```
In [0]:
```

```
from pyspark.sql.types import *
from pyspark import SparkContext
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("my project 1").getOrCreate()
sc = spark.sparkContext

# Read a CSV into a dataframe
# There is a smarter version, that will first check if there is a Parquet file and use it
def load_PD_file(filename_or_dir, schema) :
    dataPath = "/mnt/ddscoursedatastorage/fwm-stb-data/" + filename_or_dir
    df = spark.read.format("csv")\
        .option("header", "false")\
        .option("delimiter", "|")\
        .schema(schema)\
        .load(dataPath)
    return df
```

In [0]:

```
df1 = spark.read.csv("/mnt/ddscoursedatastorage/dds-students/test.csv")
```

Reading the data files

Reading Refrence Data

In [0]:

```
# Reading the Reference Parguet files
ref data = spark.read.parquet('/ref data raw').withColumnRenamed(" device-id", "device id"
) \
                                                  .withColumnRenamed(" dma", "dma") \
                                                  .withColumnRenamed("dma-code", "dma cod
e")\
                                                  .withColumnRenamed(" household-id", "hou
sehold id") \
                                                  .withColumnRenamed(" household-type", "h
ousehold type") \
                                                  .withColumnRenamed(" system-type", "syst
em type") \
                                                  .withColumnRenamed(" zipcode", "zipcode"
ref data count = ref data.count()
print(ref data count)
ref data
```

203581233

Out[6]: DataFrame[device_id: string, dma: string, dma_code: bigint, household_id: bigint, household_type: string, system_type: string, zipcode: bigint]

Reading Daily Program Data

```
daily_prog_data.show()
+----+
```

Reading Program Viewing Data

```
In [0]:
```

```
#Reading the 1% sample of the viewing data from a Parquet file
viewing_data = spark.read.parquet('/sample_viewing_2_5percent')
print(f'There are {viewing_data.count():,} entries in viewing_data dataframe!')
```

There are 130,289,194 entries in viewing_data dataframe!

Reading Demographic Data

```
# Reading the Demographic CSV file
demographic schema = StructType([StructField('household id',StringType()),
                      StructField('household size', IntegerType()),
                      StructField('num adults', IntegerType()),
                      StructField('num generations', IntegerType()),
                      StructField('adult_range', StringType()),
                      StructField('marital status', StringType()),
                      StructField('race code', StringType()),
                      StructField('presence children', StringType()),
                      StructField('num children', IntegerType()),
                      StructField('age_children', StringType()), #format like range - 'bi
twise'
                      StructField('age range children', StringType()),
                      StructField('dwelling type',StringType()),
                      StructField('home_owner_status',StringType()),
                      StructField('length_residence', IntegerType()),
                      StructField('home market value',StringType()),
                      StructField('num vehicles', IntegerType()),
                      StructField('vehicle_make', StringType()),
                      StructField('vehicle model', StringType()),
                      StructField('vehicle_year', IntegerType()),
                      StructField('net_worth', IntegerType()),
                      StructField('income', StringType()),
                      StructField('gender_individual', StringType()),
                      StructField('age individual', IntegerType()),
                      StructField('education highest', StringType()),
                      StructField('occupation highest', StringType()),
                      StructField('education 1', StringType()),
                      StructField('occupation 1', StringType()),
                      StructField('age 2', IntegerType()),
                      StructField('education 2', StringType()),
                      StructField('occupation 2', StringType()),
                      StructField('age_3',IntegerType()),
                      StructField('education 3',StringType()),
                      StructField('occupation_3',StringType()),
                      StructField('age_4', IntegerType()),
```

```
StructField('education_4',StringType()),
StructField('occupation_4',StringType()),
StructField('age_5',IntegerType()),
StructField('education_5',StringType()),
StructField('occupation_5',StringType()),
StructField('polit_party_regist',StringType()),
StructField('polit_party_input',StringType()),
StructField('household_clusters',StringType()),
StructField('insurance_groups',StringType()),
StructField('financial_groups',StringType()),
StructField('green_living',StringType())
])
demographic_data = load_PD_file("demographic/" , demographic_schema)
```

