# SYMBIOSIS INSTITUTE OF TECHNOLOGY, NAGPUR

# Analyzing the Impact of Flood Events on Direct Benefit Transfer (DBT) Distribution in India

**Data Science**
**Mini Project Report**

**Submitted By:**

**Shrishti Prasad**
**22070521035 , VII A**

**Submitted To:**
**Dr. Piyush Chauhan**

**Department of Computer Science and Engineering**
**Symbiosis Institute of Technology, Nagpur**

# INDEX

# 1. Introduction

The reason I chose this topic is that I want to see and learn about the process of how DBT works in India, not only at the state level but also at the district level.

The DBT is basically a DIRECT BENEFIT TRANSFER, which means that direct funds come from government like subsidiaries, pension, and from other sources for different age groups of people, like 0- 5 for kids, 11-18, 18-25, and old age people, etc, directly to a person's account without any middlemen to avoid delays and corruption. The NIA takes this initiative. Later than floods, which is one of the problems affecting many lives in the coastal region of India. Many houses, crops, and animals were affected by floods. I want to link these two datasets, DBT and Floods, and see how much DBT came from the government during those flood years in different states. My project differs from existing ones, as previous ones only show data at the national level. However, my project also displays data at the district level for different years, showing the amount of DBT allocated to various districts and states. Do they receive more DBT in those years when floods occurred more frequently, or do they receive more debt when floods occurred rarely? Do natural disasters, such as floods, affect DBT? Can we identify gaps where we are lacking? How to improve it?

So basically my problem statement is to analyze how natural disasters like DBT  affect floods in those years by comparing transaction , flood amount etc. By identifying the trends in affected areas and non-affected areas.

# 2. Literature Review

The Government of India launched the DBT system in 2013. Its objective is to ensure that government subsidies, welfare funds, and financial aid are credited directly into the bank accounts of beneficiaries, reducing middlemen

and, consequently, leaks in the transfers. Various studies have looked at how DBT improves transparency, efficiency, and financial inclusion, particularly in rural areas.

Other scholars remarked that DBT has significantly enhanced efficiency in the distribution of welfare at a faster and more precise scale. They also stressed that the success of DBT relies on the corresponding infrastructure: internet connectivity, banking access, and good local governance.

Delays, which often characterize these paper-based approaches, are reduced and transparency enhanced.

Yet, few studies have quantitatively analyzed how flood intensity impacts the volume of DBT transfers across different states. This creates a research gap. There is also limited evidence on whether DBT payments increase during flood years and how relief allocation correlates with the severity of disasters.

This project tries to fill this gap by linking the flood data with the DBT performance data. In this project, exploratory data analysis, regression, and clustering will be performed to understand:

- Whether DBT acts as a responsive welfare mechanism during floods, and

Which states have stronger associations between flood events and DBT fund transfers. This analysis forms part of the growing area of data-driven public policy and disaster management. It gives insight into how welfare systems based on technology can make timely relief distribution possible.

## 3. Abstract

DBT is one of the major initiatives of the Government of India to ensure efficient and effective delivery of various subsidies and welfare benefit funds directly into the accounts of the citizens. In this respect, the effectiveness of DBT could be influenced by meteorological events like flooding. Flooding generally disrupts livelihood and creates immediate needs for financial assistance. This project examines the flood event impact on the flow of DBT funds across Indian states from 2017 to 2021.

The work involved cleaning and combining two datasets: district-wise DBT records and state-wise flood relief data. We analyzed them using various techniques from data science comprising EDA, Regression, Classification, and Clustering. In the study, we investigated if heavy floods in any state are also associated with a rise in DBT transfers in that period.

The result indicates that the flood intensity and DBT disbursement are positively related. This implies that in disaster years, DBT acts as a useful welfare tool.

## 4. Methodology

1. Description of the Approach/ Model/ System Implemented

It involves a data-driven analytical approach in understanding how flood events shape the distribution of DBT across Indian states. The methodology involves the integration of EDA, statistical modeling, and machine learning to identify trends, correlations, and clustering among states based on the severity of flooding and DBT activity.

Overall, the workflow includes:

Data: Collected data of Direct Benefit Transfer and flood dataset from kaggle and india data portal.

Data Cleaning & Preprocessing : Remove missing values , null values , duplicate values and organize the data by year.

