# Netflix Case study

This case study focuses on leveraging Apache Cassandra to manage a Netflix dataset containing user subscriptions, which includes details such as user demographics, subscription types, revenue, and device usage. It explores key aspects like creating a keyspace with appropriate replication strategies, designing a schema with suitable data types, and defining primary and clustering keys for efficient data retrieval. The study also covers CRUD operations (insert, update, delete), batch operations for optimizing data handling, and deriving actionable analytics.

**Key Analytics**

1. <u>Total Revenue by Country:</u> Analyzing total revenue contribution from different countries.
2. <u>User Count by Device:</u> Understanding which devices are most popular among users.
3. <u>Subscription Type Popularity:</u> Tracking the popularity of different subscription plans.
4. <u>Churn Rate Analysis:</u> Identifying users who haven't made a recent payment (using last_payment_date).
5. <u>Demographic Insights:</u> Understanding the age distribution and gender ratios of users.

| User ID | Subscription Type | Monthly Revenue | Join Date | Last Payment Date | Country | Age | Gender | Device | Plan Duration |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Basic | 10 | 15-01-22 | 10/6/2023 | United States | 28 | Male | Smartphone | 1 Month |
| 2 | Premium | 15 | 5/9/2021 | 22-06-23 | Canada | 35 | Female | Tablet | 1 Month |
| 3 | Standard | 12 | 28-02-23 | 27-06-23 | United Kingdom | 42 | Male | Smart TV | 1 Month |
| 4 | Standard | 12 | 10/7/2022 | 26-06-23 | Australia | 51 | Female | Laptop | 1 Month |
| 5 | Basic | 10 | 1/5/2023 | 28-06-23 | Germany | 33 | Male | Smartphone | 1 Month |
| 6 | Premium | 15 | 18-03-22 | 27-06-23 | France | 29 | Female | Smart TV | 1 Month |
| 7 | Standard | 12 | 9/12/2021 | 25-06-23 | Brazil | 46 | Male | Tablet | 1 Month |
| 8 | Basic | 10 | 2/4/2023 | 24-06-23 | Mexico | 39 | Female | Laptop | 1 Month |
| 9 | Standard | 12 | 20-10-22 | 23-06-23 | Spain | 37 | Male | Smartphone | 1 Month |
| 10 | Premium | 15 | 7/1/2023 | 22-06-23 | Italy | 44 | Female | Smart TV | 1 Month |
| 11 | Basic | 10 | 16-05-22 | 22-06-23 | United States | 31 | Female | Smartphone | 1 Month |
| 12 | Premium | 15 | 23-03-23 | 28-06-23 | Canada | 45 | Male | Tablet | 1 Month |
| 13 | Standard | 12 | 30-11-21 | 27-06-23 | United Kingdom | 48 | Female | Laptop | 1 Month |
| 14 | Basic | 10 | 1/8/2022 | 26-06-23 | Australia | 27 | Male | Smartphone | 1 Month |
| 15 | Standard | 12 | 9/5/2023 | 28-06-23 | Germany | 38 | Female | Smart TV | 1 Month |
| 16 | Premium | 15 | 7/4/2022 | 27-06-23 | France | 36 | Male | Tablet | 1 Month |
| 17 | Basic | 10 | 24-01-22 | 25-06-23 | Brazil | 30 | Female | Laptop | 1 Month |
| 18 | Standard | 12 | 18-10-21 | 24-06-23 | Mexico | 43 | Male | Smartphone | 1 Month |
| 19 | Premium | 15 | 15-02-23 | 23-06-23 | Spain | 32 | Female | Smart TV | 1 Month |
| 20 | Basic | 10 | 27-05-23 | 22-06-23 | Italy | 41 | Male | Tablet | 1 Month |
| 21 | Premium | 15 | 10/6/2023 | 22-06-23 | United States | 26 | Female | Laptop | 1 Month |
| 22 | Basic | 10 | 22-07-22 | 28-06-23 | Canada | 34 | Male | Smartphone | 1 Month |
| 23 | Standard | 12 | 5/12/2021 | 27-06-23 | United Kingdom | 49 | Female | Smart TV | 1 Month |
| 24 | Standard | 12 | 3/4/2022 | 26-06-23 | Australia | 31 | Male | Tablet | 1 Month |
| 25 | Basic | 10 | 14-03-23 | 28-06-23 | Germany | 40 | Female | Laptop | 1 Month |

# 1. Keyspace Creation

A keyspace in Cassandra is similar to a database in other RDBMS. It defines the scope for data replication and provides a namespace for tables.