Question 1.1

In [0]:

```
from pyspark.sql.functions import *
# Removing non-relevant columns and duplicated entries
ref data Q1 = ref data.distinct().select('device id', 'household id', 'dma')
daily prog data Q1 = daily prog data.distinct().drop('title')
viewing data Q1 = viewing data.distinct().select('prog code', 'device id', 'event date')
# changing string to a number according the requries we've been given
# changing null strings to None so that when converting the 'income' column values to int
they will still considered null
demographic data = demographic data.withColumn('income', when(demographic data.income ==
'A', '10')
                                                         .when(demographic data.income =
= 'B', '11')
                                                         .when(demographic data.income =
= 'C', '12')
                                                         .when(demographic data.income =
= 'D', '13')
                                                         .when(demographic data.income =
= 'null', None)
                                                         .otherwise(demographic data.inc
ome)
           ) \
                                   .withColumn('income', col('income').cast("int"))
# Removing non-relevant columns and duplicated entries
demographic data Q1 = demographic data.distinct().select('household id', 'household size'
, 'num_adults', 'net_worth', 'income')
```

```
from pyspark.sql.column import *
from pyspark.sql.types import StringType
from pyspark.sql.functions import *
# Removing entries that contain null values or null strings in rellevant fields/columns (
fields that are used in the queries)
# prepering helpful dataframes that will help us later in the queries
#1st condition
viewing_data_Q1_no_null = viewing_data_Q1.dropna().filter((col('prog_code') != 'null') |
                                                           (col('device id') != 'null')
                                                           (col('event date') != 'null')
).cache()
viewing_data_Q1_no_null1 = viewing data Q1.filter((col('device id') != 'null') |
                                                           (col('event date') != 'null')
) \
                                            .filter((col('device id').isNotNull()) |
                                                         (col('event date').isNotNull())
).cache()
```

```
avg_daily_events_per_device = spark.sql("""SELECT device_id, COUNT(*)/COUNT(DISTINCT even
t date) AS avg daily
                                           FROM viewing data Q1 no null1
                                           GROUP BY device_id, num_dates
avg daily events per device = avg daily events per device.join(viewing data Q1 no null, '
device id', 'inner').cache()
#2nd condition
ref data Q1 no null = ref data Q1.filter((col('device id') != 'null') |
                                         (col('dma') != 'null') ) \
                                 .filter((col('device id').isNotNull()) |
                                         (col('dma').isNotNull()) ).cache()
ref data is z in dma = ref data Q1 no null.withColumn("dma contains z", ref data Q1 no nu
ll.dma.ilike('%z%')).cache()
# 3rd condition
viewing data Q1 no null3 = viewing data Q1.filter((col('prog code') != 'null') |
                                                (col('device_id') != 'null') ) \
                                         .filter((col('prog_code').isNotNull()) |
                                                (col('device id').isNotNull()) ).cache(
viewing data Q1 no null3.createOrReplaceTempView("viewing data Q1 no null3")
ref_data_Q1_no_null = ref_data Q1.filter((col('household id') != 'null') |
                                        (col('device id') != 'null') ) \
                                 .filter((col('household id').isNotNull()) |
                                        (col('device id').isNotNull()) ).cache()
ref data Q1 no null.createOrReplaceTempView("ref data Q1 no null")
demographic data Q1 no null = demographic data Q1.filter((col('household id') != 'null')
                                                         (col('num adults') != 'null')
                                                          (col('net worth') != 'null') )
                                                 .filter((col('household id').isNotNull
())
                                                          (col('num adults').isNotNull()
) |
                                                          (col('net worth').isNotNull())
).cache()
demographic data Q1 no null.createOrReplaceTempView("demographic data Q1 no null")
# 4th condition
daily prog data Q1 no null = daily prog data Q1.filter((col('prog code') != 'null') |
                                                         (col('air date') != 'null') |
                                                         (col('air time') != 'null') ) \
                                                .filter((col('prog code').isNotNull())
                                                         (col('air date').isNotNull()) |
                                                         (col('air time').isNotNull()) )
.cache()
prog with dayofweek = daily_prog_data_Q1_no_null.withColumn('day_of_week',
                                                    dayofweek(to date(daily prog data Q
1_no_null['air_date'], 'yyyyMMdd'))).cache()
demographic_data_Q1_no_null4 = demographic_data_Q1.filter((col('household_id') != 'null')
                                                         (col('household size') != 'null
'))\
                                                .filter((col('household id').isNotNull(
) ) |
                                                         (col('household size').isNotNul
1()) ).cache()
# 5th condition
demographic data Q1 no null5 = demographic data Q1.filter((col('household id') != 'null')
                                                         (col('income') != 'null' )) \
                                                .filter((col('household id').isNotNull(
```

```
)) |
                                                         (col('income').isNotNull()) ).c
ache()
ref data Q1 no null5 = ref data Q1 no null.drop('dma')
num devices per household = ref data Q1 no null5.distinct() \
                                                 .groupBy('household id').agg(countDisti
nct('device id').alias('num devices')).cache()
ref data with count = ref_data_Q1_no_null5.join(num_devices_per_household, on='household_
id', how='inner').cache()
household data with count = ref data with count.join(demographic data Q1 no null5, 'house
hold id', 'inner').cache()
# 6th condition
daily prog data Q1 no null6 = daily prog data Q1.filter((col('prog code') != 'null') |
                                                         (col('genre') != 'null') |
                                                         (col('Duration') != 'null') ) \
                                                 .filter((col('prog code').isNotNull())
                                                         (col('genre').isNotNull()) |
                                                         (col('Duration').isNotNull()) )
.cache()
prog genres = daily prog data Q1 no null6.select('prog code', 'genre', 'Duration')\
                                .withColumn("single genres", explode(split(col("genre"),
','))).cache()
In [0]:
joined 1 = household data with count.join(ref data is z in dma, on=['household id', 'devi
ce_id'], how='inner')
joined 2 = prog with dayofweek.join(prog genres, on=['prog code', 'genre', 'Duration'],
how='inner')
joined_3 = avg_daily_events_per_device.join(joined_1, on='device id', how='inner')
joined = joined 3.join(joined 2, on='prog code', how='inner')
joined.display()
Question 1.2
Part A - in the PDF attached
Part B
In [0]:
viewing num dates.show()
In [0]:
# taking the requred programs as asked in condition 1
cond 1 = avg daily events per device.where(col('avg daily') < 5)\</pre>
                                     .select('prog code').distinct().cache()
print(cond 1.count())
9138
In [0]:
# taking the requred programs as asked in condition 2
cond_2 = ref_data_is_z_in_dma.where(ref_data_is_z_in_dma['dma_contains_z'] == True) \
                             .select('device id')\
                             .join(viewing data Q1 no null3, 'device id', 'inner')
```

