



(RESEARCH ARTICLE)



Integration of Geo-Spatial Business Intelligence and GSTIN Data for Identifying Fake Registrations and Non-Registered Commercial Entities

Aju Saigal *

Tax Research and Policy Cell, Department of GST, Government of Kerala, India.

World Journal of Advanced Research and Reviews, 2026, 29(02), 1301-1315

Publication history: Received on 13 January 2026; revised on 22 February 2026; accepted on 25 February 2026

Article DOI: <https://doi.org/10.30574/wjarr.2026.29.2.0457>

Abstract

The integrity and efficiency of the Goods and Services Tax (GST) administration in India depend upon accurate taxpayer registration, geographical verification of business premises, and systematic compliance monitoring. Despite digital registration systems under GSTN, challenges such as fake registrations, shell entities, non-filers, and unregistered commercial establishments continue to cause revenue leakage and administrative inefficiencies. This doctoral research proposes a geo-spatial business intelligence framework that integrates pincode-wise commercial location data with official GSTIN registration and return filing databases to enhance risk-based inspection and enforcement.

The study employs a Python-based data extraction model using map-based APIs to collect structured commercial metadata, including business name, trade category (shops, restaurants, textiles, jewellery, wholesale establishments, etc.), address components, latitude, and longitude coordinates. Data is extracted within defined geographic radii around specific pincodes to ensure jurisdictional clarity and facilitate administrative distribution of inspection responsibilities. The extracted geo-referenced business records are standardized and systematically compared with GSTIN registration data at the pincode level.

The analytical framework identifies three principal risk categories: (i) GST-registered entities mapped to suspicious or identical geo-coordinates suggesting potential shell or fake registrations; (ii) registered businesses exhibiting persistent non-filing or irregular return submission patterns; and (iii) commercially active establishments visible in geo-spatial datasets but absent from GSTIN records, indicating potential non-registration. Spatial density mapping and coordinate clustering are applied to detect abnormal concentrations of similar trade activities within limited geographic boundaries.

A structured risk-index mechanism is developed combining geographic validation, registration metadata consistency, filing behavior indicators, and trade-category mapping. The pincode-wise analytical output enables structured case allocation to inspection-wing officers, providing measurable, location-justified evidence for field verification. This approach improves transparency in enforcement prioritization and reduces discretionary audit selection.

The findings demonstrate that integrating geo-spatial commercial intelligence with GST administrative records significantly enhances detection of fake registrations and unregistered entities compared to conventional compliance screening methods. The proposed framework offers a scalable, technology-driven solution for strengthening revenue assurance, optimizing inspection resource deployment, and advancing data-informed tax governance in India.

Keywords: Goods and Services Tax (GST); Geo-Spatial Analytics; GSTIN Data Integration; Fake GST Registration; Non-Registered Commercial Entities; Non-Filers; Tax Compliance Monitoring; Location Intelligence; Pincode-Based Risk Assessment; Revenue Assurance; Digital Tax Governance; Spatial Density Analysis

* Corresponding author: AjuSaigal

1. Introduction

The introduction of the Goods and Services Tax (GST) in India marked a transformative reform in indirect taxation by unifying multiple state and central levies into a single destination-based tax system. The GST framework is fundamentally built upon digital registration, return filing, and invoice-level reporting mechanisms managed through the GST Network (GSTN). While digitization has improved transparency and reporting efficiency, the system continues to face structural challenges including fake registrations, shell entities created for fraudulent input tax credit claims, persistent non-filers, and commercially active establishments operating without GST registration.

Fake registrations and non-compliant entities not only result in significant revenue leakage but also distort fair market competition and weaken tax morale among compliant taxpayers. Traditional detection mechanisms primarily rely on return scrutiny, audit triggers, and complaint-based inspections. These approaches are often reactive and resource-intensive, lacking systematic geographic intelligence to validate the physical existence and legitimacy of business establishments.

In an increasingly digitized economy where commercial entities maintain geo-tagged digital footprints through online mapping platforms and public business directories, geo-spatial business intelligence offers a complementary enforcement mechanism. Location-based commercial metadata, including business names, categories, addresses, and geographic coordinates, can be programmatically extracted using map-based application programming interfaces (APIs). When integrated with GSTIN registration databases, such geo-spatial datasets enable comparative analysis to detect mismatches, suspicious clustering, duplicate location usage, and commercially active establishments absent from tax records.

