data Generation

```
1  import pandas as pd
2  import random
3  from faker import Faker
4  from sqlalchemy import create_engine
5  import urllib
```

```
 6  import uuid
 7  from datetime import datetime, timedelta
 8
 9  # Configuration (Masked Credentials)
10  server = 'sql-fitflow-team8-2025.database.windows.net'
11  database = 'fitflow_oltp'
12  username = 'fitflowadmin'
13  password = '*****'
14
15  # Connection setup
16  params = urllib.parse.quote_plus(
17      f'Driver={{ODBC Driver 18 for SQL Server}};'
18      f'Server=tcp:{server},1433;Database={database};'
19      f'Uid={username};Pwd={password};Encrypt=yes;TrustServerCertificate=no;'
20  )
21  engine = create_engine(f'mssql+pyodbc:///?odbc_connect={params}')
22  fake = Faker()
23
24  # Data Generation Logic (Simplified for brevity)
25  # 1. Locations
26  locations = [{'location_id': f'LOC-{i+1:03d}', 'city': city}
27              for i, city in enumerate(['SF', 'Berkeley', 'Oakland', 'SJ', 'Palo Alto'])]
28  pd.DataFrame(locations).to_sql('locations', engine, if_exists='append', index=False)
29
30  # ... (Additional generation logic for all entities) ...
31  print(Data Generation Complete.)
```

Listing 3: Synthetic Data Generation

## A.3 Python Code for ETL Process

```
 1  import pandas as pd
 2  from sqlalchemy import create_engine, text
 3  from datetime import datetime, timedelta
 4  import urllib
 5
 6  # Configuration
 7  server = 'sql-fitflow-team8-2025.database.windows.net'
 8  username = 'fitflowadmin'
 9  password = '*****'
10
11  def get_engine(db_name):
12      params = urllib.parse.quote_plus(
13          f'Driver={{ODBC Driver 18 for SQL Server}};Server=tcp:{server},1433;'
14          f'Database={db_name};Uid={username};Pwd={password};Encrypt=yes;'
15      )
16      return create_engine(f'mssql+pyodbc:///?odbc_connect={params}')
17
18  oltp_engine = get_engine('fitflow_oltp')
19  olap_engine = get_engine('fitflow_olap')
20
21  # --- STEP 1: EXTRACT (Incremental & Lock-Free) ---
22  last_run = datetime.now() - timedelta(days=90)
23  df_payments = pd.read_sql(text(SELECT * FROM payments WHERE payment_date > :d),
24                            oltp_engine, params={d: last_run})
25  df_subs = pd.read_sql(SELECT * FROM subscriptions, oltp_engine)
26  df_plans = pd.read_sql(SELECT * FROM membership_plans, oltp_engine)
```

```
27
28  # --- STEP 2: TRANSFORM (Staging Area) ---
29  # Joins are performed in-memory (Pandas)
30  df_subs_clean = df_subs[['subscription_id', 'plan_id']]
31  df_rev = df_payments.merge(df_subs_clean, on='subscription_id', how='left')
32  df_rev = df_rev.merge(df_plans, on='plan_id', how='left')
33
34  df_rev['date_key'] = pd.to_datetime(df_rev['payment_date']).dt.strftime('%Y%m%d').astype(int)
35
36  fact_revenue = df_rev.groupby(['location_id', 'plan_id', 'date_key']).agg({
37      'amount': 'sum',
38      'payment_id': 'count',
39      'member_id': 'nunique'
40  }).rename(columns={'amount': 'total_revenue'}).reset_index()
41
42  # --- STEP 3: LOAD (OLAP) ---
43  fact_revenue.to_sql('Fact_Monthly_Revenue', olap_engine, if_exists='append', index=False)
44  print(ETL Completed Successfully.)
```