**Example:**

CREATE KEYSPACE netflix_analytics

WITH replication = {'class': 'SimpleStrategy', 'replication_factor': '1'};

**Explanation:**

class: Defines the replication strategy (e.g., SimpleStrategy or NetworkTopologyStrategy).

replication_factor: The number of replicas for each data point.

# 2. Create a Table

Creating a table involves defining the data types for each column.

**Example:**

```
CREATE TABLE netflix_analytics.subscriptions (
    user_id UUID PRIMARY KEY,
    subscription_type TEXT,
    monthly_revenue DECIMAL,
    join_date DATE,
    last_payment_date DATE,
    country TEXT,
    age INT,
    gender TEXT,
    device TEXT,
    plan_duration INT
);
```

**Explanation:**

UUID: Universally Unique Identifier, A 128-bit unique identifier ideal for globally unique values.

TEXT: Stores variable-length strings for text-based data. Used for string data (subscription_type, country, gender, device).

DECIMAL: Represents fixed-point numbers, suitable for precise financial values. For monetary values (monthly_revenue).

DATE: Stores only the date component (year, month, day). For date fields (join_date, last_payment_date).

INT: A 32-bit signed integer for numerical data without decimals. For integers (age, plan_duration).

# 3. Defining Primary Key, Partition Key, Clustering

Primary Key: A unique identifier for rows in a table.

Partition Key: Determines the distribution of data across nodes.

Clustering Columns: Define the order of data storage within a partition.

**Example:**

```
CREATE TABLE netflix_analytics.subscriptions (
    user_id UUID,
    country TEXT,
    device TEXT,
    subscription_type TEXT,
    monthly_revenue DECIMAL,
    join_date DATE,
    last_payment_date DATE,
    age INT,
    gender TEXT,
    plan_duration INT,
    PRIMARY KEY ((country, device), user_id)
);
```

**Explanation:**

Partition Key: (country, device) ensures data distribution based on country and device.

Clustering Key: user_id defines the ordering within the partition.

# 4. Insert Command and Batch Insert

Inserts data into the table.

**Single Insert Example:**

```
INSERT INTO netflix_analytics.subscriptions (
    user_id, country, device, subscription_type, monthly_revenue,
    join_date, last_payment_date, age, gender, plan_duration
) VALUES (
    uuid(), 'USA', 'Mobile', 'Premium', 15.99, '2023-01-01', '2024-01-01', 25, 'M', 12
```

);

**Batch Insert Example:**

BEGIN BATCH

INSERT INTO netflix_analytics.subscriptions (user_id, country, device, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration)

VALUES (uuid(), 'USA', 'Mobile', 'Premium', 15.99, '2023-01-01', '2024-01-01', 25, 'M', 12);


INSERT INTO netflix_analytics.subscriptions (user_id, country, device, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration)

VALUES (uuid(), 'Canada', 'Laptop', 'Standard', 12.99, '2023-02-01', '2024-02-01', 30, 'F', 6);

APPLY BATCH;


# 5. Update Command and Batch Update

Updates existing records.


## Single Update Example:

UPDATE netflix_analytics.subscriptions

SET monthly_revenue = 19.99

WHERE country = 'USA' AND device = 'Mobile' AND user_id = some-uuid;


## Batch Update Example:

BEGIN BATCH

UPDATE netflix_analytics.subscriptions

SET monthly_revenue = 19.99

WHERE country = 'USA' AND device = 'Mobile' AND user_id = some-uuid;


UPDATE netflix_analytics.subscriptions

SET plan_duration = 18

WHERE country = 'Canada' AND device = 'Laptop' AND user_id = another-uuid;

APPLY BATCH;

# 6. Delete Command

Removes records from the table.