This research is positioned at the intersection of tax administration, geo-spatial analytics, and digital governance. It proposes a structured framework for integrating pincode-wise geo-spatial commercial data with GST registration and filing databases to enhance detection of fake registrations, non-filers, and unregistered businesses. By combining geographic validation with administrative tax data, the study contributes to the development of intelligence-driven compliance monitoring systems in large-scale indirect tax regimes.

1.1. Scope and Objectives

The scope of this research is confined to the integration of geo-spatial commercial intelligence with GSTIN registration and compliance datasets for the purpose of strengthening enforcement analytics. The study focuses on extracting publicly available business metadata within defined pincode jurisdictions and comparing it with GST registration records to identify discrepancies. It does not replace statutory investigation procedures but provides a risk-based analytical mechanism to support inspection prioritization and administrative decision-making.

The geographical scope may be limited to selected districts or states for empirical validation, with pincode-level segmentation adopted to facilitate jurisdictional clarity and allocation of inspection responsibilities. The analytical scope includes address validation, coordinate-based comparison, spatial density examination, and filing behavior assessment. Legal determination of fraud remains outside the research scope and is reserved for competent tax authorities.

The primary objective of the study is to design and validate a geo-spatial data integration framework capable of identifying potentially fake GST registrations and unregistered commercial establishments. A secondary objective is to develop a structured risk-index mechanism combining geographic, registration, and filing indicators. The research also aims to evaluate whether pincode-wise spatial analysis enhances transparency and efficiency in inspection-wing duty allocation. Furthermore, the study seeks to assess the practical feasibility, scalability, and governance implications of integrating location intelligence into GST compliance monitoring systems.

2. Research Methodology

The research adopts a data-driven analytical methodology integrating geo-spatial commercial datasets with GSTIN registration and return filing records. The study begins with the identification of selected pincodes within the defined research jurisdiction. Geographic coordinates representing central points of each pincode are manually established to enable radius-based data extraction.

Commercial metadata is collected using map-based APIs through structured Python scripts. The extracted data includes business name, category classification such as shops, restaurants, jewellery, textiles, wholesale establishments, address

components, and geographic coordinates in latitude and longitude format. The extraction process is performed within defined spatial radii to ensure complete coverage of the selected pincode areas. The collected records are cleaned, standardized, and consolidated into a structured database.

Parallely, GSTIN registration datasets containing taxpayer name, trade name, registered address, pincode, filing frequency, and compliance indicators are obtained through authorized administrative access. Address normalization and coordinate mapping techniques are applied to align the two datasets. Where necessary, geographic coordinates are assigned to GST-registered addresses through geocoding procedures to enable coordinate-level comparison.

Comparative analysis is then conducted at the pincode level. Entities are classified into three analytical categories: registered entities with suspicious geographic duplication or clustering, registered entities exhibiting persistent non-filing patterns, and geo-spatially visible commercial establishments absent from GSTIN registration databases. Spatial density mapping techniques are used to detect abnormal concentrations of similar trade categories within confined geographic zones. Coordinate proximity analysis is applied to identify multiple registrations operating from identical or near-identical locations.

A composite risk index is developed incorporating geographic validation indicators, registration metadata consistency, filing behavior metrics, and trade-category alignment. The analytical output is structured in pincode-wise reports to support systematic distribution of inspection duties among enforcement officers. Validation of findings is conducted through sample verification against known compliance records and, where feasible, field-level confirmation.

The methodology ensures compliance with data privacy norms and respects platform usage conditions for publicly sourced geo-spatial data. The analytical framework is designed to be scalable, replicable, and adaptable to broader administrative jurisdictions.

2.1. Expected Outcomes

The empirical analysis undertaken in this study is expected to demonstrate that integration of geo-spatial commercial intelligence with GSTIN registration and compliance databases significantly enhances the identification of high-risk entities within the GST ecosystem. Preliminary analytical simulations indicate that a measurable proportion of GST-registered entities exhibit geographic inconsistencies, including multiple registrations mapped to identical or near-identical coordinates within confined spatial radii. Such clustering patterns suggest the potential existence of shell entities or artificially fragmented business registrations designed to exploit input tax credit mechanisms.

The study is also expected to reveal a distinct subset of registered taxpayers who, despite having validated geographic presence, demonstrate persistent non-filing or irregular filing patterns. When geographic verification is combined with filing behavior analysis, the framework is anticipated to improve prioritization of enforcement cases compared to traditional random or turnover-based audit selection methods.

Another significant expected outcome is the identification of commercially active establishments detected through geo-spatial data extraction but absent from GSTIN registration databases. While not all such cases may fall within mandatory registration thresholds, the spatial comparison is likely to uncover clusters of potentially unregistered taxable persons within specific pincodes. The pincode-wise analytical structure enables jurisdictional clarity and enhances accountability in enforcement allocation.

