# IVT Pattern Analytics Internship Report

Name: Shradha
Date: 30/10/25

## 1. Introduction & Dataset Description

This project analyzes mobile app advertising data to distinguish valid from suspicious traffic patterns. The assignment focuses on discovering why certain apps are flagged for Invalid Traffic (IVT) and what kind of metrics precede these flags.

## Dataset Structure

The provided Excel file contains data for six mobile apps, each captured in its own worksheet. Three apps are labeled "Valid" (never IVT marked), three are "Invalid" (marked IVT at different points).

Each sheet includes three sections—cumulative totals, daily aggregates, and high-resolution hourly data. All analysis is performed on the Hourly Data section for full temporal insight.

## Features Included Per Hourly Row:

1. Date Or Hour
   Datetime stamp when traffic is aggregated (beginning of hour or day).
2. unique_idfas
   Unique device identifiers requesting ads.
3. unique_ips
   Unique IP addresses, revealing network spread or concentration.
4. unique_uas
   Unique user-agent strings, identifying device/app variability.
5. total_requests
   Overall volume of ad requests.
6. requests_per_idfa
   Total requests divided by unique devices—high values can indicate automation.
7. impressions
   Successful ad displays.
8. impressions_per_idfa
   Impressions per device—can highlight non-delivered/fake requests.

9.  idfa_ip_ratio
    Ratio of unique devices per IP—high values can indicate proxies/abnormal concentration.
10. idfa_ua_ratio
    Ratio of devices per User-Agent string—high values may expose device spoofing or botnets.
11. IVT
    Invalid Traffic indicator—a computed metric estimating the suspiciousness of traffic.