Listing 4: ETL Process

## A.4   Dashboard Implementation Code (Dash)

```
1   # =====================================================
2   # IMPORTS
3   # =====================================================
4   import pandas as pd
5   import urllib
6   from sqlalchemy import create_engine
7
8   import dash
9   from dash import dcc, html
10  from dash.dependencies import Input, Output
11  import plotly.express as px
12
13  # =====================================================
14  # DB CONNECTION
15  # =====================================================
16  server = 'sql-fitflow-team8-2025.database.windows.net'
17  database = 'fitflow_olap'
18  username = 'fitflowadmin'
19  password = '*****'
20
21  params = urllib.parse.quote_plus(
22      fDriver={{ODBC Driver 18 for SQL Server}};
23      fServer=tcp:{server},1433;
24      fDatabase={database};
25      fUid={username};Pwd={password};
26      fEncrypt=yes;TrustServerCertificate=no;Connection Timeout=30;
27  )
28
29  engine = create_engine(fmssql+pyodbc:///?odbc_connect={params})
30
31  # =====================================================
32  # LOAD DATA  VISUALIZATION 1 (REVENUE)
33  # =====================================================
34  query_revenue =
```

```
35  SELECT
36      d.year,
37      d.month,
38      CONCAT(d.year, '-', RIGHT('00' + CAST(d.month AS VARCHAR), 2)) AS year_month,
39      l.city AS location,
40      SUM(f.total_revenue) AS revenue
41  FROM Fact_Monthly_Revenue f
42  JOIN Dim_Date d ON f.date_key = d.date_key
43  JOIN Dim_Location l ON f.location_key = l.location_key
44  GROUP BY d.year, d.month, l.city
45  ORDER BY d.year, d.month;
46
47
48  df_rev = pd.read_sql(query_revenue, engine)
49  months_rev = sorted(df_rev[year_month].unique())
50
51  # =====================================================
52  # LOAD DATA  VISUALIZATION 2 (CLASS ATTENDANCE)
53  # =====================================================
54  query_attendance =
55  SELECT
56      d.year,
57      d.month,
58      CONCAT(d.year, '-', RIGHT('00' + CAST(d.month AS VARCHAR), 2)) AS year_month,
59      l.city AS location,
60      c.name AS class_type,
61      SUM(f.attendance_count) AS total_attendance,
62      AVG(f.attendance_count) AS avg_attendance
63  FROM Fact_Class_Attendance f
64  JOIN Dim_Date d ON f.date_key = d.date_key
65  JOIN Dim_Location l ON f.location_key = l.location_key
66  JOIN Dim_Class_Type c ON f.class_type_key = c.class_type_key
67  GROUP BY d.year, d.month, l.city, c.name
68  ORDER BY d.year, d.month;
69
70
71  df_att = pd.read_sql(query_attendance, engine)
72  months_att = sorted(df_att[year_month].unique())
73
74  # =====================================================
75  # LOAD DATA  VISUALIZATION 3 (HEATMAP)
76  # =====================================================
77  query_heatmap =
78  SELECT
79      DATENAME(month, DATEFROMPARTS(d.year, d.month, 1)) AS month_name,
80      l.city AS location,
81      SUM(f.attendance_count) AS total_attendance
82  FROM Fact_Class_Attendance f
83  JOIN Dim_Date d ON f.date_key = d.date_key
84  JOIN Dim_Location l ON f.location_key = l.location_key
85  GROUP BY d.year, d.month, l.city
86
87
88  df_heat = pd.read_sql(query_heatmap, engine)
89
90  month_order = [
91      January,February,March,April,May,June,
```

```
 92        July,August,September,October,November,December
 93  ]
 94
 95  df_heat[month_name] = pd.Categorical(
 96      df_heat[month_name],
 97      categories=month_order,
 98      ordered=True
 99  )
100
101  # =======================================================
102  # DASH APP (ONE APP ONLY)
103  # =======================================================
104  app = dash.Dash(__name__)
105
106  app.layout = html.Div([
107
108      html.H1(FitFlow Management Dashboard, style={textAlign: center}),
109
110      # ================================================
111      # VISUALIZATION 1  REVENUE
112      # ================================================
113      html.H2(Monthly Revenue Trends),
114
115      dcc.Dropdown(
116          id=rev-location-filter,
117          options=[{label: i, value: i} for i in sorted(df_rev[location].unique())],
118          multi=True,
119          placeholder=Select Location(s)
120      ),
121
122      dcc.RangeSlider(
123          id=rev-time-filter,
124          min=0,
125          max=len(months_rev) - 1,
126          value=[0, len(months_rev) - 1],
127          marks={i: months_rev[i] for i in range(len(months_rev))},
128          step=1
129      ),
130
131      dcc.Graph(id=revenue-line),
132
133      html.Hr(),
134
135      # ================================================
136      # VISUALIZATION 2  CLASS ATTENDANCE
137      # ================================================
138      html.H2(Attendance by Class Type),
139
140      dcc.Dropdown(
141          id=att-location-filter,
142          options=[{label: i, value: i} for i in sorted(df_att[location].unique())],
143          placeholder=Select Location,
144          clearable=True
145      ),
146
147      dcc.RadioItems(
148          id=metric-toggle,
```

```
149          options=[
150              {label: Total Attendance, value: total},
151              {label: Average Attendance, value: average}
152          ],
153          value=total,
154          inline=True
155      ),
156
157      dcc.RangeSlider(
158          id=att-time-filter,
159          min=0,
160          max=len(months_att) - 1,
161          value=[0, len(months_att) - 1],
162          marks={i: months_att[i] for i in range(0, len(months_att), max(1, len(months_att)//6))},
163          step=1
164      ),
165
166      dcc.Graph(id=attendance-bar),
167
168      html.Hr(),
169
170      # ================================================
171      # VISUALIZATION 3  HEATMAP
172      # ================================================
173      html.H2(Seasonality in Class Attendance),
174
175      dcc.Dropdown(
176          id=heat-location-filter,
177          options=[{label: loc, value: loc} for loc in sorted(df_heat[location].unique())],
178          placeholder=Select Location (optional),
179          clearable=True,
180          style={width: 40%}
181      ),
182
183      dcc.Graph(id=attendance-heatmap)
184 ])
185
186 # ====================================================
187 # CALLBACKS
188 # ====================================================
189
190 @app.callback(
191     Output(revenue-line, figure),
192     Input(rev-location-filter, value),
193     Input(rev-time-filter, value)
194 )
195 def update_revenue(locations, time_range):
196     dff = df_rev.copy()
197     dff = dff[dff[year_month].between(
198         months_rev[time_range[0]], months_rev[time_range[1]]
199     )]
200
201     if locations:
202         fig = px.line(dff[dff[location].isin(locations)],
203                       x=year_month, y=revenue,
204                       color=location, markers=True)
205     else:
```

```
206            dff = dff.groupby(year_month, as_index=False)[revenue].sum()
207            fig = px.line(dff, x=year_month, y=revenue, markers=True)
208
209        fig.update_xaxes(categoryorder=array, categoryarray=months_rev)
210        fig.update_layout(xaxis_title=Month, yaxis_title=Revenue)
211        return fig
212
213
214    @app.callback(
215        Output(attendance-bar, figure),
216        Input(att-location-filter, value),
217        Input(metric-toggle, value),
218        Input(att-time-filter, value)
219    )
220    def update_attendance(location, metric, time_range):
221        dff = df_att.copy()
222        dff = dff[dff[year_month].between(
223            months_att[time_range[0]], months_att[time_range[1]]
224        )]
225
226        if location:
227            dff = dff[dff[location] == location]
228
229        if metric == total:
230            dff = dff.groupby(class_type, as_index=False)[total_attendance].sum()
231            y = total_attendance
232        else:
233            dff = dff.groupby(class_type, as_index=False)[avg_attendance].mean()
234            y = avg_attendance
235
236        fig = px.bar(dff, x=y, y=class_type, orientation=h)
237        return fig
238
239
240    @app.callback(
241        Output(attendance-heatmap, figure),
242        Input(heat-location-filter, value)
243    )
244    def update_heatmap(location):
245        dff = df_heat.copy()
246        if location:
247            dff = dff[dff[location] == location]
248
249        fig = px.density_heatmap(
250            dff,
251            x=month_name,
252            y=location,
253            z=total_attendance,
254            color_continuous_scale=Blues
255        )
256
257        fig.update_coloraxes(colorbar_title=Total Class Attendance)
258        return fig
259
260
261    # ====================================================
262    # RUN
```

```
263  # ======================================================
264  if __name__ == __main__:
265      app.run(debug=True)
```