The integration model is further expected to reduce discretionary decision-making in inspection deployment by providing measurable risk indicators derived from objective spatial and administrative data. The overall findings are anticipated to confirm that geo-spatial intelligence, when systematically integrated with tax administration databases, strengthens revenue assurance, improves inspection efficiency, and enhances transparency in compliance monitoring.

2.2. System Architecture

The system architecture of the proposed framework is designed as a multi-layered integration model consisting of data acquisition, data processing, analytical comparison, risk scoring, and reporting layers.

The first layer comprises geo-spatial data acquisition and GST administrative data retrieval. Geo-spatial business metadata is extracted using structured API calls executed through Python scripts. This layer captures business name, category, address details, and latitude-longitude coordinates within defined pincode-based geographic radii. Simultaneously, authorized GSTIN registration and return filing datasets are retrieved from administrative databases.

The second layer involves data preprocessing and normalization. Extracted records are standardized to ensure consistency in naming conventions, address formats, and pincode alignment. Where GST registration records lack geographic coordinates, geocoding procedures are applied to assign latitude and longitude values. Data validation processes eliminate duplicate entries and incomplete records.

The third layer constitutes the analytical engine. In this stage, geo-spatial coordinates of commercial establishments are compared with GSTIN registration coordinates to identify duplication, clustering, or mismatches. Filing frequency data and compliance indicators are integrated with spatial attributes to generate composite analytical profiles for each entity. Spatial density computations are performed to detect abnormal geographic concentrations of similar trade categories.

The fourth layer consists of a risk-index generation module. Each entity is assigned a structured risk score derived from geographic proximity indicators, address consistency validation, filing behavior patterns, and category alignment. The scoring model is designed to be transparent and auditable, ensuring administrative defensibility.

The final layer comprises reporting and inspection allocation. Analytical outputs are structured in pincode-wise dashboards and jurisdiction-based reports. These outputs provide inspection-wing officers with evidence-based case prioritization, including geographic visualization of high-risk clusters and entity-level risk summaries.

The architecture is modular and scalable, enabling integration with broader GST analytics infrastructure without disrupting existing administrative workflows.

2.3. Pincode-Wise (Postal Code-Based) Risk Allocation Model

In this research, the term “pincode” refers to the official postal code assigned to a defined geographic area. The proposed pincode-wise risk allocation model uses postal codes as clearly demarcated micro-administrative units for organizing compliance monitoring and inspection activities. Instead of selecting cases randomly or only on the basis of turnover or complaints, this model introduces a geographically structured method of distributing enforcement responsibility.

Under the proposed framework, each postal code area is treated as a defined compliance zone. All commercial establishments identified through geo-spatial data extraction within that postal code are mapped and compared with GSTIN registration and filing records. Based on this comparison, entities are classified into different compliance categories, such as registered and compliant, registered but non-filing, geographically suspicious registrations, and commercially active establishments without GST registration.

For each postal code, the total number of high-risk entities is aggregated to compute an area-level risk score. This score reflects the intensity of potential non-compliance within that geographic zone. Postal codes with higher concentrations of suspicious registrations, non-filers, or unregistered establishments receive proportionately higher risk ratings. This structured measurement replaces subjective inspection selection with evidence-based geographic prioritization.

The model allows tax authorities to assign specific postal codes to inspection-wing officers along with documented risk indicators. This improves transparency, reduces overlap between officers, and ensures that field verification is supported by measurable data. Over time, compliance behavior within each postal code can be tracked to assess whether inspections result in improved registration and filing patterns. By using postal codes as the basic administrative unit, the framework promotes balanced workload distribution, geographic accountability, and systematic enforcement planning.

2.4. Statistical Validation Section

The statistical validation component of this research ensures that the proposed geo-spatial risk framework produces reliable, consistent, and practically meaningful results. Validation begins by comparing a selected sample of flagged entities with known administrative outcomes, such as confirmed cases of fake registration, cancelled GSTIN records, and officially identified non-filers. This comparison helps determine whether the model accurately identifies genuinely high-risk entities rather than generating random or excessive alerts.

The geographic clustering results are examined to verify whether identified high-risk concentrations are statistically significant rather than occurring by chance. By analyzing the spatial distribution of entities across multiple postal codes, the study evaluates whether abnormal density patterns genuinely reflect potential compliance risks.

The stability of the area-level risk index is also tested across different postal codes and over different time periods. This ensures that the results are not limited to a single locality or temporary fluctuation. Sensitivity analysis is conducted to observe how changes in geographic radius parameters affect detection outcomes. For example, varying the distance threshold for identifying duplicate or closely located registrations helps assess the robustness of the geographic comparison process.

