

1. Problem Statement & Approach

Despite Aadhaar's objective of universal and early identity coverage, enrolment patterns across India show significant inter-state and inter-district.

Hence,

Is Aadhaar enrolment occurring proactively at early life stages, or reactively when individuals require identity for welfare, employment, or compliance?

Why this matters:

- Delayed enrolment = exclusion risk
- Reactive enrolment = system stress
- Early enrolment = efficient lifecycle identity

Analytical Approach

This study constructs a **composite Enrollment Health Index** using age-wise Aadhaar enrolment data to:

- Measure balance between child and adult enrolment
- Identify structurally delayed states and districts
- Detect systemic enrolment distortions rather than one-off gaps

We intentionally focus on **enrolment-only data** to demonstrate how maximum policy insight can be extracted from minimum datasets.

2. Datasets Used

UIDAI Aadhaar Enrolment Data

Columns used:

- State
- District
- Age-wise enrolment counts:
 - 0–5 years
 - 5–17 years
 - 18+ years

- Date

Why ONLY Enrolment Dataset Was Used

Enrolment data is the **first and most foundational signal** of identity inclusion. Before updates, biometrics, or corrections occur, **access itself must exist**.

3. Methodology

The Aadhaar enrolment dataset provided by UIDAI was subjected to a structured cleaning and preprocessing pipeline to ensure analytical reliability and policy interpretability.

3.1 Data Cleaning & Preprocessing

3.1.1 Standardization of Administrative Units

The dataset contained variations in the naming of states and districts due to differences in spelling, casing, and formatting across reporting periods. To ensure consistency:

- State and district names were normalised using string standardisation techniques (case normalization, trimming of whitespace, and removal of special characters).
- A single canonical naming convention was enforced across the dataset to enable accurate aggregation and comparison.
- This step was critical to avoid artificial fragmentation of enrolment counts caused by inconsistent administrative labels.

3.1.2 Handling of Duplicate Records

During preprocessing, a substantial number of duplicate rows were detected. These duplicates did not uniformly represent data errors; instead, they reflected distributed reporting of enrolment counts across multiple records for the same administrative unit and time period.

To address this appropriately:

- **Exact duplicate rows** (identical across all relevant columns) were removed to eliminate redundant records.
- **Logical duplicates**—where the same state, district, and time period appeared multiple times with partial enrolment counts—were aggregated rather than dropped.
- Aggregation was performed by summing age-wise enrolment counts (0–5 years, 5–17 years, and 18+ years) for each unique state–district–time combination.

This dual approach ensured that genuine enrolment volumes were preserved while preventing inflation or loss of data.

3.1.3 Age-wise Enrolment Validation

Age-category enrolment columns were validated to ensure:

- All enrolment counts were non-negative integers.
- No missing values existed in age-wise enrolment fields after aggregation.
- Records with incomplete administrative identifiers were excluded, as they could not be meaningfully mapped for policy analysis.

3.1.4 Age-wise Enrolment Validation

Age-category enrolment columns were validated to ensure:

- All enrolment counts were non-negative integers.
- No missing values existed in age-wise enrolment fields after aggregation.
- Records with incomplete administrative identifiers were excluded, as they could not be meaningfully mapped for policy analysis.

3.1.5 Analytical Readiness

The final cleaned dataset represents a **district-level, policy-ready enrolment table**, with consistent administrative identifiers and validated age-wise enrolment counts. This processed dataset served as the foundation for:

- Construction of derived enrolment ratios
- Development of the composite Enrollment Health Index
- State and district-level comparative analysis

Rather than mechanically removing duplicates or inconsistencies, the preprocessing approach prioritised **preserving administrative reality**. Given that Aadhaar enrolment data is operational in nature, aggregation-based duplicate handling ensures that the analysis reflects **actual enrolment activity**, not reporting artefacts.

3.2 Feature Engineering

Following data cleaning and consolidation, a set of analytically meaningful features was engineered to move beyond raw enrolment counts and capture **structural patterns in enrolment behaviour** at the district and state level.

The objective of feature engineering was not prediction, but **diagnostic insight** — identifying whether enrolment systems function proactively, consistently, and equitably across geography and time.

3.2.1 Age-Structure Ratios

Raw enrolment volumes alone do not reflect enrolment health. To capture enrolment balance across life stages, age-wise ratios were constructed:

- **Child Enrolment Ratio**
Ratio of enrolments in the 0–5 age group to total enrolment.
- **Adult Enrolment Ratio**
Ratio of enrolments in the 18+ age group to total enrolment.

These ratios enable identification of districts where Aadhaar enrolment is **reactive** (adult-heavy) versus **proactive** (early-life enrolment).

3.2.2 Enrolment Balance Score

To quantify the relative dominance of early versus late enrolment, a **Balance Score** was engineered as:

Child Enrolment Ratio – Adult Enrolment Ratio

This score captures the direction of enrolment behaviour:

- Higher values indicate early-life, system-driven enrolment.
- Lower values indicate delayed, need-driven enrolment (often linked to service access requirements).

To ensure comparability across districts, the balance score was **min–max normalised** to a 0–1 scale.

3.2.3 Consistency Score (Temporal Behaviour)

Administrative systems often exhibit deadline-driven behaviour, where enrolments spike at the end of months or quarters due to reporting pressures.

