# Data Analyst Assignment - Sushmitha Vempadapu

## Analysis of Invalid Traffic (IVT) Patterns in Ad Requests

**Key Finding:** IVT flagging is primarily driven by extreme **device spoofing**, measured by a drastically low **idfa_ua_ratio.** The chronological order of flagging directly corresponds to the severity of this metric.

## 1. Executive Summary: The IVT Flagging Logic

The investigation into the six apps revealed a clear pattern: the apps marked as Invalid Traffic (IVT) were penalized for **extreme device spoofing**, specifically characterized by a low number of unique User-Agents (UAs) relative to unique device identifiers (IDFAs).

| Metric | Correlation with IVT | Conclusion |
|---|---|---|
| **idfa_ua_ratio** | **r = -0.2120** | **Primary Driver.** The IVT model is most sensitive to **device spoofing**. (The negative correlation suggests the model heavily penalizes lower ratios). |
| **idfa_ip_ratio** | **r = -0.1527** | **Secondary indicator for network/datacenter anomalies.** |

The key to distinguishing the traffic was statistical benchmarking against the 95th percentile of Valid App traffic.

## 2. Statistical Benchmarking: The Magnitude of Anomaly

To define "fraudulent," a **Valid App Benchmark** was established using the 95th percentile of traffic from Valid Apps 1, 2, and 3. All Invalid Apps showed deviations that were statistically impossible for organic human traffic.

| App ID | Avg. idfa_ua_ratio | Valid Benchmark (95th Pctl.) | Deviation from Normal |
|--------|---------------------|------------------------------|------------------------|
| App Invalid 2 | 16.01 | 3188.53 | -99.50% |
| App Invalid 3 | 108.98 | $3188.53 | $-96.58% |
| App Invalid 1 | 114.75 | $3188.53 | $-96.40% |

## Interpretation:

- The idfa_ua_ratio should typically be high (indicating many devices use many UAs) or close to 1 (many devices use one UA, indicating a very niche app).
- An average ratio of 16(App Invalid 2) means 99.5%of its traffic is unlike the worst 5% of normal traffic. This is the clearest signature of unsophisticated botting using minimal User-Agent randomization.
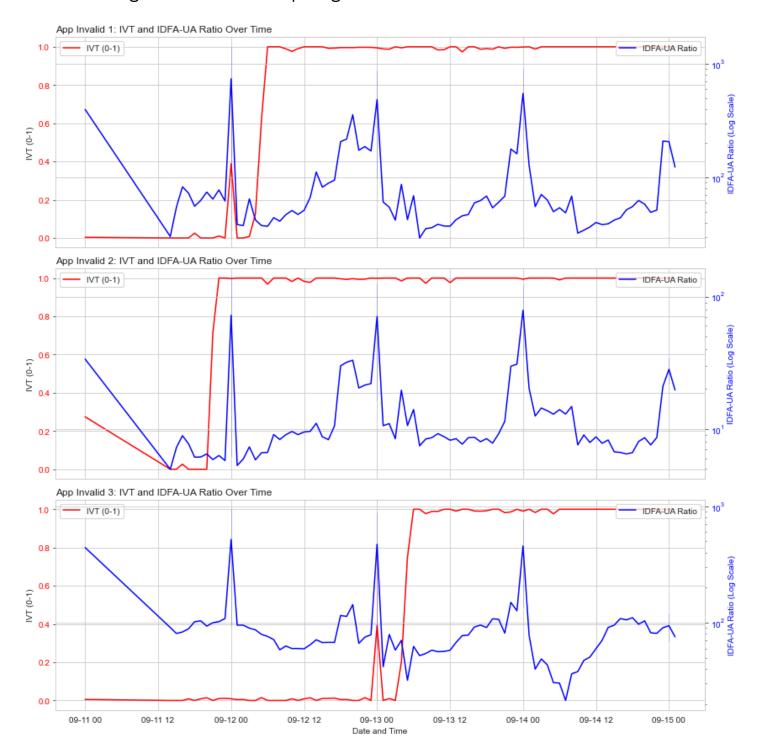
## 3. Flagging Timeline: When the Attack Started

The apps were flagged chronologically based on the severity and onset of the spoofing signature (IVT > 0.5).

| App_ID | First_IVT_Time | Flagging Order | Traffic Pattern Summary |
|--------|----------------|----------------|--------------------------|
| App Invalid 2 | 2025-09-11 21:00:00 | Earliest | Consistently catastrophic idfa_ua_ratio from the start. |
| App Invalid 1 | 2025-09-12 00:00:00 | Middle | Sharp, immediate onset of severe spoofing. |
| App Invalid 3 | 2025-09-13 00:00:00 | Latest | Delayed onset of high-intensity spoofing. |

## Visualization of the Anomaly Over Time

The image below shows the inverse relationship: as the idfa_ua_ratio (blue line, log scale) plummets, the IVT score (red line) spikes to 1.0.

- App Invalid 2 is consistently fraudulent (low blue line).
- App Invalid 1 and App Invalid 3 show a later, synchronized drop in the blue line coinciding with the red line spiking.



## 4. Final Conclusion

The core reason for IVT marking across all three apps is the presence of an extreme, non-human device spoofing pattern defined by a dangerously low idfa_ua_ratio.

- App Invalid 2 was flagged earliest because its fraudulent pattern was present from the first hours of data collection.
- App Invalid 1 and App Invalid 3 were flagged later, demonstrating that the attack traffic began on those specific dates, overwhelming the detection system.

The analysis provides a clear, quantitative defense for the IVT flagging decisions, rooted in statistical deviation from a normal traffic benchmark.

## Actionable Recommendations

Based on the quantitative findings, the following steps are recommended to improve IVT mitigation and budget protection:

1. **Strict Ratio Floor:** Implement an automated block on any traffic source (App ID/Subnet) where the idfa_ua_ratio falls below 100. This threshold is sufficient to catch the patterns observed in Invalid Apps 1, 2, and 3.
2. **Real-Time Source Blocking:** Utilize the "First IVT Time" logic to create a real-time block. Any source that maintains an IVT score > 0.5 for 3 consecutive hours should be immediately halted to prevent further wasted ad spend.
3. **Data Enrichment:** Investigate integrating device-specific parameters (e.g., screen resolution, manufacturer) into the fraud model. If multiple unique IDFAs share a single UA *and* all report identical secondary characteristics, the confidence level for device spoofing is near 100.

## Details:

**Sushmitha Vempadapu**

**9392817334**

**sushmithavempadapu3@gmail.com**

## <u>Note:</u>

**<u>Project file Done in Python Shared to mail</u>**