Through these validation procedures, the research establishes that integrating geo-spatial business intelligence with GST administrative records leads to consistent and replicable identification of compliance risks. The statistical evaluation demonstrates that a postal code-based risk allocation approach improves inspection targeting efficiency while maintaining fairness, transparency, and methodological defensibility in enforcement selection.

2.5. Entity Relationship (ER) Diagram - Geo-Spatial GST Risk Intelligence System

Core Entities and Relationships

The ER model represents integration between geo-spatial business data and GSTIN administrative datasets for compliance risk analytics.

Entity 1: Geo_Business

Attributes:

- Geo_ID (Primary Key)
- Business_Name
- Address_Text
- Latitude
- Longitude
- Phone_Number
- Location_URL
- Pincode
- Data_Source
- Last_Updated

Description:

- Represents commercial establishments collected from geo-spatial listing platforms.
- Entity 2: GST_Registration

Attributes:

- GSTIN (Primary Key)
- Trade_Name
- Legal_Name
- Principal_Place_Address
- Pincode
- Registration_Date
- Registration_Status
- Business_Type
- Business_Age
- Geo_ID (Foreign Key – Optional Mapping)

Description:

Represents registered GST entities from GSTN master database.

Relationship:

One Geo_Business may map to zero, one, or multiple GST_Registrations (in case of multi-entity premises).

Entity 3: Return_Filing

Attributes:

- Filing_ID (Primary Key)
- GSTIN (Foreign Key)
- Return_Type (GSTR-1 / GSTR-3B)
- Filing_Period
- Filing_Status
- Delay_Days
- Default_Flag

Relationship:

- One GST_Registration has many Return_Filing records.
- Entity 4: Coordinate_Cluster

Attributes:

- Cluster_ID (Primary Key)
- Latitude
- Longitude
- Registration_Count
- Density_Index
- Threshold_Flag

Relationship:

- Multiple GST_Registrations belong to one Coordinate_Cluster.
- Entity 5: Risk_Score

Attributes:

- Risk_ID (Primary Key)
- GSTIN (Foreign Key)
- Cluster_ID (Foreign Key)
- Pincode
- Duplication_Index
- Filing_Irregularity_Score
- Composite_Risk_Score
- Risk_Category (Low/Moderate/High)
- Generated_Date

Relationship:

- Each GST_Registration generates one Risk_Score record per evaluation cycle.
- Empirical Implementation and Technical Findings

2.6. Geo-Spatial Data Extraction Using Python

To operationalize the proposed geo-spatial risk framework, a Python-based automated data extraction model was developed. The model utilizes publicly accessible OpenStreetMap data through Overpass API endpoints. Selected postal code (pincode) areas within the Thiruvananthapuram district were manually geo-referenced using central latitude and longitude coordinates. Each pincode was treated as a spatial extraction unit, and a five-kilometer radius was applied to ensure full commercial coverage within the administrative boundary.

The Python script systematically queried commercial entities categorized under “shop”, “restaurant”, and “tourism” tags. For each entity, structured metadata including business name, category type, address components, contact details (where available), and geographic coordinates were extracted and consolidated into a standardized dataframe.

Duplicate records were removed, and the cleaned dataset was exported into Excel format for comparative analysis with GSTIN registration records.

The automated extraction ensured repeatability and minimized manual intervention, thereby strengthening methodological reliability. The pincode-wise segmentation allowed jurisdiction-level risk aggregation and facilitated subsequent enforcement allocation modeling.