**Example:**

DELETE FROM netflix_analytics.subscriptions

WHERE country = 'USA' AND device = 'Mobile' AND user_id = some-uuid;


## 7. Analytics Derived from the Dataset

You can derive various insights such as:

**Monthly Revenue Per Country:**

SELECT country, SUM(monthly_revenue) AS total_revenue

FROM netflix_analytics.subscriptions

GROUP BY country, device;


**Active Users by Device Type:**

SELECT device, COUNT(*) AS user_count

FROM netflix_analytics.subscriptions

GROUP BY country, device;


**Average Age of Subscribers:**

SELECT AVG(age) AS average_age

FROM netflix_analytics.subscriptions;


**Subscription Types Distribution:**

SELECT subscription_type, COUNT(*) AS subscription_count

FROM netflix_analytics.subscriptions

GROUP BY country, device;


**Revenue Trends Over Time:**

SELECT join_date, SUM(monthly_revenue) AS daily_revenue

FROM netflix_analytics.subscriptions

GROUP BY country, device;

**User-Based Analysis within Country and Device**

SELECT user_id, subscription_type, monthly_revenue, join_date

FROM netflix_analytics.subscriptions

WHERE country = 'USA' AND device = 'Mobile'

ORDER BY user_id;

**Tracking User Subscription Changes Over Time**

SELECT user_id, subscription_type, monthly_revenue, join_date

FROM netflix_analytics.subscriptions

WHERE country = 'USA' AND device = 'Mobile'

ORDER BY user_id DESC;

# More Data Types

| Data Type | Purpose | Example |
|---|---|---|
| **int** | Stores a 32-bit integer value. | age int (Stores age as an integer) |
| **bigint** | Stores a 64-bit signed integer value (2^64-1). | amount bigint (Stores large monetary values) |
| **float** | Stores a 32-bit single precision floating-point number. | price float (Stores a price like 19.99) |
| **double** | Stores a 64-bit double precision floating-point number. | temperature double (Stores temperature with higher precision) |
| **text** | Stores a string of text with variable length. | title text (Stores book title, e.g., "The Catcher in the Rye") |
| **varchar** | Stores a variable-length string. | name varchar (Stores employee name, e.g., "John Doe") |
| **uuid** | Stores a universally unique identifier (UUID). | order_id uuid (Stores order ID as UUID) |
| **timeuuid** | Stores a UUID that is time-based, useful for generating ordered UUIDs. | event_id timeuuid (Stores event ID based on time) |
| **boolean** | Stores a boolean value (true/false). | is_active boolean (Stores if a user is active or not) |
| **date** | Stores a date (year, month, day) without the time component. | meeting_date date (Stores date, e.g., "2024-11-08") |
| **timestamp** | Stores a date and time, including milliseconds. | log_time timestamp (Stores timestamp, e.g., "2024-11-08 14:30:00") |
| **time** | Stores the time of day without the date. | start_time time (Stores time, e.g., "14:30:00") |
| **inet** | Stores an IP address (either IPv4 or IPv6). | ip_address inet (Stores IP address, e.g., "192.168.1.1") |
| **decimal** | Stores arbitrary precision decimal numbers. | price decimal (Stores price with precision for financial calculations) |
| **set<type>** | Stores a collection of unique values of a specified type. | subjects set<text> (Stores unique subjects, e.g., {"Math", "Science"}) |
| **list<type>** | Stores an ordered collection of values of a specified type. | chapters list<text> (Stores an ordered list of chapters, e.g., ['Introduction', 'Chapter 1', 'Chapter 2']) |
| **map<key_type, value_type>** | Stores a collection of key-value pairs. | preferences map<text, text> (Stores key-value pairs, e.g., {"color": "blue", "language": "English"}) |
| **User-Defined Type (UDT)** | Groups related fields together as a custom data type. | address frozen<address> (Defines address UDT for users, stores street, city, state, zip_code) |

# Example

CREATE TYPE address (

    street text,

    city text,

    state text,

```
    zip_code text
);
CREATE TABLE example_table (
    id UUID,
    age int,
    amount bigint,
    price float,
    temperature double,
    title text,
    name varchar,
    order_id uuid,
    event_id timeuuid,
    is_active boolean,
    meeting_date date,
    log_time timestamp,
    start_time time,
    ip_address inet,
    price_decimal decimal,
    subjects set<text>,
    chapters list<text>,
    preferences map<text, text>,
    address frozen<address>,
    PRIMARY KEY (id, log_time)
);
```