Listing 5: Dash Application Code

## A.5    Analytical Processes Code

```
1   # ======================================================
2   # Informative Analysis 1: Monthly Revenue by Plan
3   # ======================================================
4   import pandas as pd
5   import numpy as np
6   import matplotlib.pyplot as plt
7   from statsmodels.tsa.statespace.sarimax import SARIMAX
8
9   jdbc_url = (
10      jdbc:sqlserver://sql-fitflow-team8-2025.database.windows.net:1433;
11      database=fitflow_olap;
12      encrypt=true;
13      trustServerCertificate=false;
14  )
15
16  connection_properties = {
17      user: fitflowadmin,
18      password: *****,
19      driver: com.microsoft.sqlserver.jdbc.SQLServerDriver
20  }
21
22  df_mr.createOrReplaceTempView(Fact_Monthly_Revenue)
23  df_ca.createOrReplaceTempView(Fact_Class_Attendance)
24  df_date.createOrReplaceTempView(Dim_Date)
25  df_plan.createOrReplaceTempView(Dim_Plan)
26
27  spark_df = spark.sql(
28  SELECT
29      date_format(d.full_date, 'yyyy-MM') AS Date,
30      f.total_revenue,
31      p.plan_name
32  FROM Fact_Monthly_Revenue f
33  JOIN Dim_Plan p
34      ON f.plan_key = p.plan_key
35  RIGHT JOIN Dim_Date d
36      ON f.date_key = d.date_key
37  WHERE d.full_date >= '2024-12-01'
38    AND d.full_date <= '2025-12-31'
39  )
40
41  df1 = spark_df.toPandas()
42  df1 = df1.dropna() # Remove NaNs
43
44  total_rev_per_mon = df1.groupby(['Date', 'plan_name'])['total_revenue'].sum().reset_index()
45  total_rev_per_mon['total_revenue'] = round(total_rev_per_mon['total_revenue'], 2)
46
47  total_rev_per_mon[Date] = pd.to_datetime(total_rev_per_mon[Date])
48  full_months = pd.date_range(start=2024-12-01, end=2025-12-01, freq=MS)
49
```

```python
50  plans = total_rev_per_mon[plan_name].unique()
51  full_df = (
52      total_rev_per_mon
53      .set_index([Date, plan_name])
54      .unstack(plan_name)
55      .reindex(full_months)
56      .stack(plan_name)
57      .reset_index()
58  )
59  full_df.columns = [Date, plan_name, total_revenue]
60
61  plt.figure(figsize=(10,6))
62  for plan in plans:
63      tmp = full_df[full_df[plan_name] == plan]
64      plt.plot(
65          tmp[Date],
66          tmp[total_revenue],
67          marker=.,
68          linewidth=2,
69          label=plan
70      )
71
72  plt.xticks(full_months, [d.strftime(%Y-%m) for d in full_months])
73  plt.xlabel(Month)
74  plt.ylabel(Total Revenue)
75  plt.title(Monthly Revenue Contribution by Plan)
76  plt.legend()
77  plt.grid(alpha=0.2)
78  plt.tight_layout()
79  plt.show()
80
81  # =======================================================
82  # Informative Analysis 2: Daily Utilization Forecast
83  # =======================================================
84  spark_df2 = spark.sql(
85  SELECT
86      d.full_date AS Date,
87      f.utilization_rate
88  FROM Fact_Class_Attendance f
89  RIGHT JOIN Dim_Date d
90      ON f.date_key = d.date_key
91  WHERE d.full_date >= '2024-12-01'
92    AND d.full_date <= '2025-12-14'
93  )
94
95  df2 = spark_df2.toPandas()
96  df2['utilization_rate'] = df2['utilization_rate'] * 100
97  df2[Date] = pd.to_datetime(df2[Date])
98
99  df2 = df2.groupby(Date, as_index=False).mean()
100 df2 = df2.set_index(Date).sort_index()
101
102 full_dates = pd.date_range(
103     start=df2.index.min(),
104     end=df2.index.max(),
105     freq=D
106 )
```

```python
107  df2 = df2.reindex(full_dates)
108  df2[utilization_rate] = df2[utilization_rate].interpolate(method=time)
109  df2[utilization_rate] = df2[utilization_rate].ffill().bfill().round(4)
110  df2 = df2.reset_index().rename(columns={index: Date})
111
112  # Modeling
113  df = df2.copy()
114  df[Date] = pd.to_datetime(df[Date])
115  df = df.set_index(Date).sort_index()
116  y = df[utilization_rate]
117
118  model = SARIMAX(
119      y,
120      order=(1,1,1), # ARIMA part
121      seasonal_order=(1, 0, 1, 30), # Monthly seasonality
122      enforce_stationarity=False,
