

## Methodology:

### Checking duplicates for user\_profile table:

```
1  --Checking duplicates from user profiles table
2  SELECT userid, COUNT(*) AS record_count
3  FROM user_profiles
4  GROUP BY userid
5  HAVING COUNT(*) > 1;
```

	USERID	RECORD_COUNT
Query produced no results		

User\_profile table has no duplicates.

### Checking duplicates for viewership table:

```
6  --Checking duplicates from viewership table
7  SELECT userid, channel2, recorddate2, COUNT(*) AS record_count
8  FROM viewership
9  GROUP BY userid, channel2, recorddate2
10 HAVING COUNT(*) > 1;
```

#	USERID	CHANNEL2	RECORDDATE2	RECORD_COUNT
1	808015	Channel O	2016/02/19 04:30	2
2	808330	Africa Magic	2016/01/11 17:56	2
3	809422	Africa Magic	2016/02/21 16:46	2
4	2204356	Dstv Events 1	2016/02/03 11:22	2
5	931362	ICC Cricket World Cup 2011	2016/03/30 13:49	2
6	786910	Channel O	2016/01/08 20:47	2
7	810044	Supersport Live Events	2016/02/18 07:59	2

Checking and fixing duplicates in your dataset is essential to guaranteeing data quality, accuracy, and reliable analysis as they can provide skewed results, misleading insights and wasted resources.

### Removing duplicates from viewership table and creating cleaned temp table:

```
11  --Removing viewership duplicates (keeping 1 record for eah duplicate)
12  CREATE OR REPLACE TEMP TABLE viewership_deduped AS
13  SELECT *
14  FROM (
15    SELECT *,
16    ROW_NUMBER() OVER (PARTITION BY userid, channel2, recorddate2 ORDER BY duration_2 DESC) AS rn
17  FROM viewership
18  ) AS sub
19  WHERE rn = 1;
```

status
1 Table VIEWERSHIP_DEDUPED successfully created.

```
2  DELETE FROM viewership
3  WHERE (userid, channel2, recorddate2, duration_2) NOT IN (
4    SELECT userid, channel2, recorddate2, duration_2
5    FROM (
6      SELECT *,
7      ROW_NUMBER() OVER (PARTITION BY userid, channel2, recorddate2 ORDER BY duration_2 DESC) AS rn
8    FROM viewership
9    WHERE rn = 1
10 );
11 SELECT*
```

#	number of rows deleted
2	

Cleaning data and joining the user profile and viewership table:

```
1  --CLEANING DATA and joining the tables
2  create or replace Temporary table viewership_temp_tbl AS (
3  SELECT
4  A.Userid,
5  A.Name,
6  A.Surname,
7  A.Email,
8  CASE
9  WHEN A.Gender IS NULL OR A.Gender = 'None' THEN 'other'
10 ELSE A.Gender
11 END AS gender,
12 CASE
13 WHEN A.Race IS NULL OR A.Race = 'None' THEN 'other'
14 ELSE A.Race
15 END AS Race,
16 A.Age,
17 CASE
18 WHEN A.Province IS NULL OR A.Province = 'None' THEN 'other'
19 ELSE A.Province
20 END AS Province,
21 B.Channel2 AS Channel,
22 B.Recorddate2 AS Date,
23 B.Duration_2 AS Duration
24 FROM
25 user_profiles AS A
26 INNER JOIN VIEWERSHIP AS B
27 ON A.userid=B.userid
28 )
29 --
```

Results

Chart

status
1 Table VIEWERSHIP_TEMP_TBL successfully created.

Converting UTC to SA time:

```
60  -- Converting UTC to SA time
61  UPDATE VIEWERSHIP_TEMP_TBL
62  SET date= CONVERT_TIMEZONE('UTC', 'Africa/Johannesburg', TO_TIMESTAMP(date, 'YYYY/MM/DD HH24:MI'));
63
64
```

Results

Chart

# number of rows updated	# number of multi-joined rows updated
10000	0

Adding time column and timestamp:

```
64  UPDATE VIEWERSHIP_TEMP_TBL
65  SET date = CONVERT_TIMEZONE('UTC', 'Africa/Johannesburg', TO_TIMESTAMP(date, 'YYYY-MM-DD HH24:MI:SS.FF3'));

66  --Adding timestamp to time
67  UPDATE VIEWERSHIP_TEMP_TBL
68  SET Time= TO_CHAR(TO_TIMESTAMP(date, 'YYYY-MM-DD HH24:MI:SS.FF3'), 'HH24:MI');
```

Results

Chart

# number of rows updated	# number of multi-joined rows updated
10000	0

## Adding date column and timestamp:

```
69 --Adding date column
70 ALTER TABLE VIEWERSHIP_TEMP_TBL
71 ADD COLUMN Date2 DATE;
72 --adding timestamp to date
73 UPDATE VIEWERSHIP_TEMP_TBL
74 SET Date2= TO_CHAR(TO_TIMESTAMP(date, 'YYYY-MM-DD HH24:MI:SS.FF3'), 'YYYY-MM-DD');
```