# Insert 10 Records

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 25, 5000, 19.99, 22.5, 'To Kill a Mockingbird', 'Alice Smith', uuid(), now(), true, '2024-11-08', '2024-11-08 09:30:00', '09:30:00', '192.168.1.100', 29.99, {'Math', 'English'}, ['Introduction', 'Chapter 1'], {'color': 'red', 'language': 'English'}, {street: '456 Oak St', city: 'Denver', state: 'CO', zip_code: '80202'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 40, 100000, 29.99, 35.8, '1984', 'Bob Johnson', uuid(), now(), false, '2024-11-09', '2024-11-09 10:00:00', '10:00:00', '192.168.1.101', 39.99, {'Science', 'History'}, ['Introduction', 'Chapter 1', 'Chapter 2'], {'color': 'blue', 'language': 'French'}, {street: '789 Pine St', city: 'Los Angeles', state: 'CA', zip_code: '90001'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 28, 20000, 15.49, 18.7, 'Pride and Prejudice', 'Catherine Lee', uuid(), now(), true, '2024-11-10', '2024-11-10 12:30:00', '12:30:00', '192.168.1.102', 24.99, {'History', 'Literature'}, ['Introduction', 'Chapter 1', 'Chapter 2'], {'color': 'green', 'language': 'Spanish'}, {street: '123 Maple Ave', city: 'Chicago', state: 'IL', zip_code: '60606'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 35, 300000, 99.99, 40.1, 'The Great Gatsby', 'David Green', uuid(), now(), true, '2024-11-11', '2024-11-11 14:00:00', '14:00:00', '192.168.1.103', 49.99, {'Math', 'Art'}, ['Introduction', 'Chapter 1'], {'color': 'yellow', 'language': 'German'}, {street: '321 Birch Rd', city: 'Miami', state: 'FL', zip_code: '33101'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 50, 150000, 9.99, 30.2, 'Moby Dick', 'Eve White', uuid(), now(), false, '2024-11-12', '2024-11-12 15:00:00', '15:00:00', '192.168.1.104', 59.99, {'Art', 'Philosophy'}, ['Introduction', 'Chapter 1', 'Chapter 2'], {'color': 'purple', 'language': 'Italian'}, {street: '654 Elm St', city: 'New York', state: 'NY', zip_code: '10001'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 26, 25000, 25.99, 20.0, 'Frankenstein', 'Grace Brown', uuid(), now(), true, '2024-11-13', '2024-11-13 16:30:00', '16:30:00', '192.168.1.105', 34.99, {'Biology', 'Chemistry'}, ['Introduction', 'Chapter 1', 'Chapter 2'], {'color': 'black', 'language': 'Portuguese'}, {street: '987 Cedar St', city: 'Austin', state: 'TX', zip_code: '73301'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 45, 50000, 19.49, 25.5, 'Jane Eyre', 'Helen Carter', uuid(), now(), false, '2024-11-14', '2024-11-14 17:00:00', '17:00:00', '192.168.1.106', 39.99, {'Music', 'Literature'}, ['Introduction', 'Chapter 1'], {'color': 'orange', 'language': 'English'}, {street: '234 Oak St', city: 'Boston', state: 'MA', zip_code: '02108'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 33, 70000, 17.99, 32.0, 'Wuthering Heights', 'Fiona Clark', uuid(), now(), true, '2024-11-15', '2024-11-15 18:00:00', '18:00:00', '192.168.1.107', 44.99, {'Science', 'Philosophy'}, ['Introduction', 'Chapter 1'], {'color': 'pink', 'language': 'Arabic'}, {street: '321 Pine Rd', city: 'Seattle', state: 'WA', zip_code: '98101'});

INSERT INTO example_table (id, age, amount, price, temperature, title, name, order_id, event_id, is_active, meeting_date, log_time, start_time, ip_address, price_decimal, subjects, chapters, preferences, address) VALUES (uuid(), 38, 120000, 35.99, 28.4, 'War and Peace', 'Isaac Williams', uuid(), now(), false, '2024-11-16', '2024-11-16 19:00:00', '19:00:00', '192.168.1.108', 79.99, {'Engineering', 'History'}, ['Introduction', 'Chapter 1', 'Chapter 2'], {'color': 'brown', 'language': 'Russian'}, {street: '876 Birch St', city: 'San Francisco', state: 'CA', zip_code: '94102'});

