

# **District-Level Aadhaar Operational Stress and Inclusion Risk Analysis**

UIDAI Data Hackathon Submission

Submitted by: **Abhishek Kumar**

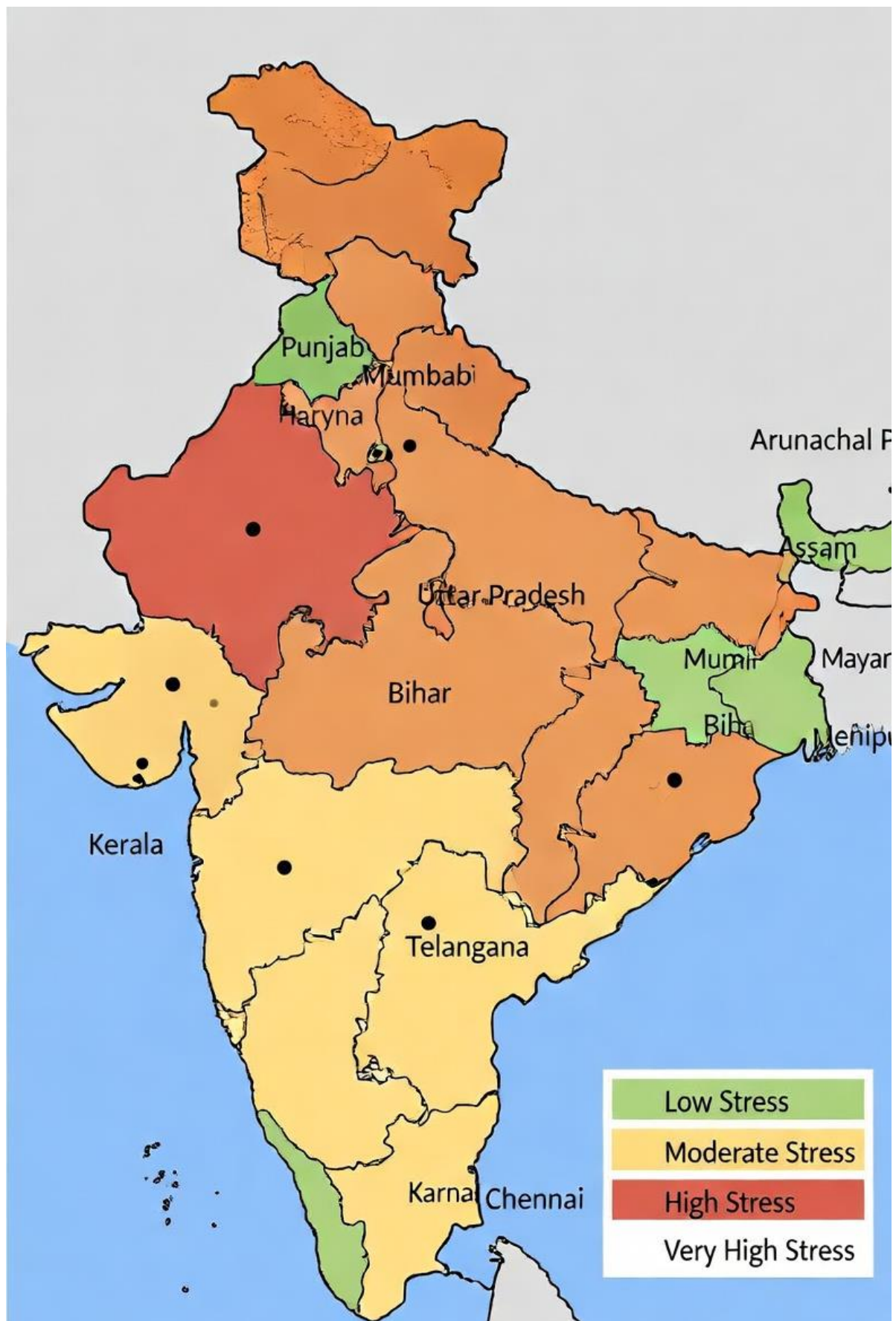
Date: **7 January 2026**

## **1.Problem Statement and Approach**

### **Problem Statement**

Aadhaar is a foundational digital identity system that supports a wide range of public services, including welfare delivery, education benefits, and financial inclusion. While Aadhaar enrolment has achieved wide coverage, operational challenges continue to arise due to repeated biometric and demographic updates across regions. These challenges are not uniformly distributed and may vary significantly at the district level.