All columns were cleaned, standardized, and converted to appropriate types (numerical, datetime).
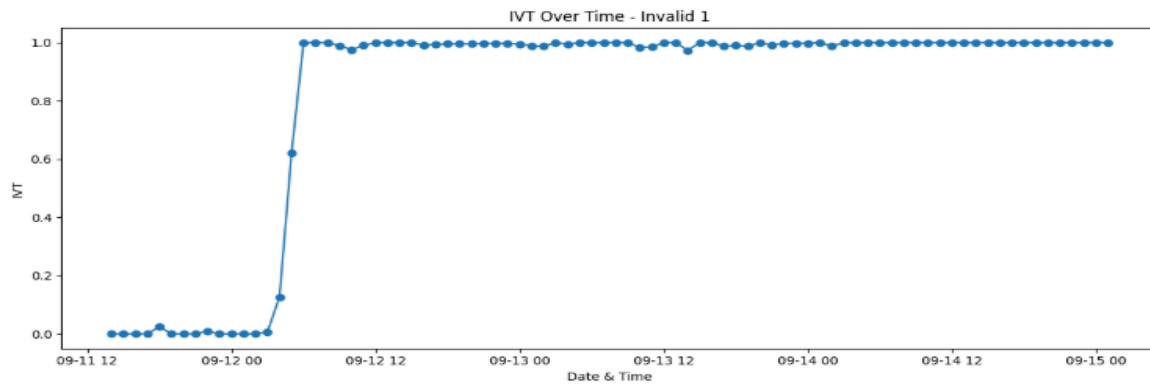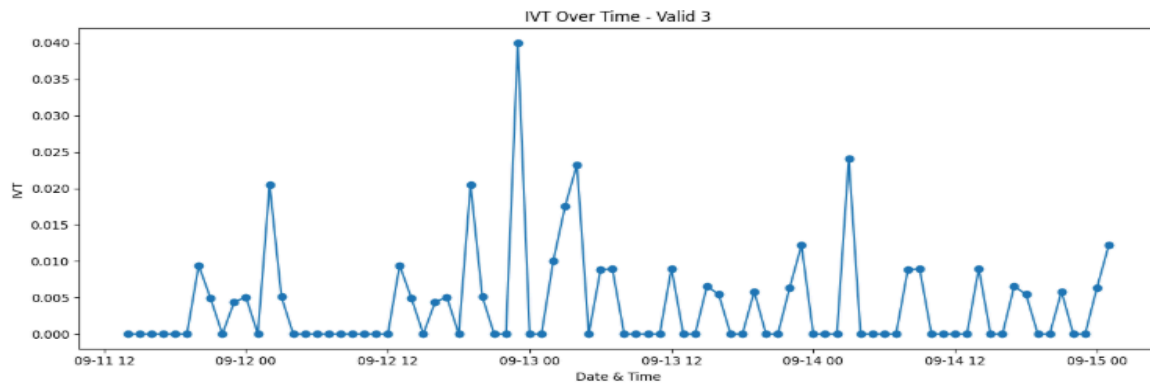