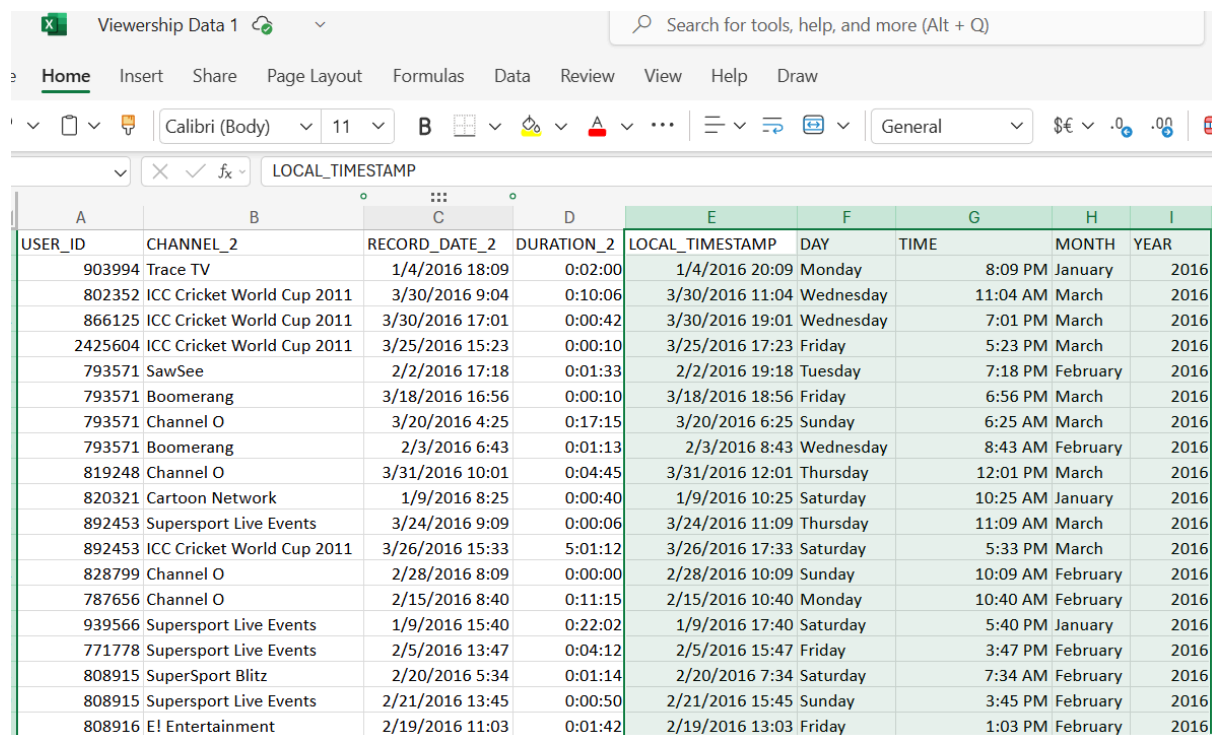# BrightTV Viewership Methodology

# Contents

# 1.Background and Introduction

BrightTV 's CEO has an objective to grow the company's subscription base for this financial year. He has approached you to provide insights that would assist CVM (Customer Value Management) team in meeting this year's objective.

The dataset provided (User_Profiles and Viewership) contains information on the user profiles and viewer transactions for BrightTV. The two files were loaded onto Databricks (SQL) for further analysis.

# 2. Date Manipulation

The viewership file contains a time stamp for each record. Times and dates in the dataset were supplied in UTC and have been converted to SA time. The Day, Time, Month, and the Year were extracted from the Timestamp using excel before loading the file onto Databricks. New columns were added as shown below:



Highlighted columns feature day, time,month and year converted to SA standard time.

# 3. Completeness of Data

Following the ingestion of the transformed Viewership table as well as the User_Profiles, using SQL queries the number of records in each file was extracted. Data cleaning was performed in case of duplicates, empty rows or missing files.