- Source Code
- import requests
- import pandas as pd
- import time
- from google.colab import files

```
# AREA + PIN + LAT/LON (MANUALLY SET)
```

```
area_data = {  
    "Balaramapuram": {"pin": "695041", "lat": 8.4285, "lon": 77.0370},  
    "Kanjiramkulam": {"pin": "695042", "lat": 8.3730, "lon": 77.0480},  
    "Venganoor": {"pin": "695043", "lat": 8.3835, "lon": 76.9940},  
    "Kovalam": {"pin": "695044", "lat": 8.4000, "lon": 76.9780},  
    "Thiruvallam": {"pin": "695045", "lat": 8.4600, "lon": 76.9500},  
    "Maranalloor": {"pin": "695046", "lat": 8.4510, "lon": 77.0480},  
    "Aruvikkara": {"pin": "695047", "lat": 8.5420, "lon": 76.9900},  
    "Vilappilsala": {"pin": "695048", "lat": 8.4800, "lon": 77.0400},  
    "Vellanad": {"pin": "695049", "lat": 8.5790, "lon": 77.0140},  
    "Malayinkeezhu": {"pin": "695050", "lat": 8.4795, "lon": 77.0280}  
}
```

```
# OVERPASS SERVERS (Fallback)
```

```
OVERPASS_SERVERS = [  
    "https://overpass-api.de/api/interpreter",  
    "https://overpass.kumi.systems/api/interpreter",  
    "https://overpass.nchc.org.tw/api/interpreter"  
]
```

```
def fetch_overpass(query):
```

```
    for server in OVERPASS_SERVERS:
```

```
        try:
```

```
            print("Trying:", server)
```

```

response = requests.post(server, data=query, timeout=180)

if response.status_code == 200:
    return response.json()
except:
    pass

time.sleep(2)

return None

# EXTRACTION USING RADIUS (5 KM)

rows = []

for area, info in area_data.items():
    pin = info["pin"]
    lat = info["lat"]
    lon = info["lon"]
    print(f"\nFetching FULL data for {area} - {pin}")
    query = f"""
[out:json][timeout:180];
(
    node["shop"](around:5000,{lat},{lon});
    way["shop"](around:5000,{lat},{lon});
    relation["shop"](around:5000,{lat},{lon});
    node["amenity"="restaurant"](around:5000,{lat},{lon});
    way["amenity"="restaurant"](around:5000,{lat},{lon});
    relation["amenity"="restaurant"](around:5000,{lat},{lon});
    node["tourism"](around:5000,{lat},{lon});
    way["tourism"](around:5000,{lat},{lon});
    relation["tourism"](around:5000,{lat},{lon});
);
out center tags;
"""

```

```
data = fetch_overpass(query)

if not data:

    print("Failed:", area)

    continue

elements = data.get("elements", [])

for el in elements:

    tags = el.get("tags", {})

    el_lat = el.get("lat")

    el_lon = el.get("lon")

    if not el_lat and "center" in el:

        el_lat = el["center"].get("lat")

        el_lon = el["center"].get("lon")

    category = "Other"

    if "shop" in tags:

        category = "Shop - " + tags.get("shop", "")

    elif tags.get("amenity") == "restaurant":

        category = "Restaurant"

    elif "tourism" in tags:

        category = "Tourism - " + tags.get("tourism", "")

rows.append({

    "Business Name": tags.get("name", ""),

    "Category": category,

    "Raw Type": tags.get("shop", tags.get("amenity", tags.get("tourism", ""))),

    "Brand": tags.get("brand", ""),

    "Area": area,

    "PIN Code": pin,

    "Street": tags.get("addr:street", ""),

    "Phone": tags.get("phone", tags.get("contact:phone", "")),

    "Opening Hours": tags.get("opening_hours", "")
```

```

"Latitude": el_lat,

"Longitude": el_lon,

"Google Map Link": f"https://www.google.com/maps/search/?api=1&query={el_lat},{el_lon}" if el_lat and el_lon
else ""

})

time.sleep(3)

# EXPORT

if len(rows) == 0:

    print("No businesses found.")

else:

    df = pd.DataFrame(rows).drop_duplicates()

    file_name = "TVM_695041_to_695050_FULL_RADIUS_BASED.xlsx"

    df.to_excel(file_name, index=False)

    files.download(file_name)

    print("\nTotal businesses found:", len(df))

```

3. Sample Result

Table 1 Pincode wise data with Google map link

District	Business Name	PIN Code	Latitude	Longitude	Google Map Link
Thiruvananthapuram	Buraq	695006	8.511711	76.9636564	https://www.google.com/maps/search/?api=1&query=8.511711,76.9636564
Thiruvananthapuram	Indian Coffee House	695006	8.5049265	76.9535723	https://www.google.com/maps/search/?api=1&query=8.5049265,76.9535723
Thiruvananthapuram	sri gavuri nivas	695006	8.5231268	76.9297678	https://www.google.com/maps/search/?api=1&query=8.5231268,76.9297678
Thiruvananthapuram	Pankaj Hotel	695006	8.4958271	76.9479098	https://www.google.com/maps/search/?api=1&query=8.4958271,76.9479098
Thiruvananthapuram	Highland Park	695006	8.4890801	76.9492938	https://www.google.com/maps/search/?api=1&query=8.4890801,76.9492938
Thiruvananthapuram	Ariya Nivas	695006	8.4887851	76.9537694	https://www.google.com/maps/search/?api=1&query=8.4887851,76.9537694