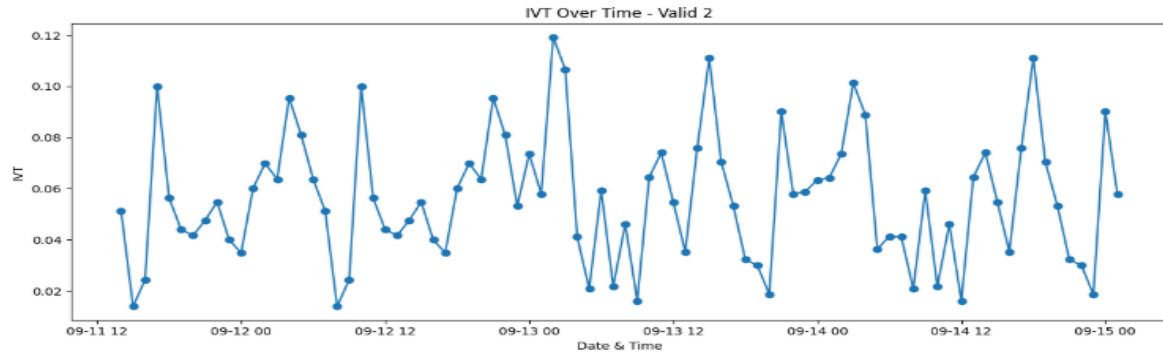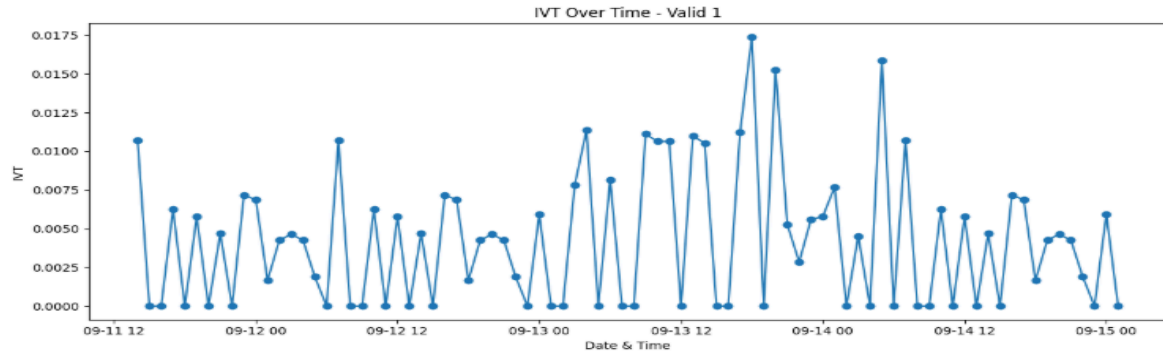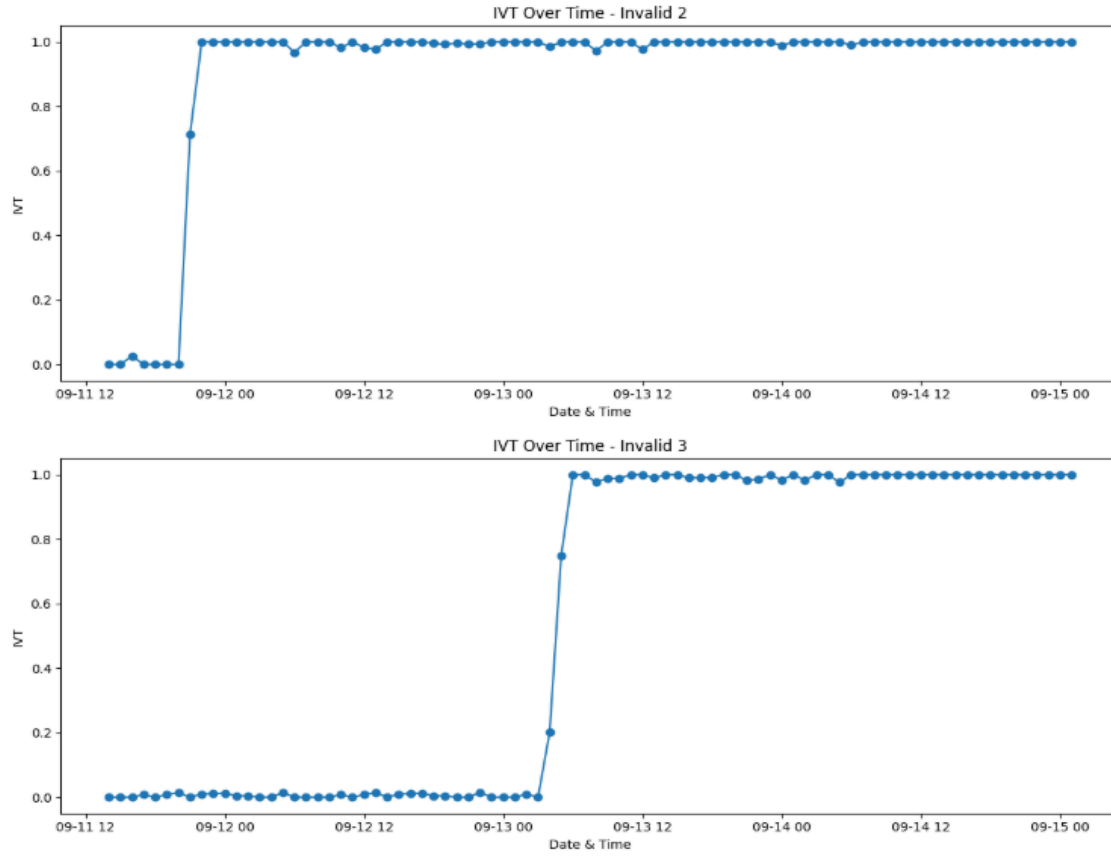
# 2. Analysis Approach

