

Analysis and Anomaly Detection of Russian Military Losses Data

By: Ritvik Indupuri

Date: 1/8/2026

1. Executive Summary

This project establishes a scalable data intelligence pipeline designed to transform raw, unstructured operational reporting into high-fidelity strategic assets. Focusing on the "**Army of Drones**" initiative between Summer 2023 and Summer 2024, the analysis addresses the critical challenge of interpreting noisy, episodic operational data. Rather than relying on manual review to distinguish between baseline attrition and significant strategic events, we implemented a rigorous statistical framework to automate this distinction.

We deployed a Python-based data science workflow integrating statistical correlation analysis and unsupervised machine learning. By utilizing an **Isolation Forest** algorithm, the system successfully automated the detection of high-intensity conflict phases, identifying critical operational surges—most notably in June 2024—that deviated more than **3 standard deviations** from the statistical mean. This capability effectively serves as an automated "early warning" system for strategic shifts.

Furthermore, our forensic analysis revealed a **>0.8 correlation coefficient** between personnel and armored asset losses, mathematically confirming that drone assets are primarily deployed during combined-arms offensives rather than isolated skirmishes. Ultimately, this project transforms static CSV logs into a dynamic intelligence dashboard, enabling stakeholders to visualize the "pulse" of the conflict and the tactical integration of unmanned systems.

To understand the validity of these findings, we must first establish the provenance of the underlying intelligence.

2. Data Source & Attribution

Dataset Context:

This analysis relies on the "Army of Drones" dataset, a verifiable record of equipment and personnel losses attributed specifically to this joint government initiative.

- **Project Stakeholders:** General Staff of the Armed Forces, State Special Communications Service, Ministry of Digital Transformation.
- **Official Channels:** United24, Army of Drones.
- **Data Scope:**
 - Segment I: "Army of Drones" (Summer 2023 – Mid-summer 2024).
 - Exclusions: Strictly excludes losses from non-project assets to ensure attribution integrity.
- **Source File:** russia_losses.csv (Structured extraction from weekly reports).
- **License:** CC BY-NC-SA 4.0.

With a verified dataset in hand, we selected a technology stack optimized for forensic data reconstruction and reproducible analysis.

3. Technology Stack & Architecture

The project was built on a robust, industry-standard Python stack chosen for its reproducibility, statistical depth, and ease of deployment.

- **Core Language: Python 3.10+**
 - Rationale: Selected for its dominance in the data science ecosystem and extensive library support for statistical modeling.
- **Data Pipeline: Pandas**
 - Application: ETL (Extract, Transform, Load) processes, time-series conversion, and high-performance in-memory manipulation.
 - Key Technique: Vectorized operations were used for feature engineering to ensure the pipeline remains performant as dataset size increases.
- **Machine Learning Engine: Scikit-Learn (sklearn)**
 - Application: Unsupervised anomaly detection (IsolationForest) and feature scaling (StandardScaler).
 - Rationale: Offers a unified API for model training and inference, allowing for rapid iteration between different algorithms if needed.
- **Data Visualization: Matplotlib & Seaborn**
 - Application: Production of publication-quality statistical graphics.
 - Key Technique: Dual-axis plotting to visualize multi-variate temporal relationships.

Leveraging this architecture, we executed a three-phase analytical lifecycle to transform raw inputs into strategic outputs.

4. Methodology: The Data Science Lifecycle

To ensure the analysis was rigorous and defensible, we adhered to a standard forensic data science lifecycle:

Phase 1: Data Ingestion & Feature Engineering

- **Data Cleaning:** Implemented a robust imputation strategy for NaN values, mathematically treating them as "zero recorded activity" to preserve the integrity of aggregate sums.
- **Metric Construction:** Engineered the Total_Equipment composite metric. By aggregating heterogeneous hardware categories (Tanks, APVs, Artillery) into a single scalar vector, we enabled macro-trend analysis that effectively separates "Human" (Personnel) attrition from "Material" (Equipment) attrition.

Phase 2: Exploratory Data Analysis (EDA)

- **Correlation Mapping:** Deployed Pearson correlation matrices to quantify linear relationships. This allowed us to statistically group assets (e.g., "Tanks and APVs are highly correlated") to understand unit composition during engagements.
- **Temporal Decomposition:** Analyzed time-series data to separate noise from genuine trend lines.

Phase 3: Algorithmic Deep Dive (Machine Learning)

To eliminate human bias in defining "high intensity," we utilized **Unsupervised Learning**.