.select('prog code').distinct().cache()

print(cond 2.count())

```
In [0]:
```

180061

In [0]:

```
# dataframe of the programs (prog code) that was aired between Friday at 6PM and Saturday
at 7PM
prog filtered date time = prog with dayofweek.filter(( (col('day of week') == 6) &
                                                       (daily prog data Q1 no null["air
time"] >= 180000) ) |
                                                     ( (col('day of week') == 7) &
                                                      (daily_prog_data_Q1_no_null["air_t
ime"] <= 190000) ))\</pre>
                                             .select("prog code").distinct().cache()
# combining each program with the device it was watched from
devices prog date time = prog filtered date time\
                        .join(viewing data Q1 no null3.select("prog code", "device id"),
on='prog code', how='inner').cache()
# combining each device with the households they are located at
households prog date time = devices prog date time.join(ref data Q1 no null, "device id",
"inner") \
                                                   .select('prog code', 'household id').
distinct().cache()
# taking households with size higher than or equal to 8
household above 8 = demographic data Q1 no null4.filter(demographic data Q1 no null4['hou
sehold size'] >= 8).select('household id').distinct().cache()
# taking the requred programs as asked in condition 4 by combining each household with th
e programs they watched
cond_4 = household_above_8.join(households_prog_date_time, on='household_id', how='inner
')\
                          .select("prog code").distinct().cache()
prog filtered date time.unpersist()
devices prog date time.unpersist()
households prog date time.unpersist()
household above 8.unpersist()
print(cond 4.count())
```