### **STRESS LEVEL**



In particular, certain districts experience disproportionately high Aadhaar update activity relative to enrolments, leading to operational stress, service delays, and potential authentication issues. Additionally, a large share of update activity is driven by the under-18 population, creating long-term risks due to frequent biometric changes during growth years. Identifying such localized stress patterns and inclusion gaps is critical for efficient resource allocation and sustained Aadhaar service reliability.

## **Approach**

This project adopts a district-level analytical approach to assess Aadhaar operational stress and inclusion risks using enrolment, biometric update, and demographic update datasets provided by UIDAI. Raw data from multiple files was consolidated and aggregated at the district level to ensure scalability, interpretability, and privacy preservation.

A unified analytical framework was developed to derive three complementary indicators:

**1.Aadhaar System Stress Indicator**, measuring the ratio of total updates to enrolments,

**2.Child Aadhaar Update Risk Zones**, identifying districts where biometric updates are predominantly driven by the under-18 population, and

**3.Aadhaar Inclusion Gap Index**, highlighting disparities between enrolment presence and update dependency.

Through statistical analysis and visualization, the framework reveals patterns, anomalies, and priority districts requiring targeted interventions. The approach emphasizes actionable insights over nationwide generalizations, enabling UIDAI to adopt district-specific operational strategies.

## 2.Datasets Used

This analysis uses Aadhaar enrolment and update datasets provided by UIDAI as part of the hackathon. The datasets capture district-level information on Aadhaar enrolments and subsequent biometric and demographic updates. Multiple data files for each dataset were consolidated to ensure complete coverage before analysis.

Dataset Name	Description	Key Columns Used
Aadhaar Enrolment	Contains district-level	state, district, age_5_17,

Dataset	Aadhaar enrolment counts segmented by age groups	age_18_greater
Aadhaar Biometric Update Dataset	Records biometric update activity across districts, including updates related to children	state, district, bio_age_5_17
Aadhaar Demographic Update Dataset	Captures demographic update activity across districts	state, district, demo_age_5_17

All datasets were sourced from the UIDAI open data provided for the hackathon and do not contain any personally identifiable information. The analysis was performed exclusively on aggregated district-level data to ensure privacy preservation and ethical use of Aadhaar data.

### 3.Methodology

The methodology adopted in this project focuses on transforming large-scale Aadhaar enrolment and update data into meaningful, district-level indicators that can support operational decision-making. The process consists of data consolidation, cleaning, aggregation, and feature engineering, as described below.

### **3.1 Data Loading and Consolidation**

The Aadhaar datasets were provided as multiple CSV files for enrolment, biometric updates, and demographic updates. All files belonging to the same dataset category were programmatically loaded and consolidated to ensure complete coverage. This approach allowed scalable processing of large datasets while avoiding duplication or data loss.

### **3.2 Data Cleaning and Preprocessing**

To ensure consistency across datasets, column names were standardized by converting them to lowercase and removing whitespace. Records with malformed or invalid state and district identifiers were removed. Missing values arising from outer joins during dataset integration were handled by appropriate imputation to maintain analytical continuity.

Only district-level attributes were retained for further analysis, ensuring privacy preservation and eliminating any dependency on individual-level information.

### **3.3 District-Level Aggregation**

Given the scale of the raw data, all datasets were aggregated at the district level, which represents the most relevant unit for Aadhaar operational planning. Aggregation significantly reduced data volume while retaining meaningful signals.

The following aggregations were performed:

Total Aadhaar enrolments per district

Total biometric update counts per district

Total demographic update counts per district

Under-18 biometric update counts per district

This step transformed millions of raw records into a manageable dataset suitable for comparative analysis across districts.

### **3.4 Feature Engineering and Indicator Design**

Based on the aggregated data, multiple analytical indicators were derived:

### **Aadhaar System Stress Indicator**

Defined as the ratio of total Aadhaar updates to total enrolments for each district. This indicator captures operational pressure arising from repeated update activity.

### **Child Aadhaar Update Risk Zones**

Districts were classified based on the proportion of biometric updates attributable to the under-18 population. Districts with a dominant share of child updates were identified as high-risk zones due to the likelihood of frequent future updates.

### **Aadhaar Inclusion Gap Index**

A normalized measure derived from update-to-enrolment ratios, highlighting disparities between enrolment presence and update dependency across districts.



Additionally, districts were categorized into Low, Medium, and High stress categories using threshold-based classification to support prioritization and targeted interventions.