To capture this phenomenon:

- A **deadline enrolment share** was calculated using month-end and quarter-end enrolment indicators.
- A **Consistency Score** was derived as the inverse of this share, bounded between 0 and 1.

Higher consistency scores indicate **steady, distributed enrolment activity**, while lower scores signal reliance on deadline-driven bursts.

3.2.4 Geographic Spread Score (Pincode Equity)

To assess whether enrolment activity is geographically concentrated or equitably distributed, a **Spread Score** was constructed:

- Ratio of active pincodes to total reported pincodes within a district.

This feature highlights:

- High spread → broad geographic access and decentralised enrolment
- Low spread → spatial concentration, potentially indicating access gaps

The score was clipped to ensure valid bounds between 0 and 1.

3.2.5 Composite Enrollment Health Index (EHI)

The engineered features were integrated into a composite **Enrollment Health Index (EHI)** to provide a single, interpretable measure of enrolment quality.

The index combines:

- Balance Score (50% weight)
- Consistency Score (30% weight)
- Spread Score (20% weight)

Weights were assigned to prioritise **early-life enrolment behaviour**, while still accounting for operational stability and geographic reach.

The final index was scaled to a **0–100 range** to facilitate:

- State and district-level comparison
- Heatmap visualisation
- Policy communication

UIDAI enrolment data is inherently operational. Feature engineering transformed this operational data into **policy-relevant signals**, enabling identification of:

- Structural enrolment risks
- Administrative behaviour patterns
- Regions requiring systemic intervention rather than volume expansion

This step was foundational in converting raw enrolment records into **actionable governance insights**.

3.3 Data Analysis & Visualisation

The analytical phase focused on translating the engineered indicators into clear, interpretable insights for administrators and policymakers. Rather than isolated metrics, the analysis emphasized **comparative patterns, structural disparities, and behavioral signals** within the Aadhaar enrolment system.

Visualizations were designed to support diagnosis, prioritization, and communication, not just exploration.

3.3.1 District-Level Enrollment Health Analysis

At the most granular level, the Enrollment Health Index (EHI) was analysed across districts to identify intra-state heterogeneity.

Key analyses included:

- Ranking districts by EHI to identify high-performing and structurally vulnerable regions
- Comparing districts with similar enrolment volumes but divergent EHI scores, highlighting the difference between **quantity** and **quality** of enrolment

- Identifying districts with extreme values on individual components (balance, consistency, spread)

This analysis revealed that high enrolment activity does not necessarily correspond to healthy enrolment systems, reinforcing the need for composite evaluation.

3.3.2 State-Level Aggregation and Comparison

District-level EHI values were aggregated to the state level using mean and median statistics to obtain a national view.

State-wise analysis enabled:

- Comparison of enrolment health across states, independent of population size
- Identification of states with internally uneven performance (high mean but low median, or vice versa)
- Segmentation of states into relative enrolment health tiers

This aggregation formed the basis for national benchmarking while preserving sensitivity to district-level variation.

3.3.3 National Enrollment Health Heatmap

A national heatmap was developed to visualise state-wise EHI at a glance.

The heatmap:

- Uses colour intensity to represent average EHI values
- Enables rapid identification of regions with comparatively weaker enrolment health
- Supports intuitive policy communication for non-technical stakeholders

By ordering states according to EHI, the heatmap highlights structural gradients rather than isolated outliers, making regional patterns immediately visible.

3.3.4 Top and Bottom State Analysis

To support prioritisation, states were ranked to identify:

- **Top-performing states**, demonstrating strong early-life enrolment balance, operational consistency, and geographic spread

- **Lower-performing states**, where enrolment systems appear reactive, deadline-driven, or spatially concentrated

Importantly, this ranking is diagnostic, not normative. Lower EHI does not imply administrative failure, but signals areas where systemic reinforcement or targeted intervention may yield high returns.

3.3.5 Behavioural Pattern Visualisation (District Scatter Analysis)

Scatter-based visual analysis was employed to examine behavioural dimensions simultaneously, such as:

- Active enrolment days versus end-period enrolment concentration
- Districts exhibiting deadline-dependent enrolment behaviour
- Outlier districts that deviate significantly from national patterns

These visuals make administrative behaviour visible, enabling a shift from volume-based monitoring to **process-aware oversight**.

3.3.6 Policy-Relevant Insights

Across analyses and visualisations, three consistent insights emerged:

1. **Early-life enrolment is unevenly institutionalised**, with several regions relying heavily on adult enrolment.
2. **Deadline-driven enrolment behaviour persists** in multiple districts, indicating reporting pressure effects.
3. **Geographic access remains inconsistent**, as reflected by uneven pincode-level activity.

These insights collectively support the argument that enrolment quality should be monitored alongside enrolment volume.

4. Core Visualisations

To ensure clarity and decision usefulness, the analysis deliberately limits visual outputs to a small set of **high-impact, governance-ready visuals**. Each visual addresses a **specific policy or administrative question**, avoiding redundancy or exploratory clutter.