123      enforce_invertibility=False
124  )
125  result = model.fit(disp=False)
126
127  # In-sample fitted values
128  fitted = result.fittedvalues
129
130  # Forecast to 2026-01-15
131  forecast_horizon = pd.date_range(start=2025-12-15, end=2026-01-15, freq=D)
132  forecast = result.get_forecast(steps=len(forecast_horizon))
133  forecast_mean = forecast.predicted_mean
134
135  plt.figure(figsize=(12,6))
136  # Actual values
137  plt.plot(y.index, y, label=Actual, color=green, linewidth=1)
138  # Fitted values
139  plt.plot(fitted.index, fitted, label=Fitted, color=blue, linestyle=solid, alpha=0.8)
140  # Forecast values
141  plt.plot(forecast_horizon, forecast_mean, label=Forecast, color=red, linewidth=2)
142
143  plt.axvline(pd.to_datetime(2025-12-14), color=gray, linestyle=:, label=Forecast Start)
144  plt.title(Daily Utilization Rate Forecast (SARIMA))
145  plt.xlabel(Date)
146  plt.ylabel(Utilization Rate (%))
147  plt.legend()
148  plt.grid(alpha=0.2)
149  plt.tight_layout()
150  plt.show()
151
152  # =======================================================
153  # Informative Analysis 3: Daily Revenue Forecast
154  # =======================================================
155  spark_df_rev = spark.sql(
156  SELECT
157      d.full_date AS Date,
158      f.total_revenue
159  FROM Fact_Monthly_Revenue f
160  RIGHT JOIN Dim_Date d
161      ON f.date_key = d.date_key
162  WHERE d.full_date >= '2024-12-01'
163    AND d.full_date <= '2025-12-14'
```

```
164  )
165
166  df_rev = spark_df_rev.toPandas()
167
168  # 1. Convert Date to datetime
169  df_rev[Date] = pd.to_datetime(df_rev[Date])
170  # 2. Aggregate revenue by day
171  df_rev = df_rev.groupby(Date, as_index=False).sum()
172  # 3. Set Date as index and sort
173  df_rev = df_rev.set_index(Date).sort_index()
174  # 4. Create full daily date range
175  full_dates = pd.date_range(start=df_rev.index.min(), end=df_rev.index.max(), freq=D)
176  # 5. Reindex to make dates continuous
177  df_rev = df_rev.reindex(full_dates)
178  # 6. Time-based linear interpolation
179  df_rev[total_revenue] = df_rev[total_revenue].interpolate(method=time)
180  # 7. Handle edge NaNs
181  df_rev[total_revenue] = df_rev[total_revenue].ffill().bfill()
182  # 8. Restore Date column
183  df_rev = df_rev.reset_index().rename(columns={index: Date})
184
185  # Modeling
186  df = df_rev.copy()
187  df[Date] = pd.to_datetime(df[Date])
188  df = df.set_index(Date).sort_index()
189  y = df[total_revenue]
190
191  model = SARIMAX(
192      y,
193      order=(2,0,2),
194      seasonal_order=(1,1,1,7),
195      enforce_stationarity=False,
196      enforce_invertibility=False
197  )
198  result = model.fit(disp=False)
199
200  fitted = result.fittedvalues
201  forecast_horizon = pd.date_range(start=2025-12-15, end=2026-01-15, freq=D)
202  forecast = result.get_forecast(steps=len(forecast_horizon))
203  forecast_mean = forecast.predicted_mean
204
205  plt.figure(figsize=(12,6))
206  plt.plot(y.index, y, label=Actual Revenue, color=green, linewidth=1)
207  plt.plot(fitted.index, fitted, label=Fitted Revenue, color=blue, linestyle=solid, alpha=0.8)
208  plt.plot(forecast_horizon, forecast_mean, label=Forecast Revenue, color=red, linewidth=1)
209  plt.axvline(pd.to_datetime(2025-12-14), color=gray, linestyle=:, label=Forecast Start)
210  plt.title(Daily Revenue Forecast (SARIMA))
211  plt.xlabel(Date)
212  plt.ylabel(Total Revenue)
213  plt.legend()
214  plt.grid(alpha=0.2)
215  plt.tight_layout()
216  plt.show()
```

Listing 6: Analytical Processes (Data Prep & SARIMA)

# FitFlow Fitness

INDENG 215 Final Project
Team 8

Changxin Shi, Arthur Fang, Da Eun Kim, Eesha Danish, Junqi Huang, Junyu Li

# Business Overview

■ FitFlow Fitness

- Multi-location gym & wellness chain (30+ locations)

- Offers in-person classes, personal training, and digital workouts

- Subscription-based membership model

- Serves busy professionals, students, and fitness enthusiasts

■ Why This Business?