# 3.1. Check the number of records

The number of total records from  the Viewership table is **10000**, including duplicates. The number of total records from User_Profiles is **5375.** The SQL queries were ran on Databricks to extract this data:

```sql
-- Query to obtain the number of records
SELECT COUNT(*)
FROM user_profiles;

SELECT COUNT(DISTINCT userid)
FROM user_profiles;

SELECT COUNT(*)
FROM viewership;

SELECT COUNT(userid)
FROM viewership;
```

# 3.2. Checking for Duplicates

The following query ran on Databricks to check rows which contained duplicated data.

```sql
-- Query to check for completely duplicates rows

SELECT *,
    COUNT(*)
FROM user_profiles
GROUP BY ALL
HAVING COUNT(*) > 1;

SELECT *,
    COUNT(*)
FROM viewership
GROUP BY ALL
HAVING COUNT(*) > 1; -- 5 records have duplicates
```

The Use_Profile table does not contain duplicates. However, Viewership contained **5 duplicated rows**. A new temporary table, **Viewership_new** has been created using the query below to retrieve a new table without duplication. The new file has **9995 unique rows**.

```
-- Query to create a temporary table with no duplicates as viewership_new
WITH viewership_new AS (
    SELECT DISTINCT *
    FROM viewership
    GROUP BY ALL)
    SELECT *
    FROM viewership;
```

# 3.3 Checking and Replacing Missing Values

The query below was used to check for missing values:

```
-- Query to check for missing values in the tables

SELECT * FROM user_profiles
WHERE userid IS NULL OR NAME IS NULL OR surname IS NULL OR email IS NULL OR gender IS NULL OR RACE IS NULL OR AGE IS
NULL OR PROVINCE IS NULL OR SOCIAL_MEDIA_HANDLE IS NULL;


SELECT * FROM viewership_new
WHERE user_id IS NULL OR channel_2 IS NULL OR record_date_2 IS NULL OR duration_2 IS NULL OR local_timestamp IS NULL OR
month IS NULL OR year IS NULL;
```

 The viewership file did not contain any missing values. However, User_Profiles contained missing values. These records were replaced with "None" using the query shown below. A new table was created to account for this transformation. Using the CASE statement, this table was bucketed into distinct age groups.

```
-- Query to replace missing records with 'None' and creating a temp table
WITH user_profiles_new AS (
    SELECT
        userid,
        age,
        IFNULL(name, 'None') AS Name,
        IFNULL(surname, 'None') AS Surname,
        IFNULL(email, 'None') AS email,
        IFNULL(gender, 'None') AS Gender,
        IFNULL(race, 'None') AS Race,
        IFNULL(province, 'None') AS Province,
        IFNULL(social_media_handle, 'None') AS social_media_handle,
        CASE
            WHEN age BETWEEN 1 AND 12 THEN 'Younger than 13'
            WHEN age BETWEEN 13 AND 25 THEN '13 to 25'
            WHEN age BETWEEN 26 AND 44 THEN '26 to 44'
            WHEN age >= 45 THEN '45 and older'
            ELSE 'Not Specified'
        END AS Age_group
    FROM user_profiles
    GROUP BY ALL)
SELECT * FROM user_profiles;
```

# 3.4. Joining The Two Working Tables

Next, we join the two tables since data completeness is completed, the tables are joined using **INNER JOIN**. This was done to create a new comprehensive file that displays the users that watched the channels as well as the programs watched. The tables  User_Profiles and Viewership_New were joined using **UserID as a common column**.