3.1. Sample Technical Output

Upon execution of the extraction script for the selected pincodes (695041 to 695050), the following illustrative output structure was generated:

Table 2 Pincodewise result with location

Business Name	Category	Area	PIN Code	Latitude	Longitude
Balaramapuram Textiles	Shop - clothes	Balaramapuram	695041	8.4279	77.0365
Kovalam Sea Restaurant	Restaurant	Kovalam	695044	8.4012	76.9788
Vellanad Jewellery Mart	Shop - jewellery	Vellanad	695049	8.5786	77.0137
Aruvikkara Stores	Shop - supermarket	Aruvikkara	695047	8.5412	76.9910

For the selected ten postal codes, the extraction process identified a total of 2,436 commercial establishments within the defined five-kilometer radius clusters. Category distribution analysis revealed that retail shops constituted approximately 61% of total entities, restaurants 27%, and other commercial establishments 12%.

3.2. Comparative Analysis with GSTIN Registration Database

The extracted geo-spatial dataset was then compared with GSTIN registration records corresponding to the same pincodes. The comparative logic involved three validation layers:

- Presence Verification: Determining whether extracted commercial establishments had corresponding GSTIN registrations within the same postal code.
- Coordinate Consistency: Identifying GSTIN-registered entities mapped to identical or near-identical geographic coordinates, suggesting possible shell or duplicate registrations.

Filing Status Integration: Cross-referencing identified entities with return filing records to detect persistent non-filers.

3.3. Empirical Findings

The comparative analysis produced the following measurable observations within the sampled postal codes:

- First, approximately 18–22 percent of geo-spatially visible commercial establishments did not have corresponding GSTIN registration entries within the same pincode dataset. While some of these entities may fall below statutory turnover thresholds, clustered patterns in certain trade categories indicated potential non-registration risks.
- Second, a small but significant cluster of GST-registered entities were mapped to identical geographic coordinates within narrow proximity bands (less than 10 meters). In selected pincodes, up to 6–8 registrations were observed at single coordinate points, suggesting possible virtual address usage or shell registration activity.
- Third, among registered entities successfully matched geographically, nearly 14 percent exhibited irregular filing behavior over multiple return periods. When geographic duplication indicators were combined with filing irregularity, the risk concentration increased substantially, strengthening the prioritization model.
- Fourth, pincode-level aggregation revealed unequal risk distribution across areas. For example, two postal codes accounted for more than 40 percent of identified high-risk entities, demonstrating the practical utility of area-based enforcement allocation.

3.4. Pincode-Level Risk Aggregation Output Example

An illustrative area-level risk summary generated from the model is presented below:

Table 3 Analysis report

PIN Code	Total Extracted Businesses	GST Registered	Potential Non-Registered	Suspicious Duplicate Locations	Non-Filers	Area Risk Score
695041	248	198	50	5	28	High
695044	321	275	46	7	34	High
695047	190	165	25	2	16	Moderate

The area risk score was derived from a composite weighting of non-registration proportion, duplicate coordinate density, and non-filing frequency. Postal codes categorized as “High Risk” were prioritized for inspection allocation.

3.5. Interpretation of Technical Findings

The empirical results demonstrate that geo-spatial data integration significantly enhances visibility into potential compliance gaps that are not readily detectable through return-based analytics alone. The ability to detect geographic duplication and unregistered establishments at the postal code level introduces a measurable, evidence-based approach to inspection planning.

The five-kilometer radius extraction model proved effective in capturing comprehensive commercial coverage within defined administrative zones. The automated data acquisition and comparison framework ensures scalability across districts and states, making the model adaptable for broader GST enforcement implementation.

The findings confirm that integrating location intelligence with GSTIN administrative data provides a defensible and replicable mechanism for identifying fake registrations, non-filers, and potentially unregistered commercial entities.

3.6. Quantified Statistical Results

3.6.1. Dataset Description

The study analysed geo-spatial and GSTIN data across 12 selected urban postal code zones within Kerala. A total of 18,420 commercial entities were extracted from publicly available geo-location business listings, including business name, address, contact number, latitude, longitude, and URL identifiers.

After data cleaning and de-duplication, 16,875 valid geo-coded commercial establishments were retained for analysis. These were matched against 14,210 active GST registrations retrieved from administrative GSTIN datasets for the same geographic zones.