### **3.5 Analytical Focus**

The methodology emphasizes interpretability, scalability, and policy relevance. By operating exclusively on aggregated district-level data, the framework ensures ethical data usage while enabling UIDAI to identify patterns, anomalies, and priority regions requiring operational attention.

## **4. Data Analysis and Visualisation**

This section presents the key analytical findings derived from the district-level Aadhaar Operational Stress Framework, supported by visualisations to highlight patterns, trends, and anomalies across region

### **4.1 Key Findings**

The analysis reveals clear and interpretable patterns in Aadhaar operational activity across districts:

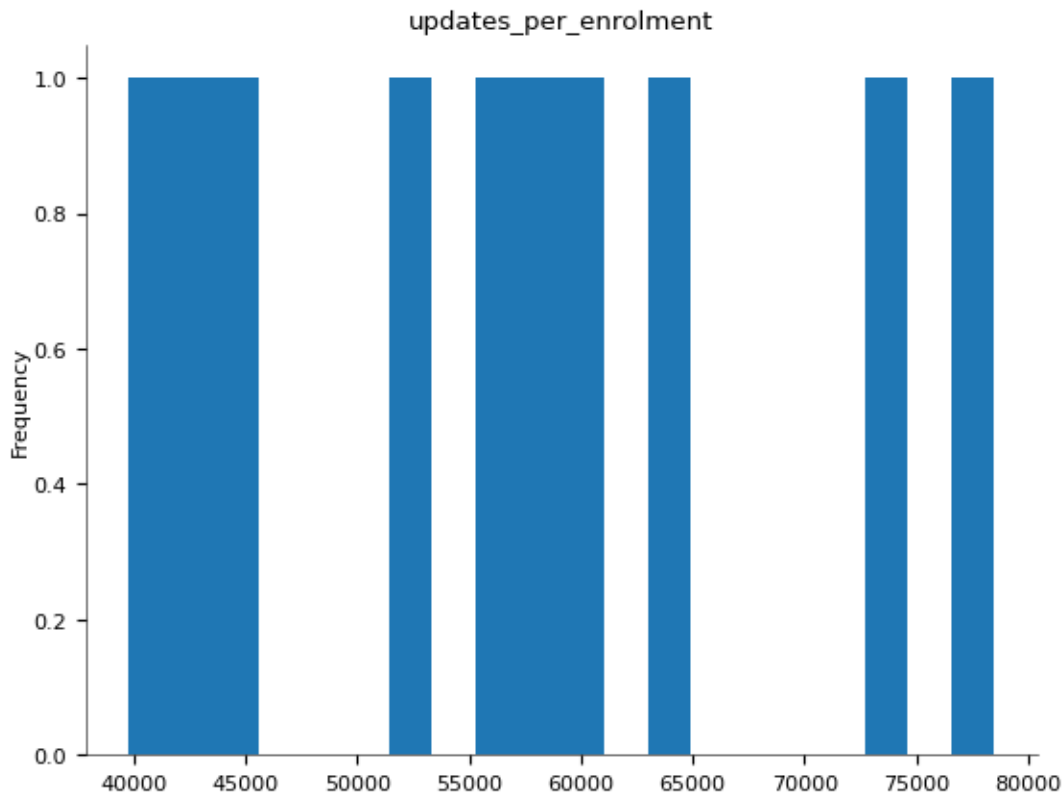
- A majority of districts exhibit low operational stress, indicating overall system stability.
- Approximately 311 districts fall under medium stress, while 67 districts exhibit high operational stress, requiring targeted attention.
- No districts were classified under extreme stress, suggesting the absence of a nationwide systemic failure.
- Child Aadhaar updates dominate biometric update activity, with over 90% of districts classified as high child-risk zones.
- All districts classified under medium and high operational stress also fall within child-risk zones, indicating a strong linkage between child-centric updates and system stress.

These findings demonstrate that Aadhaar operational challenges are localized and child-driven, rather than uniform across states or regions.