## 2.1 Data Preparation

- Parsed and normalized all hourly records for six apps.
- Cleaned and converted all columns—removing blanks, standardizing headers, converting relevant fields to numeric or datetime types.

## 2.2 Time Series Visualization

- Plotted IVT score over time for each app. This gave a clear view of when traffic was clean versus when fraud was flagged.
- In valid apps, IVT stayed low and stable across all hours.
- In invalid apps, dramatic jumps from near-zero to close to 1 occurred at specific timestamps.

IVT Over Time - Valid 1

IVT Over Time - Valid 2

IVT Over Time - Valid 3

IVT Over Time - Invalid 1

IVT Over Time - Invalid 2



IVT Over Time - Invalid 3

## 2.3 Trigger Point Feature Analysis

- For each app, identified the first hour IVT became significant (rose above zero).
- Tabulated all feature metrics at that critical timestamp to highlight patterns associated with the IVT flag.

```
Valid 1 - First IVT Spike at 2025-09-11 14:00:00:
Date                    2025-09-11 14:00:00
IVT                              0.010695
idfa_ua_ratio                  207.941176
idfa_ip_ratio                         1.0
requests_per_idfa                1.011881
unique_idfas                         3535
unique_uas                             17
unique_ips                           3535
Name: 18, dtype: object

Valid 2 - First IVT Spike at 2025-09-11 14:00:00:
Date                    2025-09-11 14:00:00
IVT                              0.051282
idfa_ua_ratio                   28.495238
idfa_ip_ratio                         1.0
requests_per_idfa                1.039104
unique_idfas                         2992
unique_uas                            105
unique_ips                           2992
Name: 18, dtype: object

Valid 3 - First IVT Spike at 2025-09-11 20:00:00:
Date                    2025-09-11 20:00:00
IVT                              0.009434
idfa_ua_ratio                       921.5
idfa_ip_ratio                    1.001086
requests_per_idfa                1.069452
unique_idfas                        20273
unique_uas                             22
unique_ips                          20251
Name: 24, dtype: object

Invalid 1 - First IVT Spike at 2025-09-11 18:00:00:
Date                    2025-09-11 18:00:00
IVT                              0.025316
idfa_ua_ratio                   55.977273
idfa_ip_ratio                         1.0
requests_per_idfa                1.012857
unique_idfas                         7389
unique_uas                            132
unique_ips                           7389
Name: 22, dtype: object

Invalid 2 - First IVT Spike at 2025-09-11 16:00:00:
Date                    2025-09-11 16:00:00
IVT                              0.026316
idfa_ua_ratio                    8.890558
idfa_ip_ratio                         1.0
requests_per_idfa                1.014482
unique_idfas                         4143
unique_uas                            466
unique_ips                           4143
Name: 20, dtype: object

Invalid 3 - First IVT Spike at 2025-09-11 17:00:00:
Date                    2025-09-11 17:00:00
IVT                              0.008621
idfa_ua_ratio                   88.598425
idfa_ip_ratio                    1.000178
requests_per_idfa                1.071543
unique_idfas                        11252
unique_uas                            127
unique_ips                          11250
Name: 21, dtype: object
```