- High transaction volume

- Strong need for both real-time operations and historical analytics

# Customers, Revenue, Data Needs

## ■ Customers

- Individual members (Basic, Plus, Premium)

- Corporate memberships

- Trainers & instructors (internal users)

## ■ Revenue Streams

- Membership subscriptions

- Class & personal training add-ons

- Retail and digital services

## ■ Data Needs

- Operational: check-ins, bookings, billing

- Analytical: revenue trends, utilization, class performance, location comparisons

# OLTP vs OLAP: Why Separate?

## ■ Transactional (OLTP)

- High-frequency writes
- Real-time consistency (ACID)
- Supports check-ins, bookings, and billing

## ■ Analytical (OLAP)

- Read-heavy, historical data
- Aggregated across time, location, and plan
- Used for trend analysis and strategic decisions

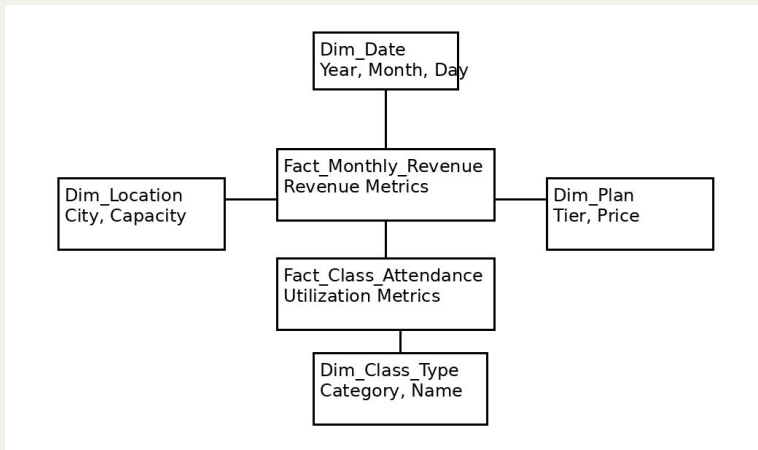Key Reason: Prevents analytical queries from slowing down live operations

# OLTP Schema Design

**Core Entities**

- Members, Locations, Membership Plans

- Subscriptions, Payments

- Class Sessions, Bookings, Check-ins, Trainers

# OLAP Schema Design (Star Schema)



■ **Dimension Tables**

- Dim Date (Year, Month, Day, Holiday flag)

- Dim Location (City, State, Capacity)

- Dim Plan (Plan name, tier, price)

- Dim Class Type (Category, name)

# Fact Table & Granularity

■ **Fact Monthly Revenue**

- Granularity: Location × Plan × Month

- Measures:

  - Total Revenue
  - Transaction Count
  - Active Members

■ **Fact Class Attendance**

- Granularity: Location × Class Type × Day

- Measures:

  - Total Bookings
  - Attendance Count
  - Utilization Rate (Attendance / Capacity)

# ETL Process Overview + Workflow

## ■ **Design Principles**

- No joins on OLTP system

- Transformations performed in a staging layer

- Incremental extraction logic supported

## ■ **Why This Matters**

- Prevents locks on live operational tables

- Preserves OLTP performance

## ■ **ETL Workflow**

1. Extract
   - Raw tables extracted individually
   - No joins or aggregations on OLTP
   - Supports incremental loads
2. Transform
   - Joins and cleaning performed in pandas
   - Invalid records filtered
   - Dates converted into surrogate date keys
3. Load
   - Dimension tables loaded first
   - Fact tables loaded after to ensure referential integrity

# Dashboard View 1: Monthly Revenue by Plan

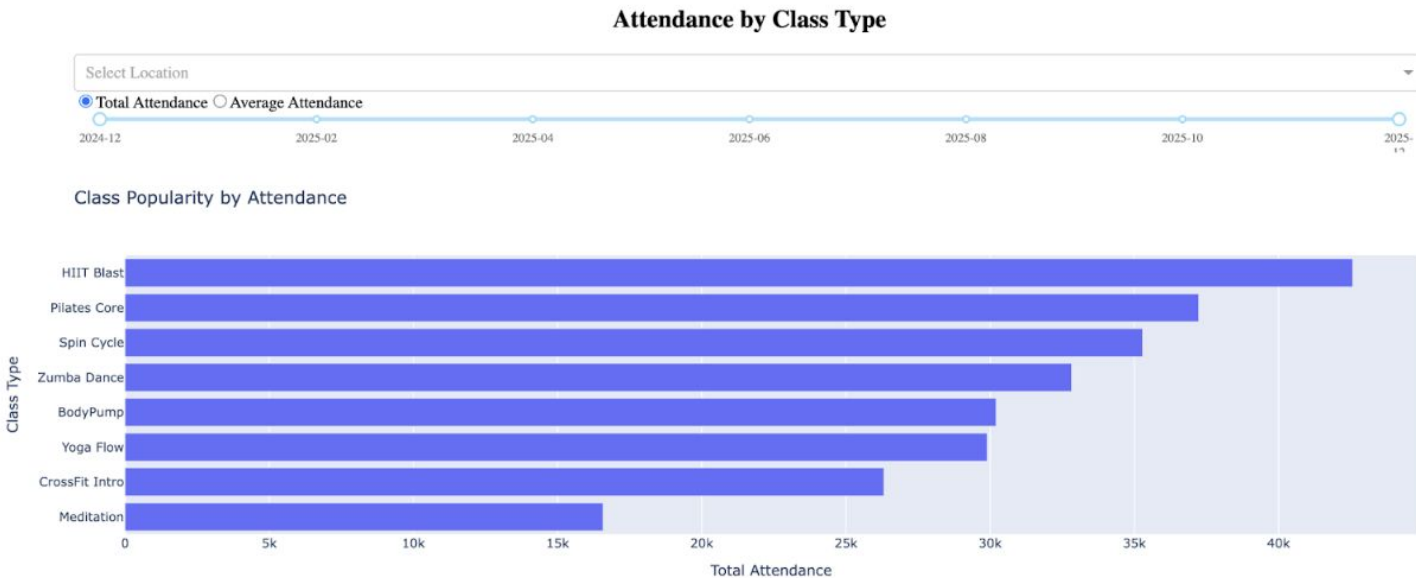# Dashboard View 1: Monthly Revenue by Plan