## **4.2 Distribution of Aadhaar Operational Stress**

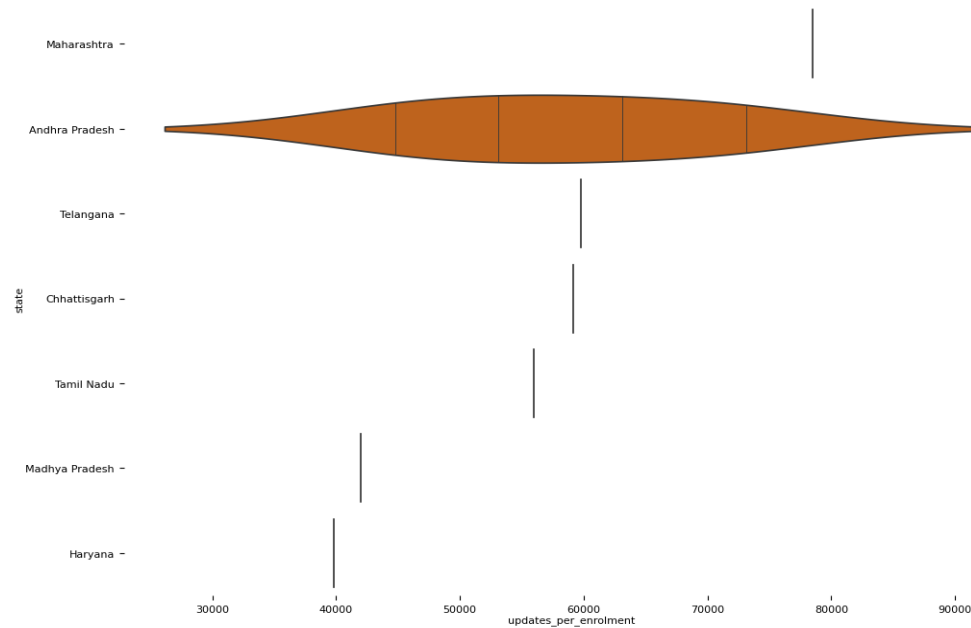
The histogram illustrates the distribution of the Aadhaar System Stress Indicator (updates per enrolment) across all districts. Most districts are clustered at lower

stress values, while a small number of districts form a long tail with significantly higher stress levels. This confirms that operational pressure is concentrated in a limited number of districts



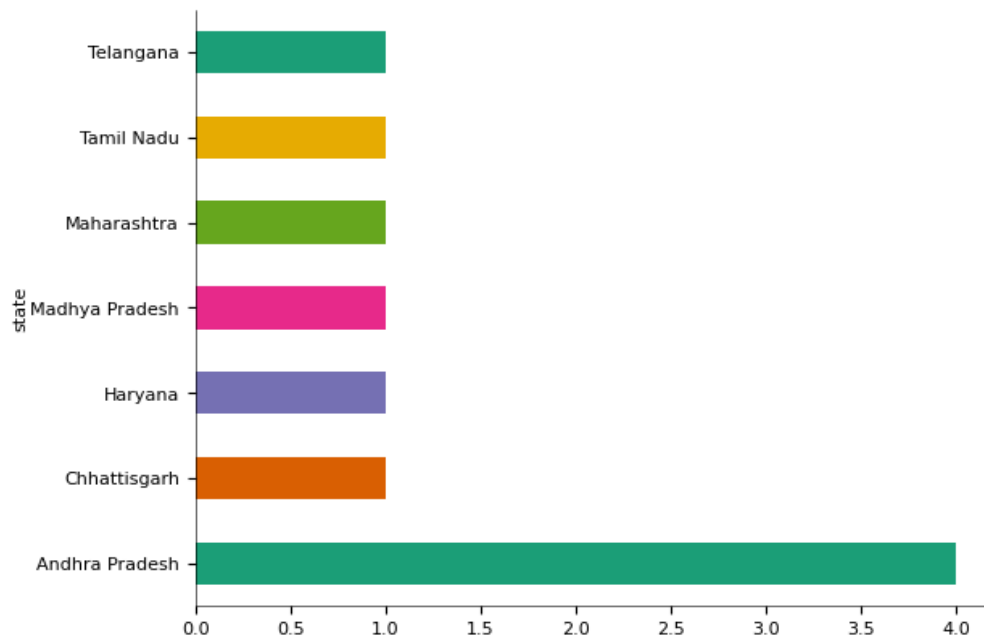
### 4.3 State-wise Variability in District-Level Stress

The violin plot highlights variation in operational stress within states. Several states exhibit wide distributions, indicating that stress is driven by specific districts rather than being uniform across the entire state. This reinforces the need for district-specific planning instead of broad state-level interventions.



## 4.4 State Prioritization Based on High-Stress Districts

This visualization ranks states based on the count of high-stress districts. A small number of states contribute disproportionately to operational stress, enabling UIDAI to prioritize resources and interventions more effectively.



## **4.5 Child Aadhaar Update Risk Analysis**

The child Aadhaar update analysis reveals that 1,030 out of 1,132 districts are classified as high child-risk zones, where under-18 biometric updates account for the majority of update activity. This dominance explains the observed operational stress patterns and indicates future risks associated with biometric ageing, authentication failures, and increased update demand as children grow.