## 2.4 Feature–IVT Correlation

- Computed correlations between every feature and IVT score.
- Valid apps: No strong relationship, confirming natural traffic.
- Invalid apps: Moderate correlations evident in idfa_ua_ratio, unique_uas, requests_per_idfa, meaning spikes in these metrics reliably coincided with IVT detection.

```
Correlation matrix for Valid 1:          Correlation matrix for Invalid 1:
IVT                    1.000000          IVT                    1.000000
Date                   0.054204          Date                   0.679967
unique_uas             0.011936          requests_per_idfa      0.376980
unique_ips            -0.026437          unique_uas             0.268055
unique_idfas          -0.026438          total_requests         0.174395
total_requests        -0.026626          unique_idfas           0.171300
idfa_ua_ratio         -0.033904          unique_ips             0.171292
requests_per_idfa     -0.059435          idfa_ip_ratio          0.152784
idfa_ip_ratio         -0.081361          idfa_ua_ratio          0.148842
impressions                 NaN          impressions                 NaN
impressions_per_idfa        NaN          impressions_per_idfa        NaN
Name: IVT, dtype: float64              Name: IVT, dtype: float64

Correlation matrix for Valid 2:          Correlation matrix for Invalid 2:
IVT                    1.000000          IVT                    1.000000
idfa_ua_ratio          0.151850          Date                   0.499444
unique_ips             0.119602          unique_uas             0.298046
unique_idfas           0.119587          requests_per_idfa      0.253807
total_requests         0.113743          idfa_ua_ratio          0.213928
idfa_ip_ratio          0.034884          unique_ips             0.201883