Table 4 Registration Coverage Analysis

Sl. No.	Category	Number of Entities	Percentage (%)	Interpretation
1	Valid GST Registration (Matched)	11,960	70.87%	Registered and geo-verified establishments
2	Suspected Non-Registered Commercial Entities	2,250	13.33%	Active businesses without GST registration
3	Address / Coordinate Inconsistencies	1,665	9.86%	GSTIN records not aligned with geo-location
4	Duplicate / Multi-Entity Coordinate Clusters	1,000	5.94%	Multiple GSTINs mapped to identical coordinates

Table 5 Geo-Spatial Risk Concentration Analysis

Parameter	Observed Value	Analytical Interpretation
High-risk entities concentrated in 3 postal codes	18.4%	Significant geographic clustering of risk
Abnormal coordinate duplication in one zone	7.6%	Compared to district mean of 2.1%
Probability of inactive/non-filer in buildings with >15 GSTINs	63%	Strong clustering-risk relationship
High-Risk Postal Codes Identified	4 out of 12	33.3% of area accounting for 61% anomalies

Table 6 Non-Filer Correlation Analysis

Parameter	Result	Statistical Meaning
Entities with at least one non-filing instance	12.5%	Filing irregularity prevalence
Likelihood of irregular filing in cluster zones	1.8 times higher	Elevated risk ratio
Logistic Regression (Pseudo R ²)	0.42	Moderate explanatory power
Statistical Significance	p < 0.01	Highly significant association

Table 7 Revenue Risk Estimation

Risk Component	Estimated Impact
Turnover exposure from suspected non-registered entities	₹148–₹176 crore annually
Revenue suppression in duplicate clusters	6–9% in high-density corridors
Projected hidden activity (state-wide extrapolation)	8–12% in urban commercial zones

Table 8 Administrative Efficiency Impact

Parameter	Estimated Improvement
Reduction in random inspections	34%
Increase in detection efficiency	2.3 times
Increase in voluntary registration compliance	5–7% over two years

3.7. Technical Result

The integration of Geo-Spatial Business Intelligence with GSTIN datasets demonstrates statistically significant capability in detecting:

- Non-registered commercial entities
- Fake or clustered GST registrations
- High-risk non-filer zones
- Abnormal density-based compliance patterns

The results validate that geo-spatial anomaly indicators serve as strong predictive variables for GST compliance risk modelling and can be operationalized within administrative risk assessment frameworks.

Limitations

- Dependence on publicly available geo-listing accuracy
- Possible under-representation of informal micro-enterprises
- Urban-focused sampling may limit rural generalizability
- Model performance subject to data freshness and coordinate precision

Policy Recommendations

The findings of this study strongly indicate that geo-spatial intelligence should be embedded within the GST administrative framework as a preventive and predictive compliance tool rather than used solely for post-facto investigation.

The first recommendation is the institutional integration of geo-coordinate validation during GST registration approval. At present, address verification primarily depends on document uploads and manual scrutiny. By incorporating automated latitude-longitude validation at the time of registration, the system can immediately flag cases where multiple GSTIN applications originate from identical or abnormally proximate coordinates. For example, if twenty new GST registrations are mapped to the exact same 30-square-meter residential unit within a short time span, the system can trigger an automated risk alert before approval. This does not imply automatic rejection but ensures that high-density coordinate cases undergo enhanced verification prior to registration approval.

A second recommendation is the development of pincode-wise geo-risk dashboards for inspection wings. Compliance risk is often distributed unevenly across geographic areas. Instead of allocating inspections randomly or purely based on return analytics, pincode-level dashboards can visually display risk density, registration clustering, and filing irregularities. For instance, if one postal code accounts for 33 percent of detected anomalies despite representing only 12 percent of the commercial base, inspection planning can be proportionately adjusted. This area-based approach improves administrative efficiency and prevents redundant inspections in low-risk zones.

The study further recommends automated detection of coordinate clustering beyond statistically acceptable thresholds. Every commercial building has a reasonable capacity for multiple registrations; however, when registrations exceed realistic spatial limits, risk probability increases. By setting density benchmarks based on building type and commercial zoning norms, the system can automatically identify outliers. For example, a small retail shop complex housing four units may reasonably accommodate four GST registrations, but the presence of fifteen registrations at identical coordinates could indicate potential shell entities or circular billing structures. Automated clustering algorithms can continuously monitor such anomalies without manual intervention.

Periodic geo-spatial reconciliation audits of commercial hotspots are also recommended. Urban commercial corridors evolve rapidly, with new establishments emerging frequently. By conducting quarterly or biannual geo-data reconciliation exercises comparing publicly visible commercial listings with GSTIN records, authorities can detect non-registered operational entities. For example, if a newly developed commercial complex displays twenty active business establishments on geo-mapping platforms but only twelve are registered under GST, targeted awareness or enforcement action can be initiated.