**Plot 1.1: Monthly Revenue Trends** *filtered by locations (Berkeley, Palo Alto, Oakland, and San Francisco) and range of months*



Note: When no location is selected, the chart aggregates revenue across all locations to provide an overall business view.
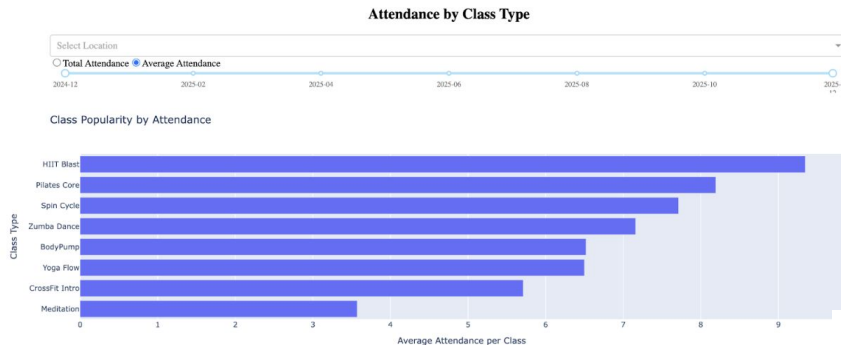
# Dashboard View 2: Class Attendance & Utilization
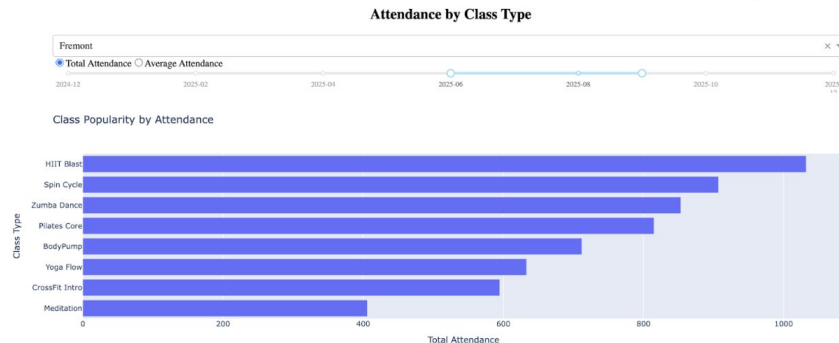


**Plot 2: Total Attendance by Class Type**

Attendance by Class Type

Select Location

⦿ Total Attendance  ○ Average Attendance

2024-12          2025-02          2025-04          2025-06          2025-08          2025-10          2025-12

Class Popularity by Attendance

HIIT Blast
Pilates Core
Spin Cycle
Zumba Dance
BodyPump
Yoga Flow
CrossFit Intro
Meditation

Class Type

0        5k        10k        15k        20k        25k        30k        35k        40k

Total Attendance

# Dashboard View 2: Class Attendance & Utilization



Plot 2.1: Average Attendance by Class Type



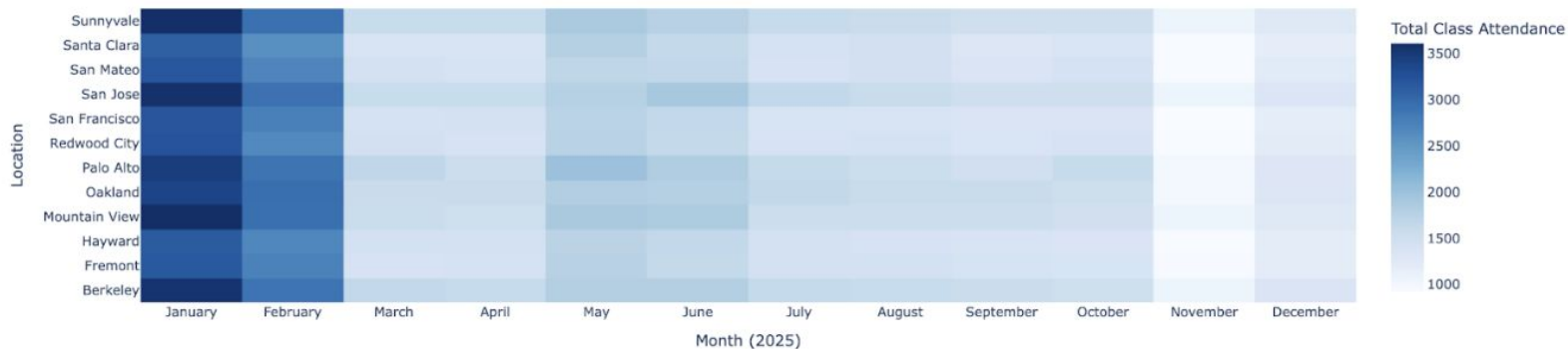Plot 2.2: Total Attendance by Class Type *filtered by location (Fremont) and months (Q3, 2025)*