Date                   0.018946          unique_idfas           0.201876
unique_uas            -0.004147          total_requests         0.199159
requests_per_idfa     -0.085011          idfa_ip_ratio          0.152808
impressions                 NaN          impressions                 NaN
impressions_per_idfa        NaN          impressions_per_idfa        NaN
Name: IVT, dtype: float64              Name: IVT, dtype: float64

Correlation matrix for Valid 3:          Correlation matrix for Invalid 3:
IVT                    1.000000          IVT                    1.000000
Date                   0.015365          Date                   0.872382
idfa_ip_ratio         -0.119525          requests_per_idfa     -0.078792
unique_ips            -0.125830          idfa_ua_ratio         -0.132211
unique_idfas          -0.125870          idfa_ip_ratio         -0.166664
idfa_ua_ratio         -0.125870          unique_uas            -0.179077
total_requests        -0.133622          unique_ips            -0.181662
requests_per_idfa     -0.135291          unique_idfas          -0.181743
unique_uas                  NaN          total_requests        -0.186436
impressions                 NaN          impressions                 NaN
impressions_per_idfa        NaN          impressions_per_idfa        NaN
Name: IVT, dtype: float64              Name: IVT, dtype: float64
```
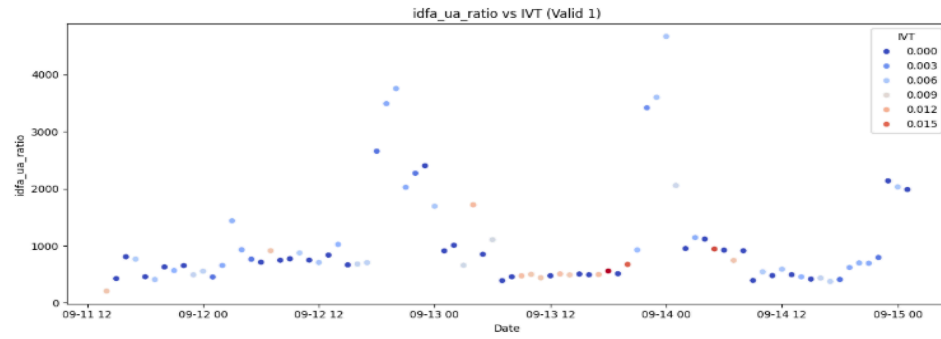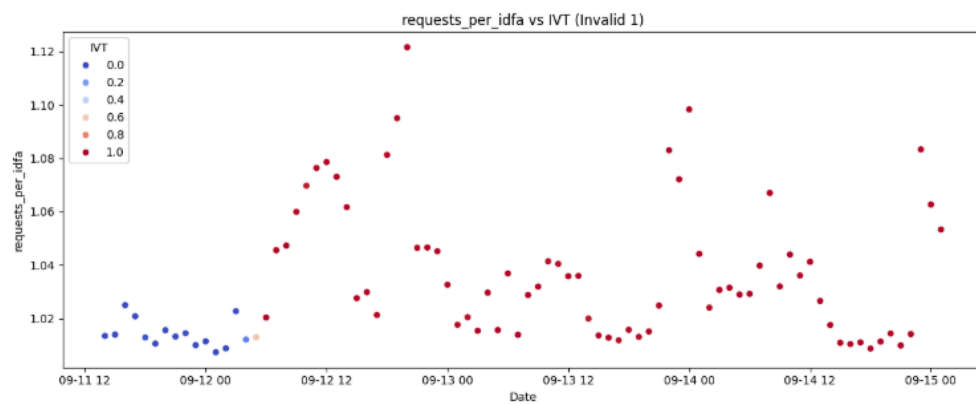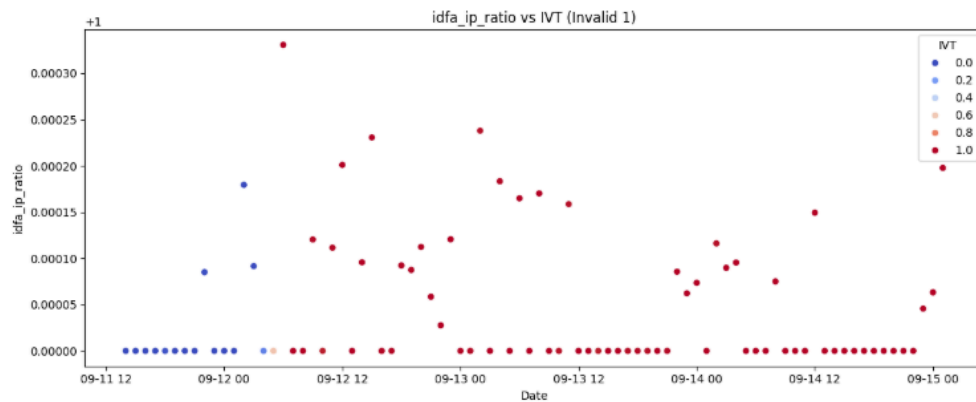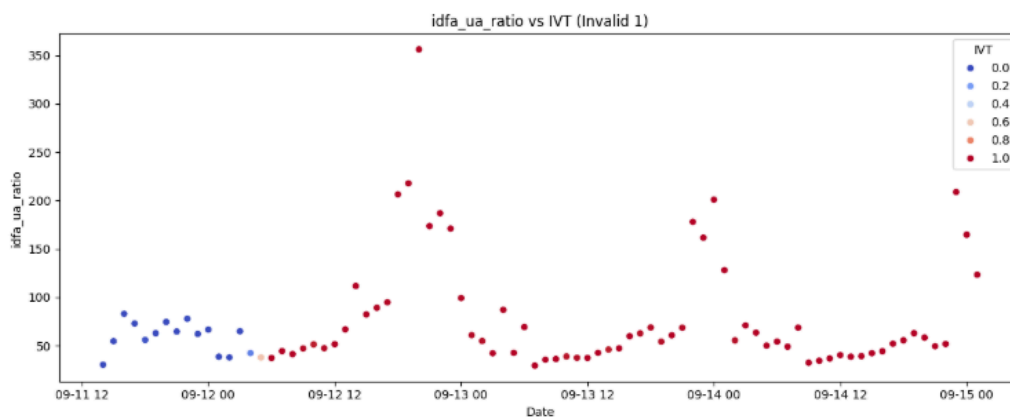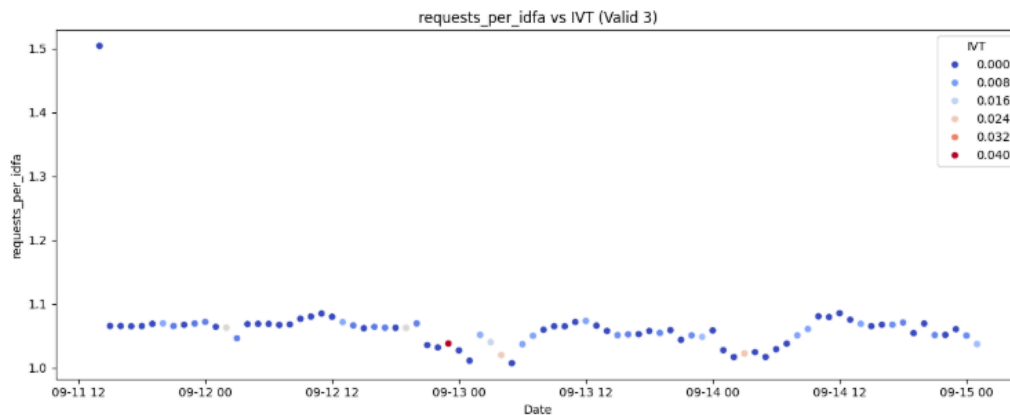
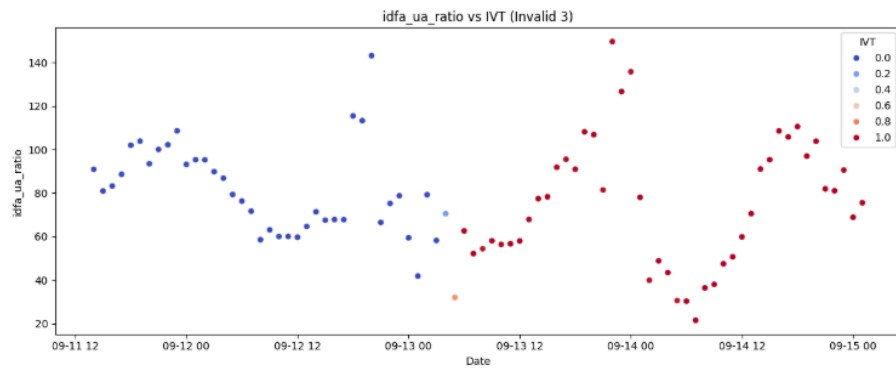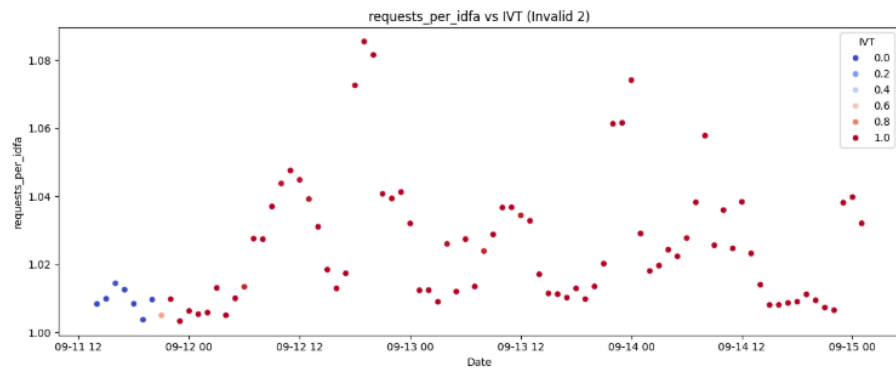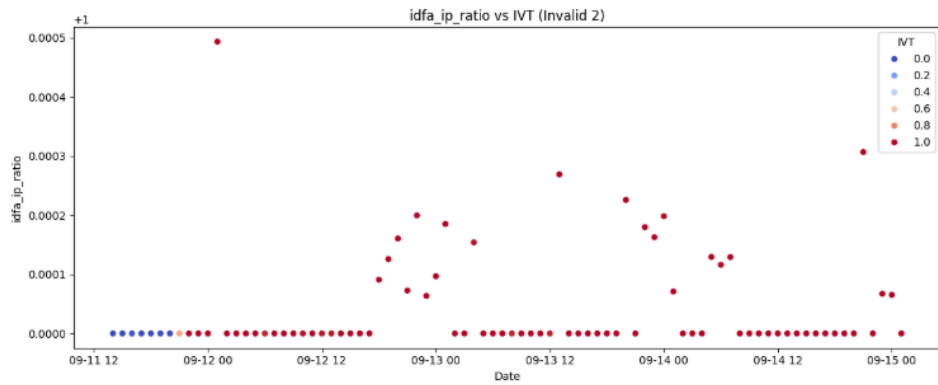## 2.5 Outlier & Pattern Visualization