4.2.1 Enrollment Health Index (EHI) – State Ranking

Governance Question:

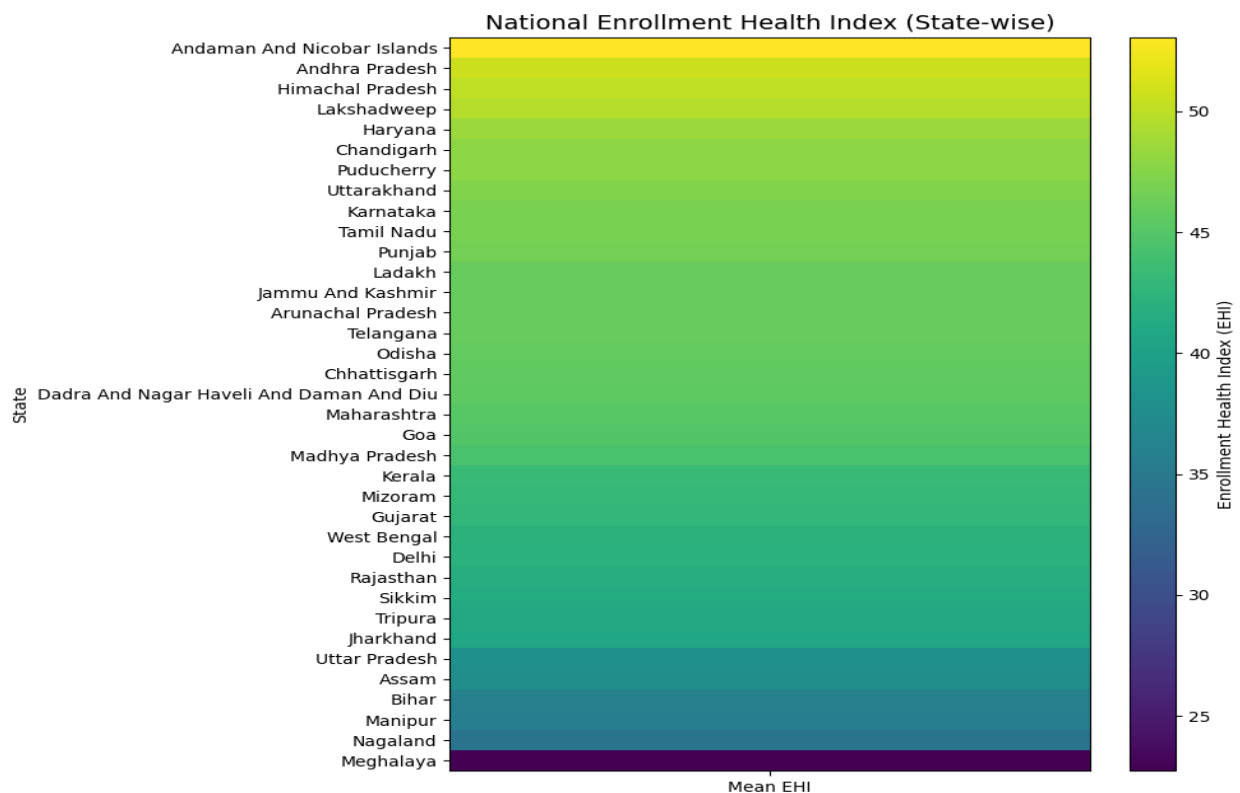
Which states exhibit comparatively stronger or weaker Aadhaar enrolment system health?

Description:

A state-wise ranking of the Enrollment Health Index (EHI), derived from district-level scores and aggregated using mean and median values.

Why this matters:

- Enables national benchmarking across states
- Shifts evaluation from enrolment volume to enrolment quality
- Supports high-level prioritisation without population bias



Policy Use:

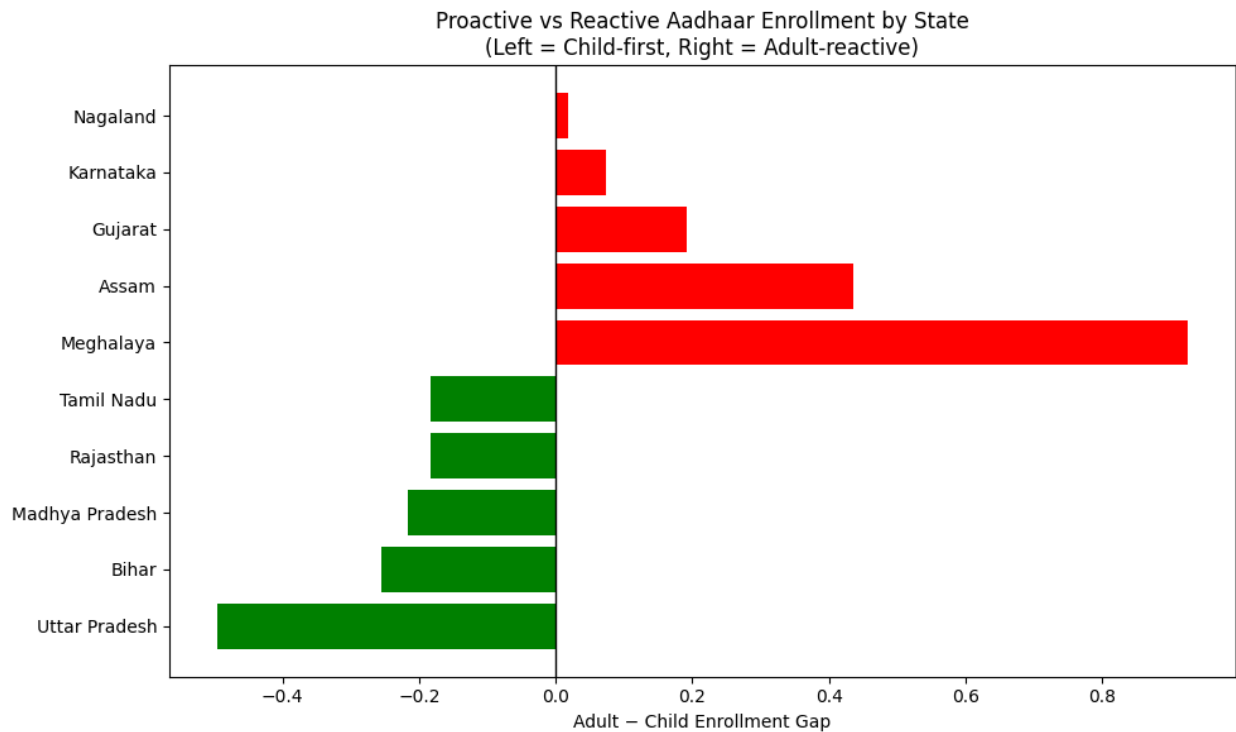
This ranking can inform:

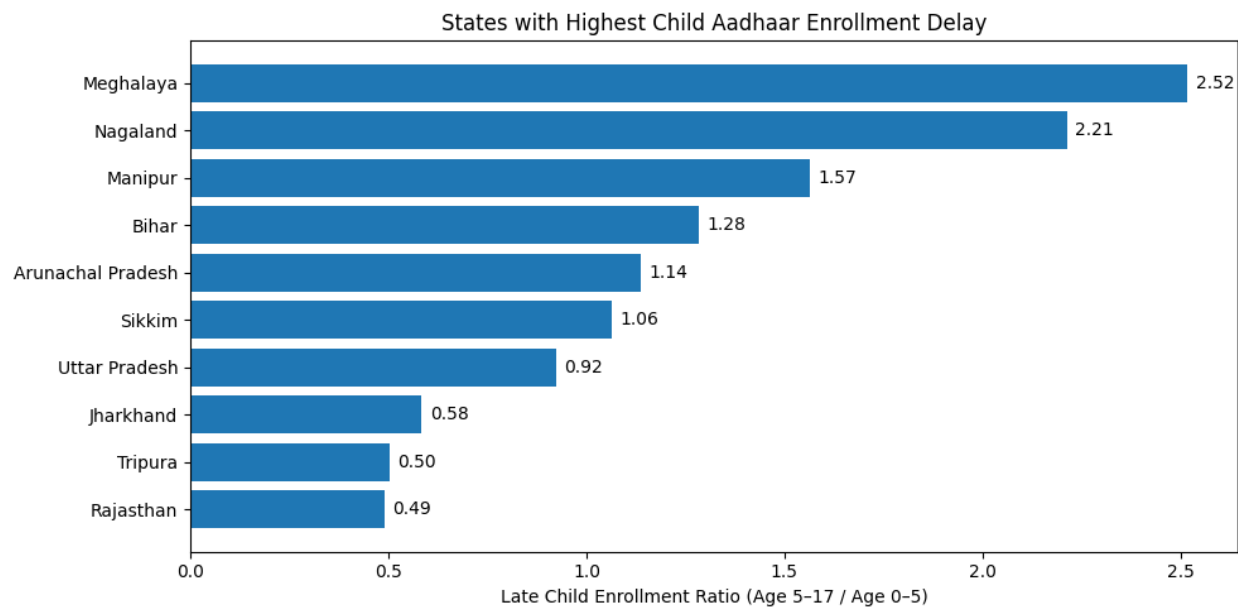
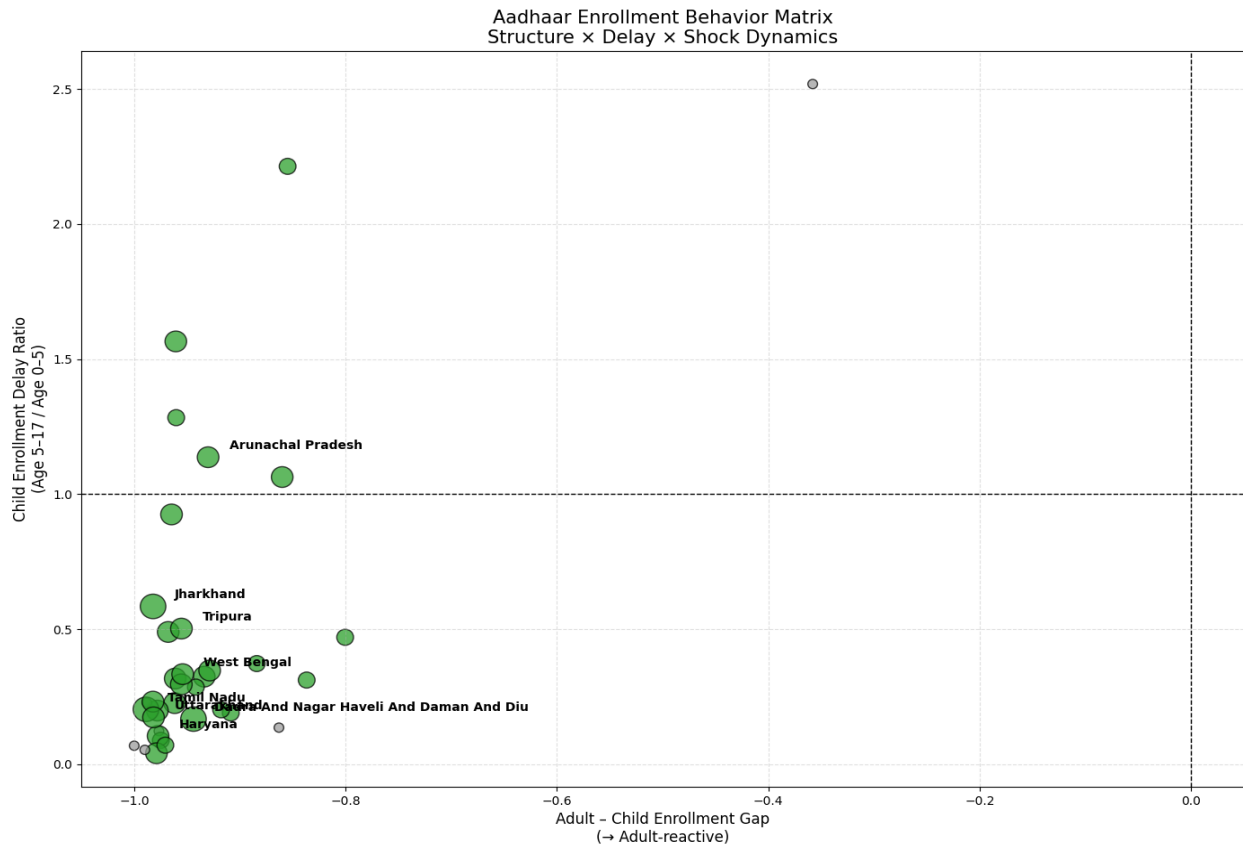
- Targeted administrative reviews
- Inter-state knowledge sharing
- Identification of best practices from higher-performing states