# Dashboard View 3: Seasonality Analysis



**Plot 3: Seasonality Analysis of Class Attendance**

**Seasonality Analysis: Attendance Heatmap**
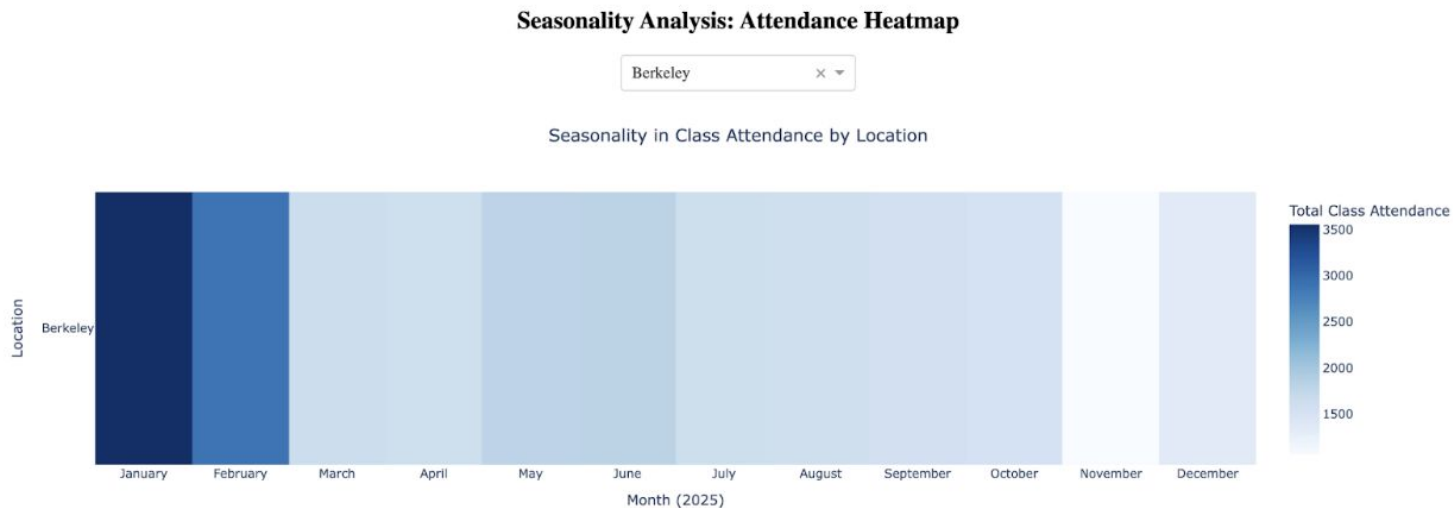
Select a location (optional) ▼

Seasonality in Class Attendance by Location

Darker colors indicate higher attendance. January spikes highlight new-year participation behavior.
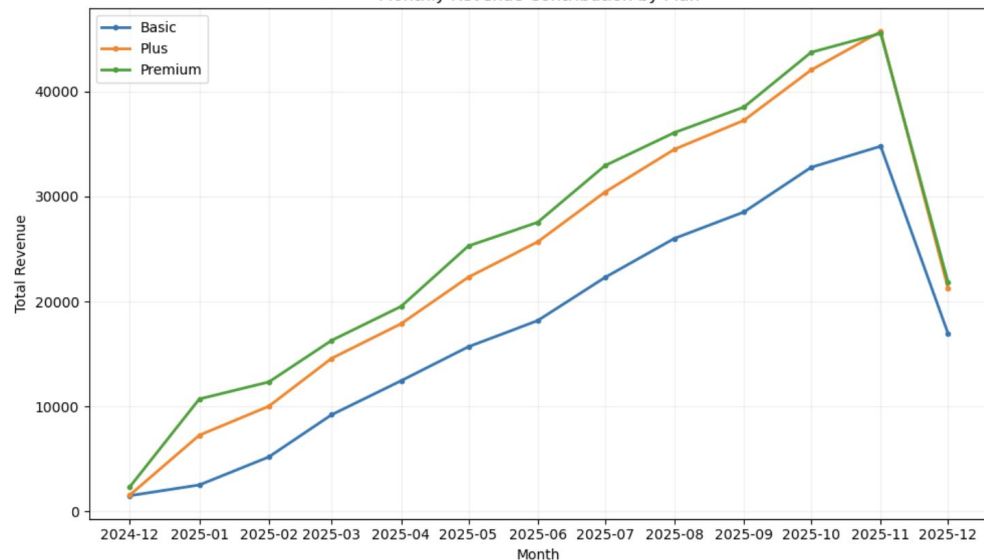
13

# Dashboard View 3: Seasonality Analysis



**Plot 3.1: Seasonality Analysis of Class Attendance** *filtered by location (Berkeley)*

**Seasonality Analysis: Attendance Heatmap**

Berkeley ✕ ▾

Seasonality in Class Attendance by Location

Total Class Attendance
3500
3000
2500
2000
1500

Location — Berkeley

January  February  March  April  May  June  July  August  September  October  November  December

Month (2025)

Darker colors indicate higher attendance. January spikes highlight new-year participation behavior.