67646

```
from pyspark.sql.functions import *
# list of the genres that at least one of them should be contained in the program genres
genres = ['Talk', 'Politics', 'News', 'Community', 'Crime']
# creating a dataframe of these genres (from the list above)
genres df = spark.createDataFrame(genres, StringType())
# renaming the column of these genres
genres_df = genres_df.withColumnRenamed('value', 'single genres')
# taking the requred programs as asked in condition 6
cond 6 = prog genres.join(genres df, on='single genres', how='inner')\
                    .select('prog_code', 'Duration')
                    .where(col('Duration') > 35)\
                    .select('prog code').distinct() \.cache()
print(cond 6.count())
33648
In [0]:
from pyspark.sql.functions import *
# Adding a counting of conditions column and adding 1 to all entries that are true for co
ndition 1.
counter_1 = daily_prog_data_Q1.select('prog_code').distinct() \
                              .filter((col('prog code') != 'null') | (col('prog code').i
sNotNull())).cache()
counter 1 = counter 1.alias('all programs') \
                     .join(cond 1.alias('cond 1'), col("all programs.prog code") == col(
"cond 1.prog code"), 'left')\
                     .select("all_programs.*", when(col("cond_1.prog_code").isNotNull(),
1)\
                                                  .otherwise(0)
                                                  .alias("num conditions met"))\
                     .distinct().cache()
In [0]:
# A function that automates the counter of conditions column.
def conditions counter(cond, curr counter):
    counter = curr counter.alias('current programs') \
                          .join(cond.alias('cond'), col("current programs.prog code") ==
col("cond.prog_code"), 'left')\
                          .select("current programs.prog code",
                                  when(col("cond.prog code").isNotNull(), curr counter.n
um conditions met + 1) \
                                 .otherwise(curr counter.num conditions met) \
                                 .alias("num conditions met"))\
                          .distinct().cache()
    return counter
In [0]:
# Adding 1 in the counter column for each entry that is true for any conditions.
counter 2 = conditions counter(cond 2, counter 1)
counter 3 = conditions counter(cond 3, counter 2)
counter 4 = conditions counter(cond 4, counter 3)
counter 5 = conditions counter(cond 5, counter 4)
counter 6 = conditions counter(cond 6, counter 5)
counter 1.unpersist()
counter 2.unpersist()
counter 3.unpersist()
counter 4.unpersist()
```

print(cond 5.count())