4.2.2 Adult vs Child Enrolment Scatter

Governance Question:

Is enrolment occurring proactively at birth/early childhood, or reactively during adulthood?





Why this matters:

- Highlights structural dependence on adult enrolment
- Reveals gaps in early-life Aadhaar institutionalisation
- Differentiates awareness-related gaps from access-related gaps

Policy Use:

Supports decisions on:

- Awareness campaigns for early enrolment
- Integration with birth registration systems
- Differentiated outreach strategies

4.2.3 Bottom 10 States by Enrollment Health Index

Governance Question:

Which states require immediate analytical or operational attention?

Description:

A tabular presentation of the bottom 10 states by mean EHI, ordered from **lowest (most structurally vulnerable)** to relatively higher values.

Why this matters:

- Prevents misinterpretation caused by unordered or descending tables
- Enables precise identification of highest-risk regions
- Maintains analytical integrity and interpretability

Bottom 10 States (Enrollment Risk Zones)			
	mean_EHI	median_EHI	districts
state			
Meghalaya	22.731818	22.290	11
Nagaland	34.415385	33.950	13
Manipur	35.842000	35.925	10
Bihar	35.918684	35.795	38
Assam	37.714545	38.840	33
Uttar Pradesh	37.906400	37.810	75
Jharkhand	40.804583	40.965	24
Tripura	41.076250	41.325	8
Sikkim	41.330000	37.305	4
Rajasthan	41.675152	41.980	33

Top 10 Healthiest States by EHI

state	mean_EHI	median_EHI	districts
Andaman And Nicobar Islands	53.046667	50.680	3
Andhra Pradesh	50.796061	48.200	33
Himachal Pradesh	50.087500	49.755	12
Lakshadweep	49.650000	49.650	1
Haryana	48.430000	48.830	22
Chandigarh	47.950000	47.950	1
Puducherry	47.885000	47.885	2
Uttarakhand	47.274615	47.440	13
Karnataka	46.959032	47.680	31
Tamil Nadu	46.913333	46.280	36

Policy Use:

This table supports:

- Priority flagging for monitoring
- Resource allocation discussions
- Further diagnostic deep-dives at the district level

4.2.4 District-Level Outlier Visualisation

Governance Question:

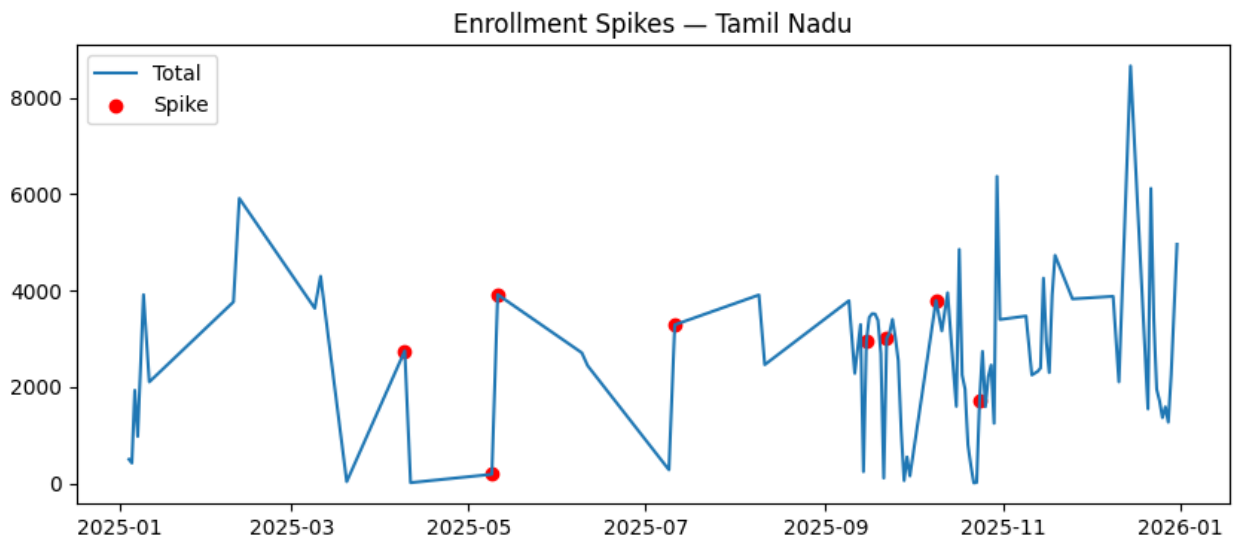
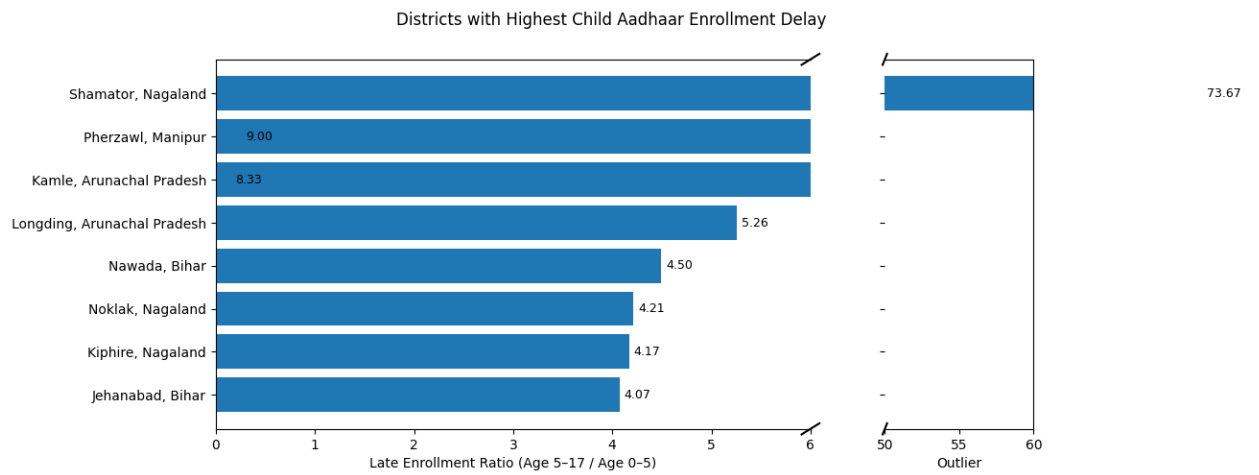
Which districts deviate significantly from expected enrolment behaviour?