# Analytical Process 1: Monthly Revenue by Plan
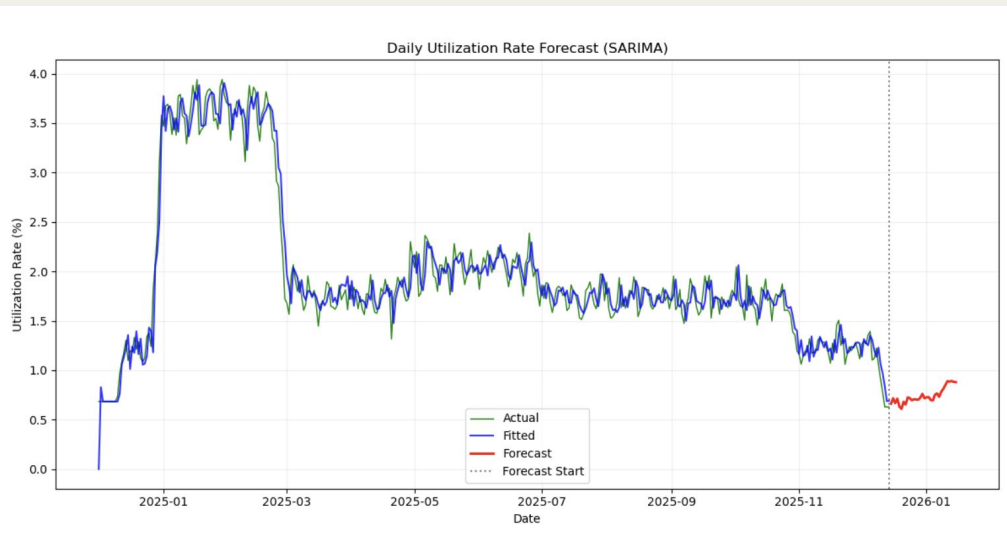


Monthly Revenue Contribution by Plan

- This plot forecasts short-term revenue trends based on historical revenue patterns.

Why it matters

- Helps management anticipate future revenue fluctuations and plan budgets accordingly.

- Useful for short-term financial planning and performance monitoring.
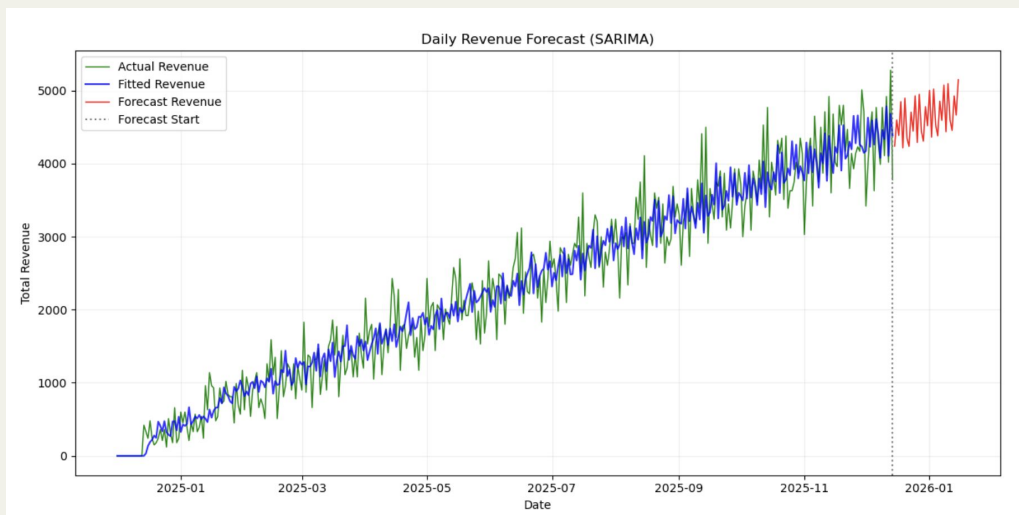
# Analytical Process 2: Utilization Forecasting



Daily Utilization Rate Forecast (SARIMA)

- This plot forecasts class attendance and utilization over time.

Why it matters

- Helps identify periods of high or low demand.

- Supports staffing decisions, class scheduling, and capacity planning.

# Analytical Process 3: Daily Revenue Trend & Forecasting



- This plot analyzes historical daily total revenue and forecasts short-term future revenue trends using time-series modeling.

Why it matters

- Daily revenue reflects real-time customer activity and overall financial performance.

- Helps management monitor trends and anticipate short-term revenue changes.

# Conclusion & Key Takeaways

## ■ What We Built

- Designed a production-style OLTP system to support real-time gym operations

- Implemented a star-schema OLAP system for scalable analytics

- Built a Python-based ETL pipeline that preserves OLTP performance

- Delivered an interactive management dashboard connected to the OLAP database

## ■ Why It matters

- Separating transactional and analytical workloads prevents performance bottlenecks

- Star schema enables fast, intuitive reporting across time, location, and plans

- Dashboard allows managers to make data-driven decisions without writing SQL

## ■ Project Outcome

- Demonstrates a complete data lifecycle: operations → ETL → analytics (dashboard & forecasting) → decision support