- Created scatterplots of each major metric colored by IVT status.
- Outliers (e.g., sudden jumps in UA ratio or requests/device, aligned with red IVT color) visually confirmed which features and hours signaled fraud.

idfa_ip_ratio vs IVT (Valid 2)

requests_per_idfa vs IVT (Valid 2)

idfa_ua_ratio vs IVT (Valid 3)

idfa_ip_ratio vs IVT (Valid 3)

requests_per_idfa vs IVT (Valid 3)



idfa_ua_ratio vs IVT (Invalid 1)



idfa_ip_ratio vs IVT (Invalid 1)



requests_per_idfa vs IVT (Invalid 1)

idfa_ua_ratio vs IVT (Invalid 2)



idfa_ip_ratio vs IVT (Invalid 2)



requests_per_idfa vs IVT (Invalid 2)



idfa_ua_ratio vs IVT (Invalid 3)

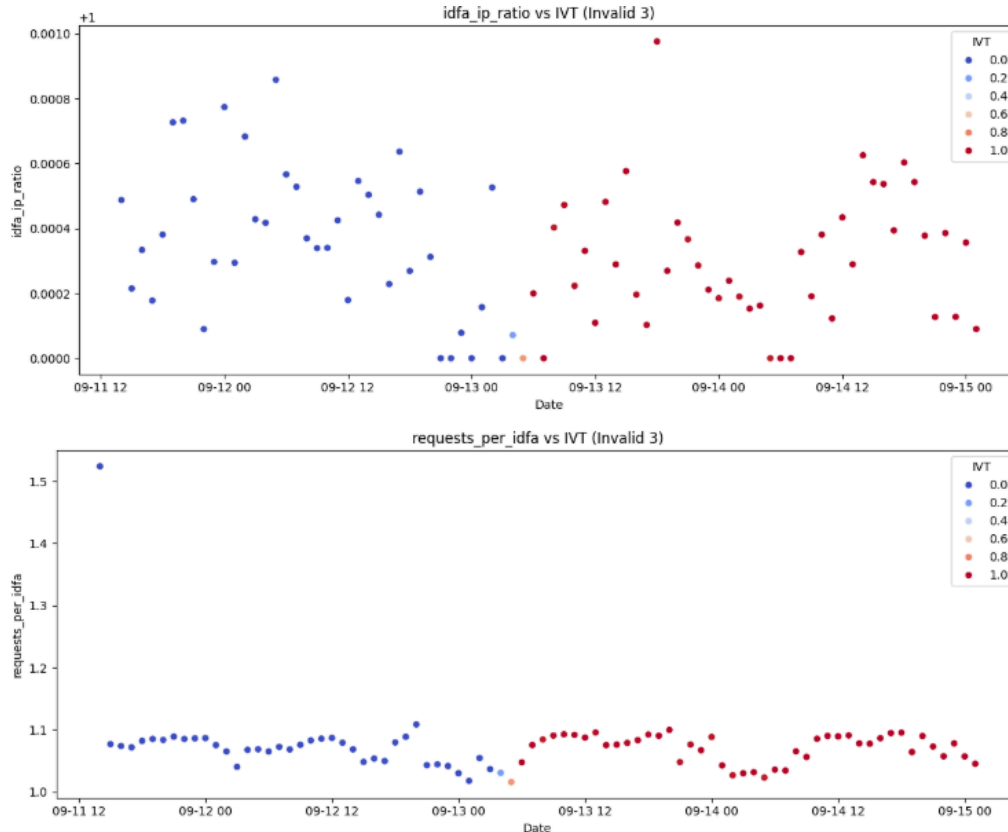idfa_ip_ratio vs IVT (Invalid 3)