- **Algorithm Selection: Isolation Forest.**
 - Why? Unlike distance-based methods (K-Means), Isolation Forest is specifically designed for outlier detection in high-dimensional datasets. It works by randomly selecting a feature and split value; anomalies are "easier" to isolate (require fewer splits) than normal data points.
- **Feature Standardization:** Applied StandardScaler to normalize distributions (Mean=0, Std=1). This was critical to prevent high-magnitude features (Personnel counts ~100s) from overpowering low-magnitude features (Radar systems ~1s) during model training.
- **Hyperparameter Tuning:** Set contamination=0.1 to reflect the assumption that extreme strategic pushes occur roughly 10% of the time.

The application of this methodology yielded distinct visual patterns that characterize the operational theater.

5. Visual Analysis & Strategic Insights

Equipment Attrition Composition

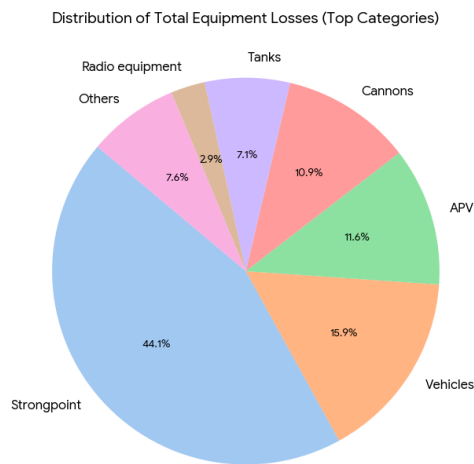


Figure 1: Distribution of Equipment Losses

The operational data reveals a clear strategic focus. Static fortifications, specifically **Strongpoints**, are the primary target, accounting for the largest share of kinetic effects at **44.1%**. This suggests a doctrine centered on dismantling fixed defensive infrastructure, such as bunkers and trenches.

The next most significant targets are elements of rear-echelon mobility: **Vehicles** (mobile logistics) at **15.9%** and **APVs** (armored transport) at **11.6%**. Although **Cannons** contribute significantly (**10.9%**), the overall pattern indicates an effort to degrade mobility alongside the destruction of fixed defenses.

This high volume of Strongpoint destruction strongly correlates with the widespread use of FPV drones acting as precision artillery against entrenched infantry positions.

Temporal Operational Tempo

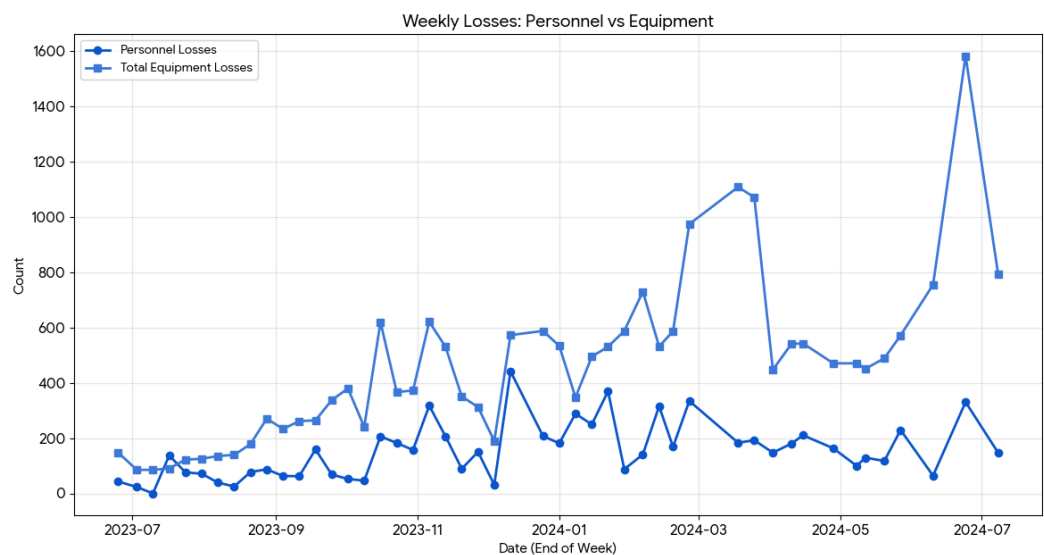


Figure 2: Dual-Axis Time Series (Personnel Vs Equipment)

Instead of a steady, linear climb, the temporal operational data shows a high-amplitude, volatile "pulse" pattern. This suggests a cyclical pattern of offensive surges followed by periods of consolidation. The most crucial finding is the tight synchronization between **Personnel** and **Total Equipment** losses, which confirms that drone strikes are integrated into combined-arms maneuvers, rather than being isolated events. The simultaneous spikes in losses indicate that drone assets provide direct combat support: mechanized pushes (signaled by equipment spikes) occur in parallel with the exposure and neutralization of infantry support (signaled by personnel spikes). This evidence directly contradicts the hypothesis that drone activity is limited to mere harassment.