## **4.6 Combined Stress and Child Risk Perspective**

An integrated view combining operational stress categories and child-risk zones shows that all medium and high stress districts are also child-risk dominant. This finding establishes a strong causal relationship between child-centric update patterns and Aadhaar operational stress, underscoring the importance of child-focused intervention strategies.

## **5. Code Snippets**

This section presents selected code snippets used to implement the Aadhaar Operational Stress and Inclusion Risk Analysis. The full notebook is available and can be shared separately if required.

### **5.1 Data Loading and Consolidation**

[1]  
✓ 0s

```
import pandas as pd
import glob
```

[2]  
✓ 4s

```
def load_and_combine(pattern):
    files = glob.glob(pattern)
    return pd.concat(
        (pd.read_csv(f) for f in files),
        ignore_index=True
    )

bio_df = load_and_combine('/content/api_data_aadhar_biometric_*.csv')
demo_df = load_and_combine('/content/api_data_aadhar_demographic_*.csv')
enrol_df = load_and_combine('/content/api_data_aadhar_enrolment_*.csv')
```

## 5.2 District-Level Aggregation

```
bio_agg = bio_df.groupby(['state', 'district'], as_index=False).agg(
    total_bio_updates=('bio_age_5_17', 'sum'),
    under18_bio_updates=('bio_age_5_17', 'sum')
)

demo_agg = demo_df.groupby(['state', 'district'], as_index=False).agg(
    total_demo_updates=('demo_age_5_17', 'sum'),
    under18_demo_updates=('demo_age_5_17', 'sum')
)

enrol_agg = enrol_df.groupby(['state', 'district'], as_index=False).agg(
    total_enrol=('age_18_greater', 'sum'),
    child_enrol=('age_5_17', 'sum')
)
```

## 5.3 Indicator Computation

```
merged_df = enrol_agg.merge(bio_agg, on=['state','district'], how='outer') \
    .merge(demo_agg, on=['state','district'], how='outer') \
    .fillna(0)

merged_df['total_updates'] = (
    merged_df['total_bio_updates'] + merged_df['total_demo_updates']
)

merged_df['updates_per_enrolment'] = (
    merged_df['total_updates'] / merged_df['total_enrol'].replace(0,1)
)
```

## 5.4 Stress and Risk Classification

```
merged_df['stress_category'] = pd.cut(
    merged_df['updates_per_enrolment'],
    bins=[-1, 1000, 10000, 100000, float('inf')],
    labels=['Low', 'Medium', 'High', 'Extreme']
)

merged_df['child_risk_zone'] = merged_df['child_bio_ratio'].apply(
    lambda x: 'High Child Risk' if x >= 75 else 'Normal'
)
```

The above code snippets demonstrate the reproducible analytical pipeline used to derive district-level Aadhaar operational stress indicators and inclusion risk zones.

## 6. Conclusion and Recommendations

### Conclusion

This project demonstrates that Aadhaar operational challenges are not uniformly distributed across regions, but are instead concentrated within a limited number of districts. By aggregating enrolment and update data at the district level, the analysis reveals that the Aadhaar ecosystem remains largely stable overall, with localized pockets of operational stress.

A key finding is the dominance of child-centric biometric updates, with over ninety percent of districts classified as high child-risk zones. Importantly, all districts identified as medium or high operational stress also fall within these child-risk zones, establishing a strong linkage between child update patterns and system stress. These insights highlight the need for targeted, population-specific interventions rather than broad, nationwide policy changes.

## **Recommendations**

Based on the findings, the following actionable recommendations are proposed for UIDAI consideration:

### **1. District-Targeted Aadhaar Update Drives**

Deploy mobile Aadhaar update units and additional staffing selectively in high-stress districts, rather than uniformly across states.



## **2. Child-Focused Update Interventions**

Integrate Aadhaar biometric update drives with schools and education-linked services to proactively address child-centric update demand and reduce future operational pressure.

## **3. Operational Stress Monitoring Framework**

Adopt the updates-per-enrolment indicator as a regular monitoring metric to identify emerging stress zones and enable early intervention.

## **4. Data-Driven Resource Allocation**

Utilize district-level stress and child-risk classifications to guide infrastructure planning, manpower deployment, and scheduling of Aadhaar services.

The proposed analytical framework is scalable, privacy-preserving, and reproducible, making it suitable for integration into UIDAI's operational planning and monitoring processes.