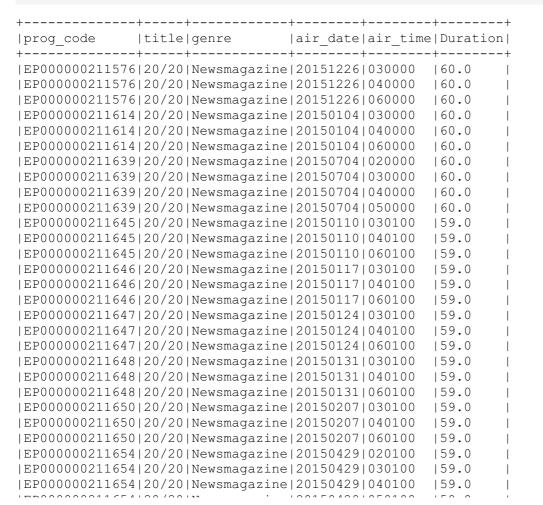
255489

In [0]:

```
malicious.count()
```

Out[30]: 71765

prog_code	title	genre	air_date	air_time	Duration
EP000000211576	20/20	Newsmagazine	20151226	030000	60.0
EP000000211576	20/20	Newsmagazine	20151226	040000	60.0
EP000000211576	20/20	Newsmagazine	20151226	060000	60.0
EP000000211614	20/20	Newsmagazine	20150104	030000	60.0
EP000000211614	20/20	Newsmagazine	20150104	040000	60.0
EP000000211614	20/20	Newsmagazine	20150104	060000	60.0
EP000000211639	20/20	Newsmagazine	20150704	020000	60.0
EP000000211639	20/20	Newsmagazine	20150704	030000	60.0
EP000000211639	20/20	Newsmagazine	20150704	040000	60.0
EP000000211639	20/20	Newsmagazine	20150704	050000	60.0



```
|EPUUUUUUZ11654|ZU/ZU|Newsmagaz1ne|ZU15U4Z9|U5U1UU
|EP000000211659|20/20|Newsmagazine|20150411|020100
|EP000000211659|20/20|Newsmagazine|20150411|030100
|EP000000211659|20/20|Newsmagazine|20150411|040100
                                                  159.0
|EP000000211659|20/20|Newsmagazine|20150411|050100
                                                  159.0
|EP000000211661|20/20|Newsmagazine|20150418|020100
                                                   |59.0
|EP000000211661|20/20|Newsmagazine|20150418|030100
                                                   |59.0
|EP000000211661|20/20|Newsmagazine|20150418|040100
                                                    |59.0
|EP000000211661|20/20|Newsmagazine|20150418|050100
                                                    |59.0
|EP000000211665|20/20|Newsmagazine|20150502|020100
|EP000000211665|20/20|Newsmagazine|20150502|030100
|EP000000211665|20/20|Newsmagazine|20150502|040100
                                                    |59.0
|EP000000211665|20/20|Newsmagazine|20150502|050100
                                                    |59.0
|EP000000211666|20/20|Newsmagazine|20150509|020100
                                                    159.0
|EP000000211666|20/20|Newsmagazine|20150509|030100 |59.0
|EP000000211666|20/20|Newsmagazine|20150509|040100 |59.0
|EP000000211666|20/20|Newsmagazine|20150509|050100 |59.0
|EP000000211667|20/20|Newsmagazine|20150516|020100 |59.0
|EP000000211667|20/20|Newsmagazine|20150516|030100 |59.0
|EP000000211667|20/20|Newsmagazine|20150516|040100 |59.0
|EP000000211667|20/20|Newsmagazine|20150516|050100 |59.0
|EP000000211669|20/20|Newsmagazine|20150530|020000 |60.0
|EP000000211669|20/20|Newsmagazine|20150530|030000 |60.0
|EP000000211669|20/20|Newsmagazine|20150530|040000 |60.0
                                                  |60.0
|EP000000211669|20/20|Newsmagazine|20150530|050000
                                                  |59.0
|EP000000211670|20/20|Newsmagazine|20150610|050100
                                                   |60.0
|EP000000211670|20/20|Newsmagazine|20150606|020000
|EP000000211670|20/20|Newsmagazine|20150606|030000
                                                   |60.0
|EP000000211670|20/20|Newsmagazine|20150606|040000
|EP000000211670|20/20|Newsmagazine|20150606|050000
|EP000000211672|20/20|Newsmagazine|20150620|020000
                                                    |60.0
|EP000000211672|20/20|Newsmagazine|20150620|030000
                                                    160.0
|EP000000211672|20/20|Newsmagazine|20150620|040000
                                                   |60.0
|EP000000211672|20/20|Newsmagazine|20150620|050000
                                                  160.0
|EP000000211675|20/20|Newsmagazine|20150718|020000 |60.0
|EP000000211675|20/20|Newsmagazine|20150718|030000 |60.0
|EP000000211675|20/20|Newsmagazine|20150718|040000 |60.0
|EP000000211675|20/20|Newsmagazine|20150718|050000 |60.0
|EP000000211676|20/20|Newsmagazine|20150725|020000 |60.0
|EP000000211676|20/20|Newsmagazine|20150725|030000
|EP000000211676|20/20|Newsmagazine|20150725|040000
                                                  160.0
|EP000000211676|20/20|Newsmagazine|20150725|050000
                                                  160.0
|EP000000211679|20/20|Newsmagazine|20150815|020000
                                                  |60.0
|EP000000211679|20/20|Newsmagazine|20150815|023000
                                                   |30.0
|EP000000211679|20/20|Newsmagazine|20150815|040000
                                                    |60.0
|EP000000211679|20/20|Newsmagazine|20150815|050000
                                                    |60.0
|EP000000211679|20/20|Newsmagazine|20150815|050700
|EP000000211679|20/20|Newsmagazine|20150815|063600
|EP000000211679|20/20|Newsmagazine|20150815|080000
                                                    |60.0
|EP000000211679|20/20|Newsmagazine|20150815|020000
                                                    160.0
|EP000000211679|20/20|Newsmagazine|20150815|023000
                                                   130.0
|EP000000211679|20/20|Newsmagazine|20150815|040000 |60.0
|EP000000211679|20/20|Newsmagazine|20150815|050000 |60.0
|EP000000211679|20/20|Newsmagazine|20150815|050700 |60.0
|EP000000211679|20/20|Newsmagazine|20150815|063600 |60.0
|EP000000211679|20/20|Newsmagazine|20150815|080000 |60.0
|EP000000211679|20/20|Newsmagazine|20150816|210000 |60.0
|EP000000211680|20/20|Newsmagazine|20150823|053500
|EP000000211680|20/20|Newsmagazine|20150822|020000 |60.0
|EP000000211680|20/20|Newsmagazine|20150822|030000 |60.0
                                                  |60.0
|EP000000211680|20/20|Newsmagazine|20150822|040000
                                                  |60.0
|EP000000211680|20/20|Newsmagazine|20150822|050000
                                                   160.0
|EP000000211680|20/20|Newsmagazine|20150822|050700
                                                   |59.0
|EP000000211680|20/20|Newsmagazine|20150822|080100
|EP000000211680|20/20|Newsmagazine|20150822|020000
|EP000000211680|20/20|Newsmagazine|20150822|040000
|EP000000211680|20/20|Newsmagazine|20150822|050000
|EP000000211680|20/20|Newsmagazine|20150822|050700
                                                    |60.0
|EP000000211680|20/20|Newsmagazine|20150822|080100
                                                    159.0
|EP000000211680|20/20|Newsmagazine|20150823|053500
                                                    |60.0
|EP000000211681|20/20|Newsmagazine|20150829|020000
                                                    |60.0
|EP000000211681|20/20|Newsmagazine|20150829|040000
                                                    160.0
```