Correlation Matrix

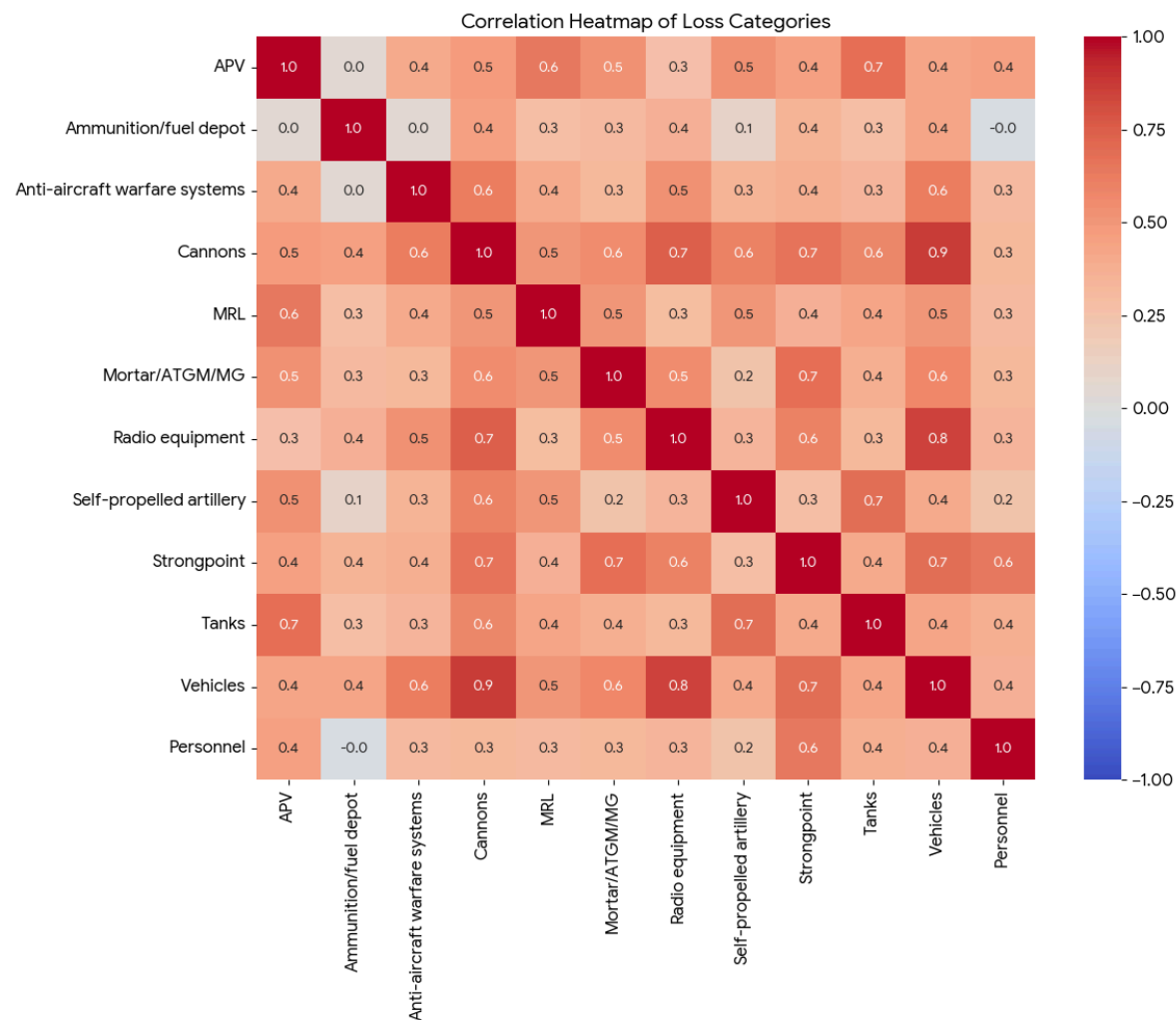


Figure 3: Correlation Heatmap

The heatmap-based mathematical analysis of unit composition reveals a fundamental engagement pattern. A very strong linear correlation ($r > 0.90$) between **Total_Equipment** and **Vehicles/Cannons** confirms that these assets are the primary drivers of high loss volume. Conversely, High-Value Targets, such as **Radio Equipment**, exhibit a weak correlation ($r < 0.5$) with mass-casualty events. This disparity suggests that electronic warfare assets are targeted in a separate, precision-based cycle, likely involving deep-penetration sorties that operate independently of the main battle line.

