

```
spark = SparkSession.builder.appName("my_pro
```

```
viewing10m_df = spark.read.format("csv")  
    .option("header", "true") \  
    .option("delimiter", ",") \  
    .schema(schemas_dict['viewing_f']) \  
    .load(dataPath)
```

mso_code	device_id	event_date	event_time	station_num	prog_code
01540	0000000050f3	20150222	193802	61812	EP009279780033
01540	0000000050f3	20150222	195314	31709	EP021056430002
01540	0000000050f3	20150222	200151	61812	EP009279780033
01540	000000005518	20150222	111139	46784	EP004891370013
01540	000000005518	20150222	190000	14771	EP012124070127
01540	000000005518	20150222	200000	14771	EP010237320166

viewing10m\_df contains 9935852 rows!

## Read reference data

Note that we removed the 'System Type' column.

```
In [0]: # Read the new parquet
ref_data_schema = StructType([
    StructField('device_id', StringType()),
    StructField('dma', StringType()),
    StructField('dma_code', StringType()),
    StructField('household_id', IntegerType()),
    StructField('zipcode', IntegerType())
])

# Reading as a Parquet
dataPath = f"dbfs:/FileStore/ddm/ref_data"
ref_data = spark.read.format('parquet') \
    .option("inferSchema", "true") \
    .load(dataPath)

display(ref_data.limit(6))
print(f'ref_data contains {ref_data.count()} rows!')
```

device_id	dma	dma_code	household_id	zipcode
0000000050f3	Toledo	547	1471346	43609
000000006785	Amarillo	634	1924512	79119
000000007320	Lake Charles	643	3154808	70634
000000007df9	Lake Charles	643	1924566	70601
000000009595	Lexington	541	1600886	40601
000000009c6a	Houston	618	1924713	77339

ref\_data contains 704172 rows!

## Part 1

### 1.1

```
In [0]: from pyspark.sql.functions import col, lower, to_date, dayofmonth, date_format, avg, count, sum as spark_sum, abs, lit, coalesce
from functools import reduce

# 1.1 Extract and add features - conditions 1, 4, 6, 7

daily_prog_df = daily_prog_df.withColumn("air_date_parsed", to_date(col("air_date"), "yyyyMMdd")) \
    .withColumn("day", dayofmonth(col("air_date_parsed"))) \
    .withColumn("weekday", date_format(col("air_date_parsed"), "E"))

# 1st Condition: Duration > avg duration
avg_duration = daily_prog_df.select(avg("Duration")).first()[0]
prog_data_cond = daily_prog_df.withColumn("cond_1", col("Duration") > avg_duration)

# 6th Condition: Genre contains specific values
genre_keywords = ['Collectibles', 'Art', 'Snowmobile', 'Public affairs', 'Animated', 'Music']
```

```
prog_data_cond = prog_data_cond.withCol

# 7th Condition: Title contains at least one of the words in word_list
word_list = ['better', 'girls', 'the', 'she', 'he', 'her', 'his']
prog_data_cond = prog_data_cond.withCol(
    "words_array",
    lambda row: any(word in row["Title"] for word in word_list))
```

```
"cond"
coale
)
display(p
```

EP000000250063	21 Jump Street	Crime drama	20151219	180000	60.0	2015-12-19	19	Sat	false	false	0	false	false	false	false
EP000003690043	Are You Being Served?	Sitcom	20151219	123000	30.0	2015-12-19	19	Sat	false	false	0	false	false	false	false
EP000003240027	The Andy Griffith Show	Sitcom	20151219	230000	30.0	2015-12-19	19	Sat	false	false	1	false	false	false	false
EP000000510180	A Different World	Sitcom	20151219	103000	30.0	2015-12-19	19	Sat	false	false	0	false	false	false	false
EP000000510180	A Different World	Sitcom	20151219	133000	30.0	2015-12-19	19	Sat	false	false	0	false	false	false	false
EP000003240109	The Andy Griffith Show	Sitcom	20151220	023000	30.0	2015-12-20	20	Sun	false	false	1	false	false	false	false
EP000002040083	All in the Family	Sitcom	20151219	213000	30.0	2015-12-19	19	Sat	false	false	1	false	false	false	false
EP000003690041	Are You Being Served?	Sitcom	20151219	113000	30.0	2015-12-19	19	Sat	false	false	0	false	false	false	false

  

```
In [0]: # 5th Condition : Household with more then 3 devices and income less than average
# Compute average income
avg_income = demo_df.select(avg(col("income"))).first()[0]

# Count number of devices per household
device_count_df = ref_data.groupBy("household_id").agg(
    countDistinct("device_id").alias("num_devices")
)
# Join income and device info
household_info_df = demo_df.join(device_count_df, on="household_id", how="inner")

# Keep households with > 3 devices and income < average
qualified_households_df = household_info_df.filter(
    (col("num_devices") > 3) & (col("income") < avg_income)
).select("household_id").distinct()

# Get device_ids from these households
qualified_devices_df = ref_data.join(
    qualified_households_df, on="household_id", how="inner"
).select("device_id").distinct()

# Get program codes watched on these devices
qualified_prog_codes_df = viewing10m_df.join(
    qualified_devices_df, on="device_id", how="inner"
).select("prog_code").distinct().withColumn("cond_5", lit(True))

# Join to main DataFrame and set cond_5 flag
prog_data_cond = prog_data_cond.join(
    qualified_prog_codes_df,
    on="prog_code",
    how="left"
).withColumn(
    "cond_5",
    coalesce(col("cond_5"), lit(False))
)
```

## 1.2

  

```
In [0]: # 1.2

condition_cols = [f"cond_{i}" for i in range(1, 8)]

# Calculate the number of conditions met per row
prog_data_cond = prog_data_cond.withColumn(
    "num_conditions_met",
    reduce(lambda a, b: a + b, [col(c).cast("int") for c in condition_cols])
)

# Flag malicious programs
prog_data_cond = prog_data_cond.withColumn(
    "is_malicious",
    when(col("num_conditions_met") >= 4, True).otherwise(False)
)

# Group by title and calculate malicious percentage
malicious_summary_df = prog_data_cond.groupBy("title").agg(
    (spark_sum(col("is_malicious").cast("int")) / count("*")).alias("malicious_percentage"),
    count("*").alias("total_records")
)
```

```
# Fi  
top_
```

title	malicious_percentage	total_records
Philomena	1.0	658
Sabata	1.0	39
Weekends With Alex Witt	1.0	160
Lone Rider	1.0	96
Fox34 Weather Nation	1.0	2299
Today's Country: NASH	1.0	701