			20120829 020000	60.0
			20150829 050700	60.0
EP000000211681	20/20	Newsmagazine	20150829 063600	60.0
EP000000211681	20/20	Newsmagazine	20150829 080000	60.0
EP000000211681	20/20	Newsmagazine	20150829 090100	59.0
EP000000211681	20/20	Newsmagazine	20150830 053500	60.0
EP000000211681	20/20	Newsmagazine	20150830 053500	60.0
EP000000211681	20/20	Newsmagazine	20150829 020000	60.0
EP000000211681	20/20	Newsmagazine	20150829 030000	60.0
EP000000211681	20/20	Newsmagazine	20150829 040000	60.0
EP000000211681	20/20	Newsmagazine	20150829 050000	60.0
EP000000211681	20/20	Newsmagazine	20150829 050700	60.0
EP000000211681	20/20	Newsmagazine	20150829 063600	60.0
EP000000211681	20/20	Newsmagazine	20150829 080000	60.0
EP000000211681	20/20	Newsmagazine	20150829 090100	59.0
EP000000211682	20/20	Newsmagazine	20150905 020000	60.0
			20150905 030000	60.0
EP000000211682	20/20	Newsmagazine	20150905 040000	60.0
EP000000211682	20/20	Newsmagazine	20150905 050000	60.0
EP000000211683	20/20	Newsmagazine	20150912 020100	59.0
EP000000211683	20/20	Newsmagazine	20150912 030100	59.0
EP000000211683	20/20	Newsmagazine	20150912 040100	59.0
EP000000211683	20/20	Newsmagazine	20150912 050100	59.0
EP000000211684	20/20	Newsmagazine	20150919 020100	59.0
EP000000211684	20/20	Newsmagazine	20150919 020200	59.0
EP000000211684	20/20	Newsmagazine	20150919 030100	59.0
			20150919 040100	59.0
			20150919 050100	59.0
EP000000211685	20/20	Newsmagazine	20150926 020100	59.0
		-	20150926 030100	59.0
EP000000211685	20/20	Newsmagazine	20150926 040100	59.0
		-	20150926 050100	59.0
			20151003 020100	59.0
			20151003 030100	59.0
		-	20151003 040100	59.0
			20151003 050100	59.0
			20151017 020100	59.0
			20151017 030100	59.0
•			20151017 040100	59.0
			20151017 050100	59.0
			20151024 020100	59.0
			20151024 030100	59.0
		-	20151024 040100	59.0
			20151024 050100	59.0
			20151031 020100	59.0
			20151031 030100	59.0
			20151031 040100	59.0
			20151031 050100	59.0
		-	20151107 030100	59.0
		_	20151107 040100	59.0
+	+	+	+	-++

only showing top 150 rows