# Display Records

select name, age, address, subjects, chapters, preferences, address from example_table;

# Alter Table

ALTER TABLE example_table Drop age;

ALTER TABLE example_table ADD age INT;

# Insert Data from External Sources and To External Sources

**Export Data to CSV**

COPY example_table TO 'example_table.csv' WITH HEADER = true;

**Export Data From CSV**

**Step1: Do the basic configuration in Docker**

**Mounting a folder in the docker**

docker run -d --name data_mount -v D:/Projects/HeroVired_Online_Teaching/Lectures/Cassandra_Training/cassandra/datasets:/data cassandra


**Execute the cqlsh command to get the terminal**

docker exec -it data_mount cqlsh

**Step2. Let's First Create the Table**

CREATE KEYSPACE netflix with replication = {'class': 'SimpleStrategy', 'replication_factor':1};

CREATE TABLE netflix .netflix_shows (

   show_id TEXT,

   type TEXT,

   title TEXT,

   director TEXT,

   cast TEXT,

   country TEXT,

   date_added TEXT,

   release_year TEXT,

   rating TEXT,

   duration TEXT,

   listed_in TEXT,

   description TEXT,

   PRIMARY KEY (show_id, type)

);

**Insert data from the external sources 'CSV File'**

COPY netflix.netflix_shows (show_id, type, title, director, cast, country, date_added, release_year, rating, duration, listed_in, description) FROM 'data/netflix_titles.csv' WITH HEADER = true;