→ Results ↗ Chart

# number of rows updated	# number of multi-joined rows updated
10000	0

## QUESTION 1:

### Age buckets

```
85 --QUESTION 1
86 --create age buckets
87 SELECT
88 CASE
89   WHEN age = 0 THEN 'Not specified'
90   WHEN age BETWEEN 1 AND 12 THEN 'Kids'
91   WHEN age BETWEEN 13 AND 19 THEN 'Teenagers'
92   WHEN age BETWEEN 20 AND 24 THEN 'Young Adults'
93   WHEN age BETWEEN 25 AND 34 THEN 'Early Adulthood'
94   WHEN age BETWEEN 35 AND 44 THEN 'Mid Adulthood'
95   WHEN age BETWEEN 45 AND 54 THEN 'Mature Adults'
96   ELSE 'Seniors'
97 END AS age_group,
98 COUNT (*) AS viewer_count
99 FROM VIEWERSHIP_TEMP_TBL
100 GROUP BY age_group
101 ORDER BY viewer_count DESC;
```

↶ Results ↗ Chart

	AGE_GROUP	# VIEWER_COUNT	
1	Early Adulthood	4148	Q
2	Mid Adulthood	2635	Q
3	Young Adults	1228	R
4	Mature Adults	927	Q

### Time buckets

```
102 --Create time of day buckets
103 SELECT
104 CASE
105   WHEN TO_TIME(Time) BETWEEN '00:00:00' AND '05:59:59' THEN 'Late Night'
106   WHEN TO_TIME(Time) BETWEEN '06:00:00' AND '11:59:59' THEN 'Morning'
107   WHEN TO_TIME(Time) BETWEEN '12:00:00' AND '17:59:59' THEN 'Afternoon'
108   WHEN TO_TIME(Time) BETWEEN '18:00:00' AND '23:59:59' THEN 'Evening'
109   ELSE 'Unknown'
110 END AS time_bucket,
111 COUNT(*) AS total_views,
112 FROM VIEWERSHIP_TEMP_TBL
113 GROUP BY time_bucket
114 ORDER BY total_views DESC;
```

↶ Results ↗ Chart

	TIME_BUCKET	# TOTAL_VIEWS	
1	Evening	3807	Q
2	Afternoon	3085	Q
3	Late Night	2352	R
4	Morning	756	Q

Group by Province:

115 --Group by province  
116 SELECT  
117 province,  
118 COUNT(\*) AS total\_views,  
119 FROM VIEWERSHIP\_TEMP\_TBL  
120 GROUP BY province  
121 ORDER BY total\_views DESC;

ResultsChart

	PROVINCE	TOTAL_VIEWS
1	Gauteng	3654
2	Western Cape	1845
3	Kwazulu Natal	1001
4	Mpumalanga	918

Most viewed channels

122 --Most viewed channels  
123 SELECT  
124 channel,  
125 COUNT(\*) AS total\_views,  
126 FROM VIEWERSHIP\_TEMP\_TBL  
127 GROUP BY channel  
128 ORDER BY total\_views DESC;

ResultsChart

	CHANNEL	TOTAL_VIEWS
1	Supersport Live Events	1638
2	ICC Cricket World Cup 2011	1465
3	Channel O	1050
4	Trace TV	952

Views per day:

129 --views per day  
130 SELECT  
131 date2,  
132 COUNT(\*) AS views,  
133 FROM VIEWERSHIP\_TEMP\_TBL  
134 GROUP BY date2  
135 ORDER BY date2;

ResultsChart

	DATE2	VIEWS
1	2016-01-01	51
2	2016-01-02	79
3	2016-01-03	63
4	2016-01-04	70

## Grouping by gender and age

136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
...

--Grouping gender and age  
SELECT  
 gender,  
 CASE  
 WHEN age = 0 THEN 'Not specified'  
 WHEN age BETWEEN 1 AND 12 THEN 'Kids'  
 WHEN age BETWEEN 13 AND 19 THEN 'Teenagers'  
 WHEN age BETWEEN 20 AND 24 THEN 'Young Adults'  
 WHEN age BETWEEN 25 AND 34 THEN 'Early Adulthood'  
 WHEN age BETWEEN 35 AND 44 THEN 'Mid Adulthood'  
 WHEN age BETWEEN 45 AND 54 THEN 'Mature Adults'  
 ELSE 'Seniors'  
 END AS age\_group,  
 COUNT(\*) AS viewer\_count,  
FROM VIEWERSHIP\_TEMP\_TBL  
GROUP BY gender, age\_group  
ORDER BY age\_group, gender;

ResultsChart

GENDER	AGE_GROUP	VIEWER_COUNT
female	Early Adulthood	417
male	Early Adulthood	3731
female	Kids	26
male	Kids	73

## Top viewers by watch time:

160  
161  
162  
163  
164  
165  
166  
167  
...

--top users by watch time  
SELECT  
 userid,  
 SUM(DATE\_PART('minute', Duration) \* 60 + DATE\_PART('second', Duration)) / 60 AS total\_minutes\_watched  
FROM VIEWERSHIP\_TEMP\_TBL  
GROUP BY userid  
ORDER BY total\_minutes\_watched DESC  
LIMIT 10;

ResultsChart

USERID	TOTAL_MINUTES_WATCHED
2415121	337.066667
2000272	226.033333
810145	197.250000
789220	196.700000