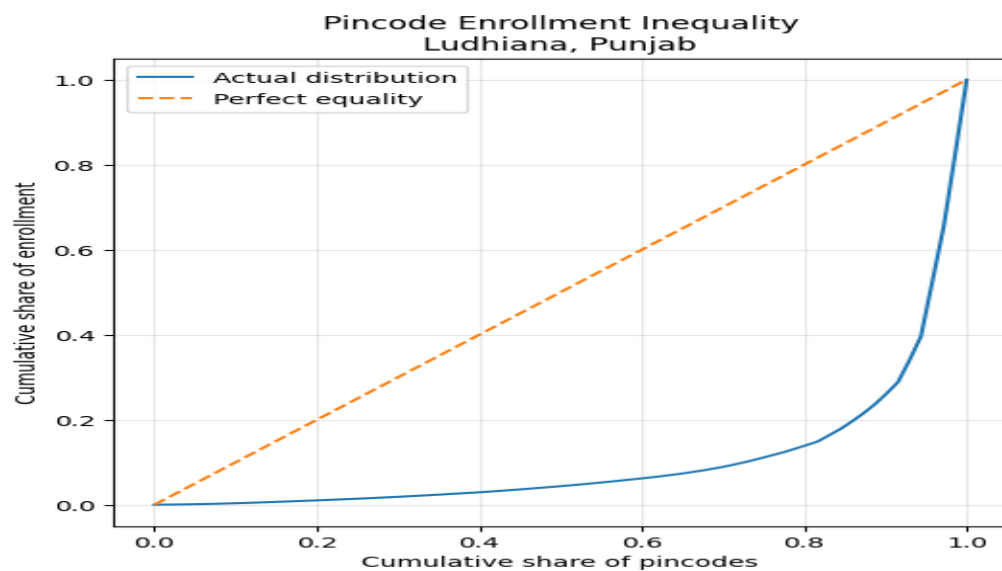
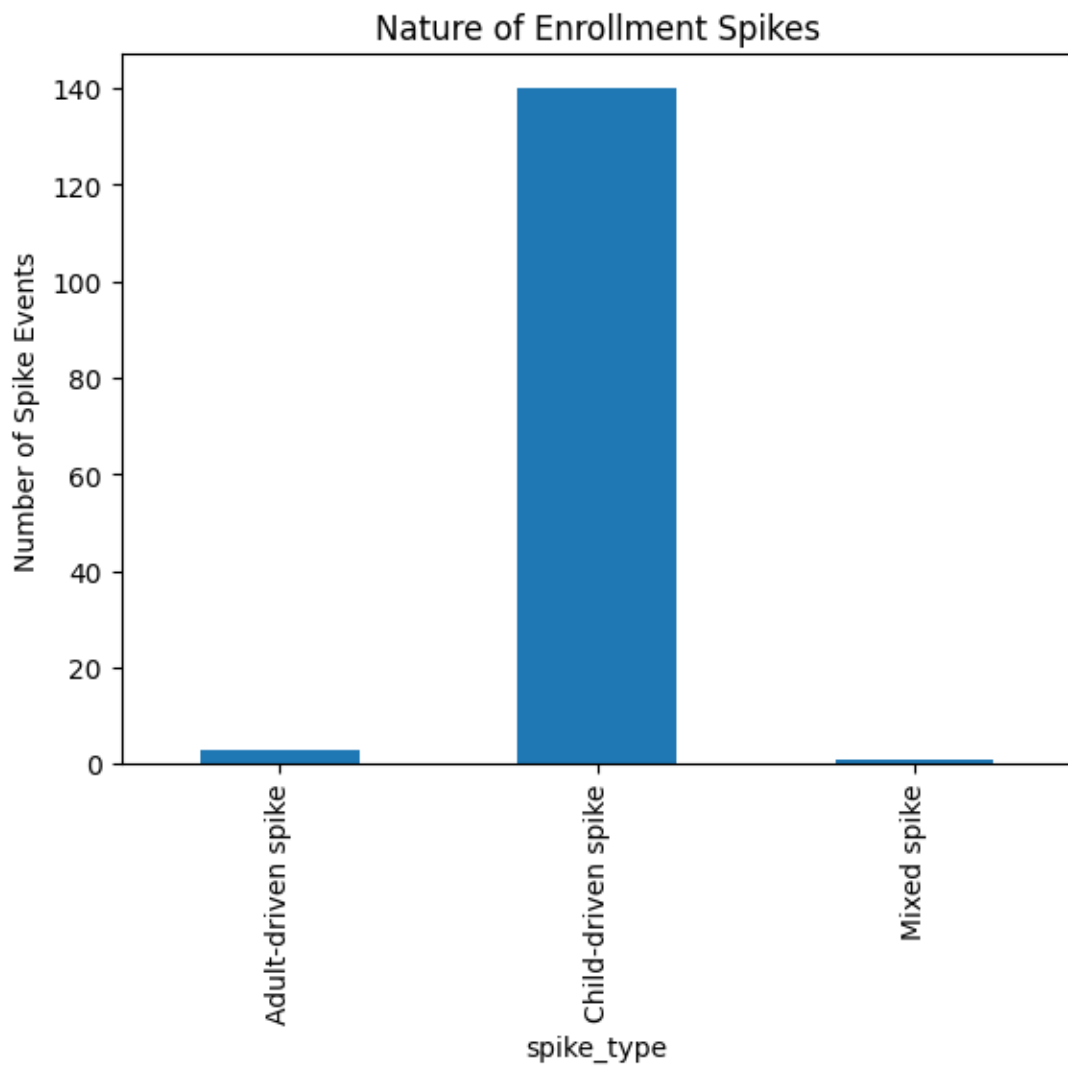
Description:

A district-level visual highlighting extreme EHI values or behavioural outliers, such as:

- Very low enrolment consistency

- Excessive deadline-driven activity
- Severe imbalance between child and adult enrolment





Why this matters:

- Makes administrative stress patterns visible
- Identifies districts needing operational review rather than policy change
- Supports district-specific corrective action

Policy Use:

Enables:

- Targeted district audits
- Field-level interventions
- Performance support rather than punitive evaluation