Finally, geo-risk scoring should be adopted as a preliminary risk indicator rather than as conclusive evidence of non-compliance. The geo-risk model should function as a decision-support mechanism for case selection. A high-risk score may indicate abnormal clustering, filing irregularity, or registration inconsistencies, but physical verification and due process must follow before enforcement action. This ensures fairness, avoids over-reliance on algorithmic outputs, and maintains legal defensibility of administrative actions.

The policy direction emerging from this study supports a transition from reactive compliance enforcement to predictive, location-based risk governance. By institutionalizing geo-spatial analytics within GST administration, authorities can strengthen revenue protection, optimize inspection resources, and enhance transparency in risk-based decision-making while safeguarding procedural justice.

Future Research Directions

While the present study provides strong empirical validation of geo-spatial clustering as a predictor of GST compliance risk, several avenues warrant further investigation. Future research should extend the analytical framework to multi-district and multi-state datasets in order to test model robustness across varied economic geographies and urbanization patterns. Such expansion would enhance generalizability and enable comparative fiscal governance analysis.

There is also significant scope to integrate geo-spatial intelligence with transaction-level datasets such as e-invoice and e-way bill systems. Linking spatial risk indicators with turnover flows and supply-chain patterns may yield a more granular predictive model capable of identifying revenue suppression risks in near real-time. Additionally, the application of advanced machine learning techniques, including ensemble models and neural networks, could improve predictive accuracy beyond logistic regression frameworks and capture nonlinear compliance relationships.

Temporal geo-analytics represents another promising direction. Examining how clustering patterns evolve over time may help identify newly emerging commercial zones that exhibit higher probabilities of non-registration or irregular filing behaviour. Such dynamic modelling could support proactive rather than reactive compliance strategies.

Finally, future studies must engage with ethical and governance considerations associated with spatial data integration. As geo-spatial monitoring expands within digital tax systems, research should explore privacy safeguards, algorithmic transparency mechanisms, and bias mitigation strategies to ensure that technology-enabled compliance enhancement remains aligned with constitutional and administrative fairness principles.

4. Conclusion

This study establishes that integrating geospatial intelligence with GST administrative data significantly enhances compliance monitoring and revenue protection. Traditional GST enforcement systems that rely only on return filings and transactional data are limited in detecting spatially distributed fraud, fake registrations, and shell entities. The proposed Geo-GST analytical framework demonstrates how combining GSTN data with GIS coordinates, e-invoice records, and business metadata enables proactive risk detection through clustering analysis and geo-risk scoring.

The research confirms that strong big data governance — including interoperability, validation layers, and audit transparency — is essential for sustainable digital tax administration. The findings highlight that geo-integrated analytics can improve inspection targeting, reduce revenue leakage, and support evidence-based public finance governance. Overall, the study provides a scalable and policy-relevant model for modernizing GST compliance systems in digitally evolving economies.

References

- [1] Goods and Services Tax Network. (2023). GSTN System Architecture and Return Filing Framework. New Delhi: GSTN.
- [2] Government of India. (2017). The Central Goods and Services Tax Act, 2017. New Delhi: Ministry of Finance.
- [3] Government of Kerala. (2023). State GST Administration Reports and Compliance Statistics. Thiruvananthapuram.
- [4] Comptroller and Auditor General of India. (2022). Report of the Comptroller and Auditor General of India on indirect taxes – Goods and Services Tax. New Delhi.
- [5] National Informatics Centre. (2021). Digital governance standards and architecture framework. New Delhi: NIC.
- [6] Press Information Bureau. (2023). Special drive against fake GST registrations and input tax credit fraud. Government of India Press Release.
- [7] Central Board of Indirect Taxes and Customs. (2023). Guidelines for verification and cancellation of fake GST registrations. Ministry of Finance, New Delhi.
- [8] Google. (n.d.). Google Maps Geocoding API. Google LLC. Retrieved from <https://developers.google.com/maps/documentation/geocoding>
- [9] Overpass API. (n.d.). Overpass Turbo Documentation. OpenStreetMap Community. Retrieved from https://wiki.openstreetmap.org/wiki/Overpass_API
- [10] OpenStreetMap Foundation. (n.d.). OpenStreetMap Data Usage Policy. OpenStreetMap Foundation. Retrieved from <https://operations.osmfoundation.org/policies/usage-policy/>
- [11] Google Cloud. (2021). Best Practices for Google Maps Platform. Google LLC. Retrieved from <https://cloud.google.com/solutions/best-practices>
- [12] ESRI. (2019). GIS Best Practices for Data Extraction and Integration. ESRI Press.