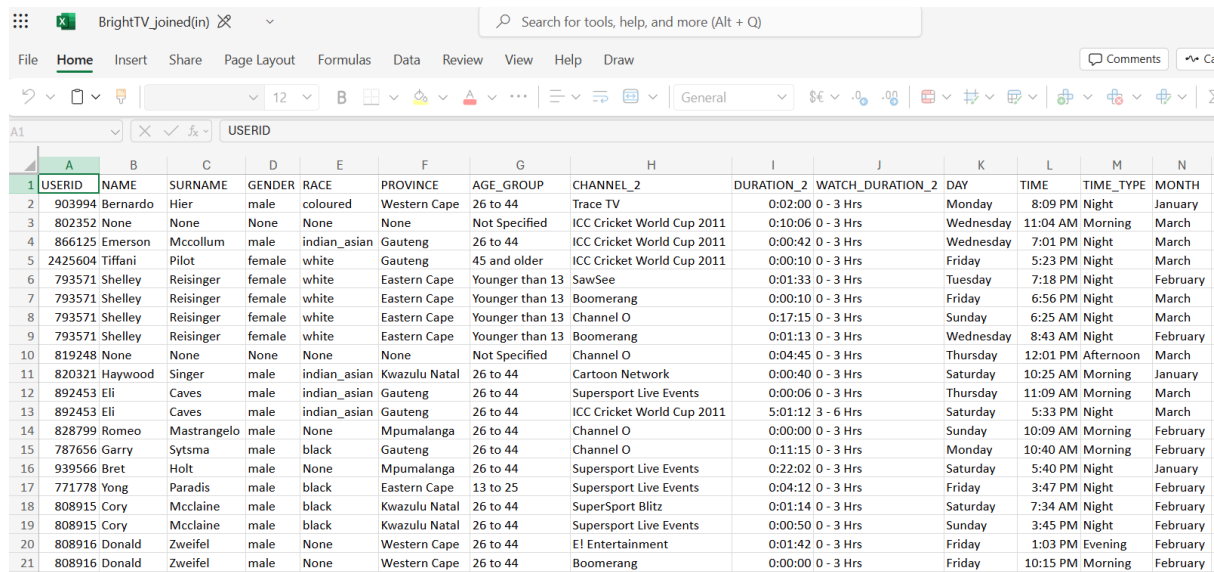Viewership Duration and Time of Day were bucketed into new columns upon the joining of the two tables, as indicated by the query below.

```sql
SELECT
    u.userid,
    u.Name,
    u.Surname,
    u.Gender,
    u.Race,
    u.Province,
    u.Age_group,
    v.channel2,
    v.duration_2,
    CASE
        WHEN v.duration_2 between '00:00:00' AND '02:59:59' THEN '0 - 3 Hrs'
        WHEN v.duration_2 between '03:00:00' AND '05:59:59' THEN '3 - 6 Hrs'
        WHEN v.duration_2 between '06:00:00' AND '08:59:59' THEN '6 - 9 Hrs'
        ELSE '9 - 12 Hrs'
    END AS Watch_Duration,
    v.day,
    v.time,
    CASE
        WHEN v.time between '06:00:00' AND '11:59:59' THEN 'Morning'
        WHEN v.time between '12:00:00' AND '17:59:59' THEN 'Afternoon'
        WHEN v.time between '18:00:00' AND '23:59:59' THEN 'Evening'
        ELSE 'Night'
    END AS Time_Type,
    v.month
FROM user_profiles_new AS u
INNER JOIN viewership_new AS v ON u.userid = v.userid;
```

The file has been converted to CSV format for further analysis, contains all columns needed for analysis.

# 4. Analysis

Display of the final exported table:



From the above table, pivot tables and visuals were developed to perform the following analysis:

Demographic Analysis
- ★ Viewership by Race
- ★ Viewership by Gender
- ★ Viewership by Age
- ★ Viewership by Province

Trend Analysis
- ★ Total Viewers by Weekday
- ★ Total Views Over Time

Channel Analysis
- ★ Top 10 most watched channels
- ★ Viewership by Day and Duration

# 5. Pivot Table Samples

| CHANNEL | ⌄ | GRAND TOTAL |
|---|---|---|
| Africa Magic | | 857 |
| Boomerang | | 714 |
| Cartoon Network | | 793 |
| Channel O | | 1048 |
| CNN | | 505 |
| E! Entertainment | | 367 |
| ICC Cricket World Cup 2011 | | 1465 |
| SuperSport Blitz | | 896 |
| Supersport Live Events | | 1637 |
| Trace TV | | 952 |
| **Grand Total** | | **9234** |

| RACE | ⌄ | Grand Total |
|---|---|---|
| | | 0.10% |
| BLACK | | 43.31% |
| COLOURED | | 16.32% |
| INDIAN_ASIAN | | 15.76% |
| NONE | | 10.58% |
| OTHER | | 1.02% |
| WHITE | | 12.92% |
| **Grand Total** | | **100.00%** |