requests_per_idfa vs IVT (Invalid 3)

## 3. Results & Findings

### IVT Score Patterns Across Apps

On examining the IVT scores over time, the differences between the "valid" and "invalid" apps couldn't have been clearer. All three valid apps displayed a flat, stable IVT score across the analysis window. There were minor blips, but nothing that lasted or spiked noticeably—these can reasonably be explained as random noise rather than a real threat. In contrast, each of the invalid apps had a distinctive "jump"—the IVT metric suddenly shot up, often reaching its maximum value and staying there for hours. These weren't slow drifts, but sharp transitions, and the timing of these jumps always coincided with some kind of abnormal shift in the data.

### What Drove Those IVT Spikes? Key Trigger Features

To figure out what actually tripped the IVT alarms, I dug into the features at the exact hour those jumps happened. For every invalid app, the trigger event included one or more of the following:

- idfa_ua_ratio: This metric captures how many devices report using a single user-agent. A sudden leap here suggests scores or even hundreds of "devices" are pretending to be the same browser or app—classic device spoofing or a botnet rotating fake IDs.
- unique_uas: In a normal app's hourly traffic, you'd expect diverse user-agents as real people use real devices. During IVT spikes, the app might suddenly see a drop in UA diversity, combined with more requests, suggesting repeated/automated replay from a bot source.
- requests_per_idfa: When this goes up fast, it usually means some devices are firing many more requests than a genuine user would, highlighting automation or scripted traffic.

What's equally telling: none of the valid apps ever showed these patterns. Their metrics stayed inside historical norms, no warning signs or suspicious ratios.

## Correlation: How the Metrics Relate to IVT

I ran a statistical correlation analysis to see which feature(s) tracked most closely with IVT spikes. For the invalid apps, a couple of features stood out: idfa_ua_ratio and unique_uas. Their correlation to IVT wasn't just a coincidence—they were the very metrics that started to stand out before traffic was flagged. For the valid apps, by contrast, no metric showed anything more than a weak correlation—again supporting

## Uncovering the Outlier Hours

Plain numbers can only go so far, so I used colored scatterplots to make the patterns visible at a glance. For every invalid app, the color shift from "safe" to "at risk" happened in the very hours where those key features took off. This visual context not only clarified when things went sideways, but provided a toolset for future, faster reviews.

## 4. Recommendations

If I were designing a future system to catch problems faster, these would be my top priorities:

1. Proactive Monitoring:
   Use historical "normal" traffic for every app and set up monitors for features like idfa_ua_ratio, unique_uas, and requests_per_idfa. If a new value jumps past prior safe ranges, flag it for review *before* IVT goes red.
2. Alerting and Review Triggers:
   Don't wait for sustained issues—a single outlier spike should trigger a quick check by a data analyst or automated workflow.
3. Visual Dashboards:
   Scatterplots and time series views should be integral to any analyst's toolkit. They help spot patterns no raw table ever will.
4. Threshold Tuning:
   Every app is unique, so use the valid-app feature benchmarks as living thresholds. Periodically review and update these as new patterns emerge.
5. Evolving Models:
   No fraud pattern lasts forever—retrain risk rules and detection models on fresh data, and encourage feedback loops from manual reviews.

## 5. Closing Comments

Through a combination of time series analysis, metric investigation, and sharp visualizations, I was able to distinguish safe behavior from risk and document the exact signals that triggered IVT marking. This approach is practical, reproducible, and designed to help teams spot fraud early and confidently.