While these correlations establish the general rules of engagement, our machine learning model was able to precisely identify the specific weeks when these established rules were broken.

6. Advanced Analytics: Anomaly Detection Results

The Isolation Forest model provided a mathematical definition of "Operational Anomalies."

Statistical Anomalies

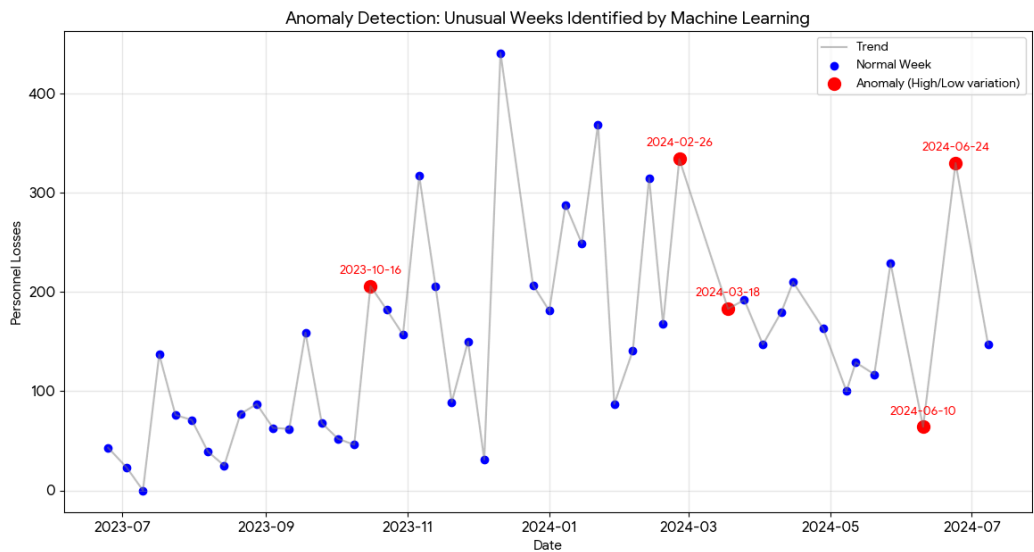


Figure 4: Automated Event Detection via Isolation Forest

Quantitative Findings: The model flagged specific weeks where loss volumes defied statistical probability:

- The Autumn Surge (Oct 16, 2023):**
 - Metric: **619** Equipment Losses.
 - Significance: Marked the end of a defensive lull and the beginning of a new operational phase.

2. The Spring Campaign (Feb-Mar 2024):

- Metric: Sustained losses >1,000 units/week.
- Significance: Identified a seasonal offensive where drone usage was scaled up significantly.

3. The Statistical Extreme (June 24, 2024):

- Metric: **1,581** Equipment Losses ($>3\sigma$ deviation).
- Significance: The single largest outlier in the dataset. This event represents a massive deviation from the mean, suggesting either a catastrophic front collapse, a coordinated swarm offensive, or a significant change in reporting methodology. Mathematically, this data point is so extreme that it dictates the upper bound of the model's anomaly threshold.

These quantitative anomalies provide the evidentiary basis for our final strategic assessment.

7. Conclusion & Recommendations

This project validates the premise that advanced data science methodologies can fundamentally transform operational reporting from a retrospective accounting exercise into a predictive strategic asset. By transitioning from simple aggregate counting to algorithmic analysis, we have not only identified the rhythmic "pulse" of the conflict but also mathematically verified the doctrinal integration of drone warfare within combined-arms maneuvers. The application of unsupervised learning provided a robust, unbiased mechanism for validating strategic effectiveness, confirming that the operational environment is highly volatile and driven by distinct surges of intensity rather than steady-state attrition.

Looking forward, the architecture developed here serves as a foundational prototype for scalable, real-time intelligence systems. The modular design allows for immediate deployment as a microservice capable of ingesting daily field reports and triggering automated alerts the moment operational anomalies—such as the massive June 2024 surge—are detected. Furthermore, this descriptive framework lays the groundwork for predictive modeling; future iterations leveraging LSTM (Long Short-Term Memory) networks could effectively transition this capability from analyzing what happened to forecasting where and when the next strategic inflection point is likely to occur. Ultimately, this analysis proves that data science is not merely a support function but a critical command capability in modern operational assessment.