4.2.5 Time Trend Visualisation**Governance Question:**

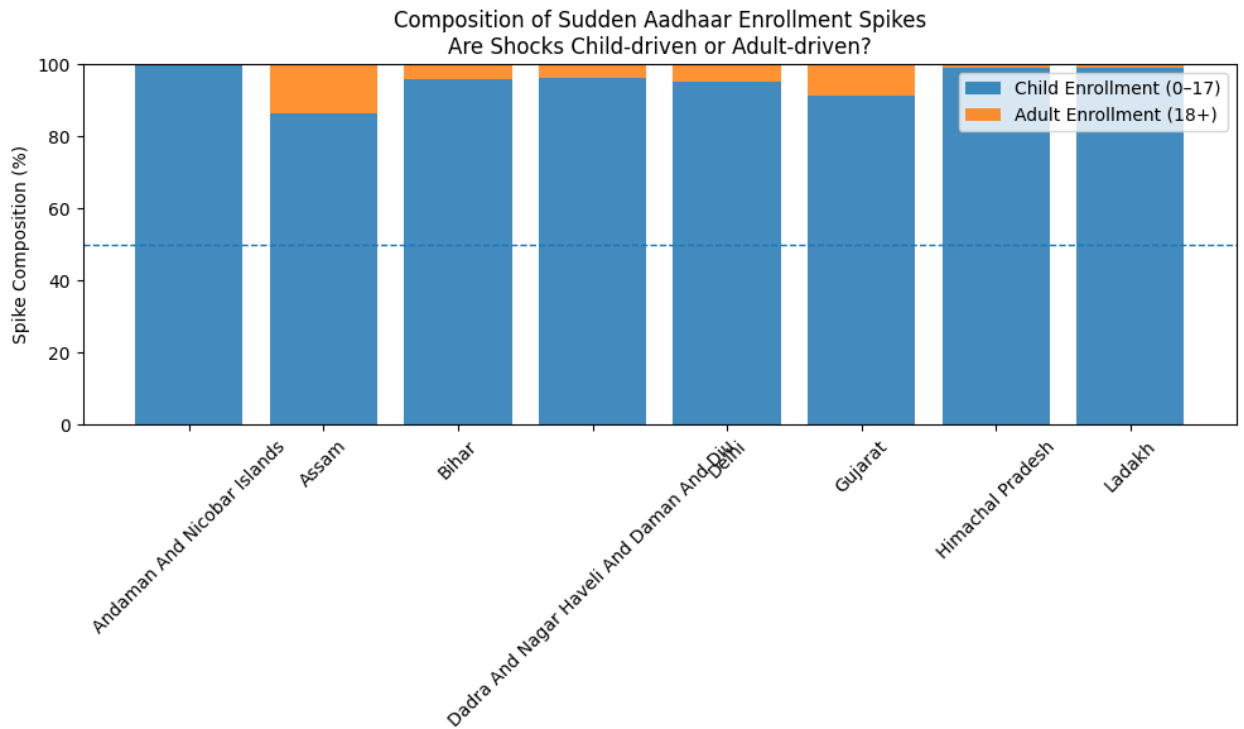
Is enrolment system health improving, stagnating, or deteriorating over time?

Description:

A simple temporal trend showing enrolment volume or behavioural indicators across months or quarters.

Why this matters:

- Captures momentum rather than static snapshots
- Helps distinguish temporary shocks from structural issues



Policy Use:

Useful for:

- Monitoring effects of policy changes
- Evaluating intervention impact
- Forward planning and early warning systems

Executive Policy Brief

India Enrollment Health Index (EHI)

Purpose: Provide a governance-ready diagnostic of Aadhaar enrollment quality across India, moving beyond volume-based metrics to assess systemic health, equity, and timeliness.

1. Background & Problem Statement

Current enrollment monitoring primarily emphasizes **absolute volumes** (total Aadhaar issued). While useful, volume-only metrics mask critical risks:

- Late-life enrollment dominating early-life registration
- Deadline-driven enrollment spikes distorting performance
- Geographic concentration excluding remote populations

These patterns can create a **false signal of success** while weakening long-term identity coverage and service delivery.

2. The Enrollment Health Index (EHI)

The **Enrollment Health Index (EHI)** is a composite governance indicator (0–100) that evaluates enrollment quality across districts and states using three policy-relevant dimensions:

A. Early-Life Coverage Balance (50%)

- Measures balance between child (0–5) and adult (18+) enrollment
- Penalizes reactive, late-stage registration

B. Operational Consistency (30%)

- Measures steady enrollment versus month-end/quarter-end surges
- Flags deadline-driven administrative behavior

C. Access Equity (20%)

- Measures enrollment spread across pincodes
- Detects geographic concentration and access gaps

Higher EHI = Healthier, more resilient enrollment systems

3. National Findings (Illustrative)

- States with **high enrollment volumes** do not always score high on EHI
- Several districts show **adult-heavy enrollment patterns**, indicating delayed identity creation
- Enrollment spikes near reporting deadlines correlate with **lower consistency scores**

These findings suggest that enrollment performance must be evaluated as a **system behavior**, not a counting exercise.

4. Policy Classification Framework

EHI Band System Condition		Governance Meaning
>48	Healthy	Proactive, equitable, stable enrollment
38–48	Fragile	Functional but structurally vulnerable
<38	Risk Zone	Reactive, inequitable, deadline-driven

5. Actionable Policy Recommendations

For High-EHI States/Districts

- Maintain current operational models
- Use as **benchmark districts** for replication
- Pilot new Aadhaar-linked services

For Medium-EHI States/Districts

- Deploy mobile enrollment units
- Rebalance operator incentives away from end-period targets
- Increase monitoring of enrollment timing patterns

For Low-EHI States/Districts (Risk Zones)

- Shift KPIs from volume to **continuity & balance metrics**
 - Integrate enrollment at Anganwadi/school entry points
 - Conduct incentive audits to reduce deadline gaming
-

6. National Indicator: India Enrollment Health Score

A national-level indicator can be derived as the **median district EHI**, representing the typical citizen's enrollment experience.

Use Cases:

- National performance dashboard
 - Year-on-year governance improvement tracking
 - Evidence-based resource allocation
-

7. Strategic Value

The EHI framework enables UIDAI and policymakers to:

- Detect hidden enrollment risks early
- Prioritize structural fixes over short-term targets
- Strengthen Aadhaar as a foundational public good

Key Shift: From *How many enrollments were done?* → *How healthy is the enrollment system?*

8. Conclusion

The Enrollment Health Index transforms operational data into a **decision-grade governance tool**. Adoption of EHI-aligned monitoring can significantly improve equity, timeliness, and resilience of Aadhaar enrollment nationwide.

Prepared for policy review and administrative decision-making.