# Few Records for the analytics

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Mobile', uuid(), 'Standard', 12.99, '2023-01-10', '2023-12-01', 28, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Mobile', uuid(), 'Basic', 8.99, '2022-11-12', '2023-11-10', 35, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Laptop', uuid(), 'Premium', 16.99, '2021-08-22', '2023-12-05', 42, 'M', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Tablet', uuid(), 'Standard', 12.99, '2023-06-17', '2023-12-03', 29, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Mobile', uuid(), 'Basic', 8.99, '2023-02-01', '2023-11-28', 25, 'F', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Smart TV', uuid(), 'Premium', 16.99, '2022-12-05', '2023-11-30', 37, 'M', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Tablet', uuid(), 'Standard', 12.99, '2023-03-10', '2023-12-04', 22, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Smart TV', uuid(), 'Basic', 8.99, '2021-07-11', '2023-11-22', 30, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Mobile', uuid(), 'Premium', 16.99, '2023-01-21', '2023-12-01', 32, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Laptop', uuid(), 'Standard', 12.99, '2022-04-10', '2023-11-29', 26, 'F', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Smart TV', uuid(), 'Basic', 8.99, '2023-02-01', '2023-12-05', 27, 'M', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Tablet', uuid(), 'Premium', 16.99, '2021-05-15', '2023-12-01', 39, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Smart TV', uuid(), 'Standard', 12.99, '2023-03-12', '2023-12-02', 33, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Laptop', uuid(), 'Basic', 8.99, '2022-09-05', '2023-11-10', 26, 'M', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Smart TV', uuid(), 'Standard', 12.99, '2023-01-19', '2023-12-03', 32, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Tablet', uuid(), 'Premium', 16.99, '2022-06-14', '2023-11-30', 40, 'M', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Mobile', uuid(), 'Basic', 8.99, '2023-08-11', '2023-12-05', 25, 'M', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Mobile', uuid(), 'Standard', 12.99, '2023-01-28', '2023-12-05', 28, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Smart TV', uuid(), 'Premium', 16.99, '2021-11-15', '2023-12-02', 34, 'M', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Laptop', uuid(), 'Basic', 8.99, '2023-07-17', '2023-11-30', 29, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Tablet', uuid(), 'Standard', 12.99, '2023-01-06', '2023-12-03', 31, 'F', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Smart TV', uuid(), 'Premium', 16.99, '2022-03-25', '2023-12-02', 41, 'M', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Mobile', uuid(), 'Standard', 12.99, '2023-01-10', '2023-12-01', 28, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Mobile', uuid(), 'Basic', 8.99, '2022-11-12', '2023-11-10', 35, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Laptop', uuid(), 'Premium', 16.99, '2021-08-22', '2023-12-05', 42, 'M', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Tablet', uuid(), 'Standard', 12.99, '2023-06-17', '2023-12-03', 29, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Mobile', uuid(), 'Basic', 8.99, '2023-02-01', '2023-11-28', 25, 'F', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Smart TV', uuid(), 'Premium', 16.99, '2022-12-05', '2023-11-30', 37, 'M', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Tablet', uuid(), 'Standard', 12.99, '2023-03-10', '2023-12-04', 22, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Smart TV', uuid(), 'Basic', 8.99, '2021-07-11', '2023-11-22', 30, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Mobile', uuid(), 'Premium', 16.99, '2023-01-21', '2023-12-02', 32, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Laptop', uuid(), 'Standard', 12.99, '2022-04-10', '2023-11-29', 26, 'F', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Smart TV', uuid(), 'Basic', 8.99, '2023-02-01', '2023-12-05', 27, 'M', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Tablet', uuid(), 'Premium', 16.99, '2021-05-15', '2023-12-01', 39, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Smart TV', uuid(), 'Standard', 12.99, '2023-03-12', '2023-12-02', 33, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Laptop', uuid(), 'Basic', 8.99, '2022-09-05', '2023-11-10', 26, 'M', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Smart TV', uuid(), 'Standard', 12.99, '2023-01-19', '2023-12-03', 32, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Tablet', uuid(), 'Premium', 16.99, '2022-06-14', '2023-11-30', 40, 'M', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Mobile', uuid(), 'Basic', 8.99, '2023-08-11', '2023-12-05', 25, 'M', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Mobile', uuid(), 'Standard', 12.99, '2023-01-28', '2023-12-05', 28, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Smart TV', uuid(), 'Premium', 16.99, '2021-11-15', '2023-12-02', 34, 'M', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Laptop', uuid(), 'Basic', 8.99, '2023-07-17', '2023-11-30', 29, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Tablet', uuid(), 'Standard', 12.99, '2023-01-06', '2023-12-03', 31, 'F', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Smart TV', uuid(), 'Premium', 16.99, '2022-03-25', '2023-12-02', 41, 'M', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Mobile', uuid(), 'Basic', 8.99, '2023-07-14', '2023-12-05', 24, 'M', 12);

```sql
INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Smart TV', uuid(), 'Premium', 16.99, '2023-05-11', '2023-11-29', 38, 'F', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Mobile', uuid(), 'Standard', 12.99, '2022-02-18', '2023-12-03', 31, 'F', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Mobile', uuid(), 'Premium', 16.99, '2023-04-22', '2023-12-01', 29, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Tablet', uuid(), 'Premium', 16.99, '2022-10-30', '2023-11-25', 33, 'M', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Mobile', uuid(), 'Standard', 12.99, '2023-02-07', '2023-11-26', 36, 'M', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Smart TV', uuid(), 'Premium', 16.99, '2023-03-09', '2023-11-15', 44, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Australia', 'Laptop', uuid(), 'Basic', 8.99, '2021-12-19', '2023-11-25', 31, 'M', 36);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('India', 'Smart TV', uuid(), 'Basic', 8.99, '2023-05-20', '2023-11-29', 28, 'M', 6);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Tablet', uuid(), 'Standard', 12.99, '2023-08-15', '2023-12-01', 26, 'F', 12);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('Canada', 'Laptop', uuid(), 'Premium', 16.99, '2021-09-13', '2023-12-02', 45, 'M', 24);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('UK', 'Tablet', uuid(), 'Standard', 12.99, '2023-03-03', '2023-12-05', 33, 'F', 18);

INSERT INTO netflix_analytics.subscriptions (country, device, user_id, subscription_type, monthly_revenue, join_date, last_payment_date, age, gender, plan_duration) VALUES ('USA', 'Smart TV', uuid(), 'Basic', 8.99, '2023-04-09', '2023-11-26', 22, 'M', 6);
```