Merging of Data: DBT and floods dataset are merged on common parameters like state and year

Feature Engineering, creating new analytical columns such as year-over-year DBT growth (dbt_change) and flood indicators (flood_flag).
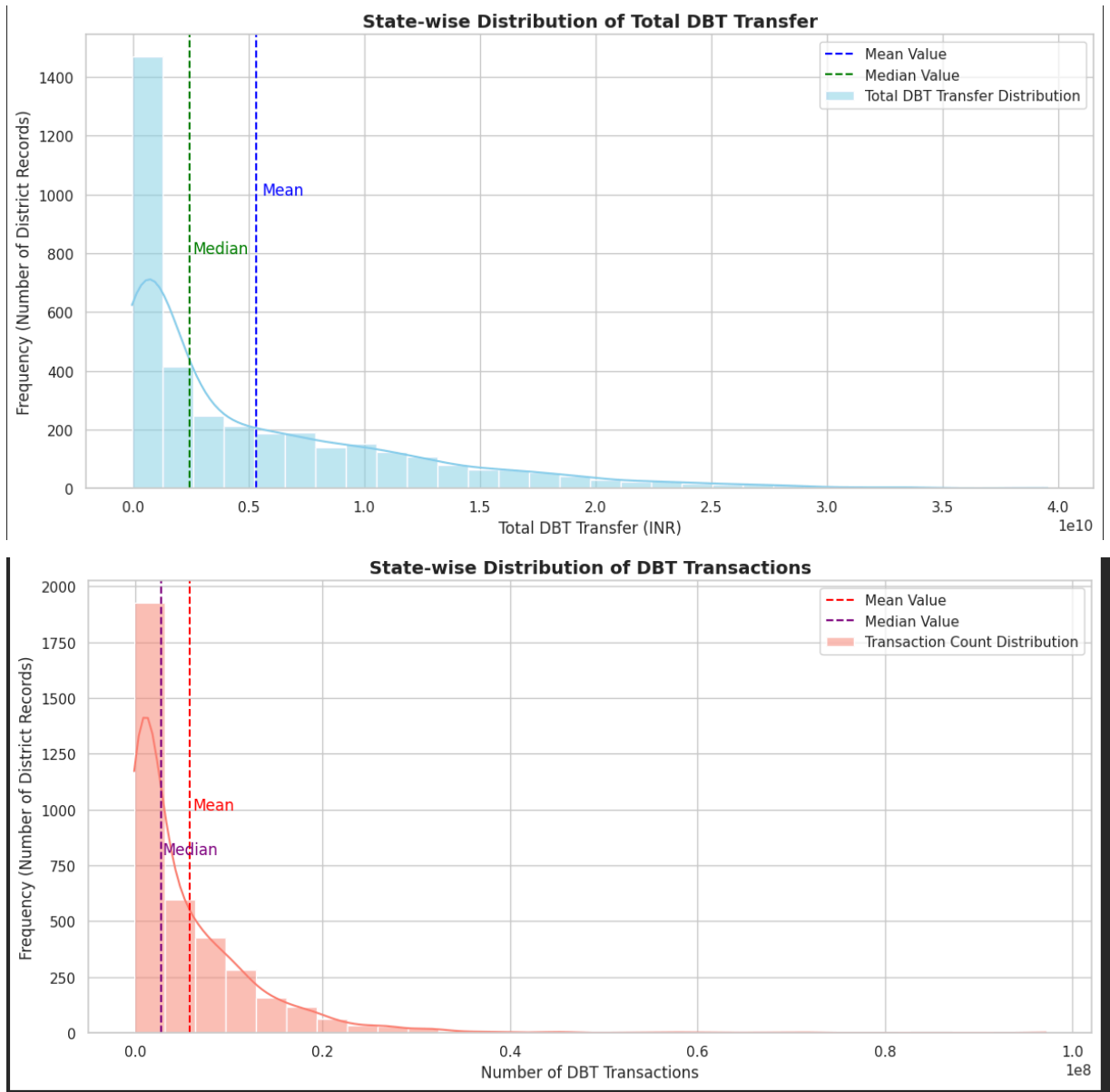
EDA; Visualization of distributions, correlations, and year-wise variations.

Modeling & Analysis: Use regression analysis , classification and clustering to get insights.

Visualization: Make a streamlit dashboard to get interactive results.

| Category | Tools & Libraries Used |
|---|---|
| Data Handling | Pandas , numpy |
| Visualization | matplotlib,seaborn,plotly,express |
| Statistical Modeling | Statsmodels.api, scikit-learn |
| Machine learning | Kmeans, StandardScaler |
| Dashboard | streamlit |
| Deployment | Cloudflared , ngrok, localtunnel |

# Implementation



**State-wise Distribution of Total DBT Transfer**
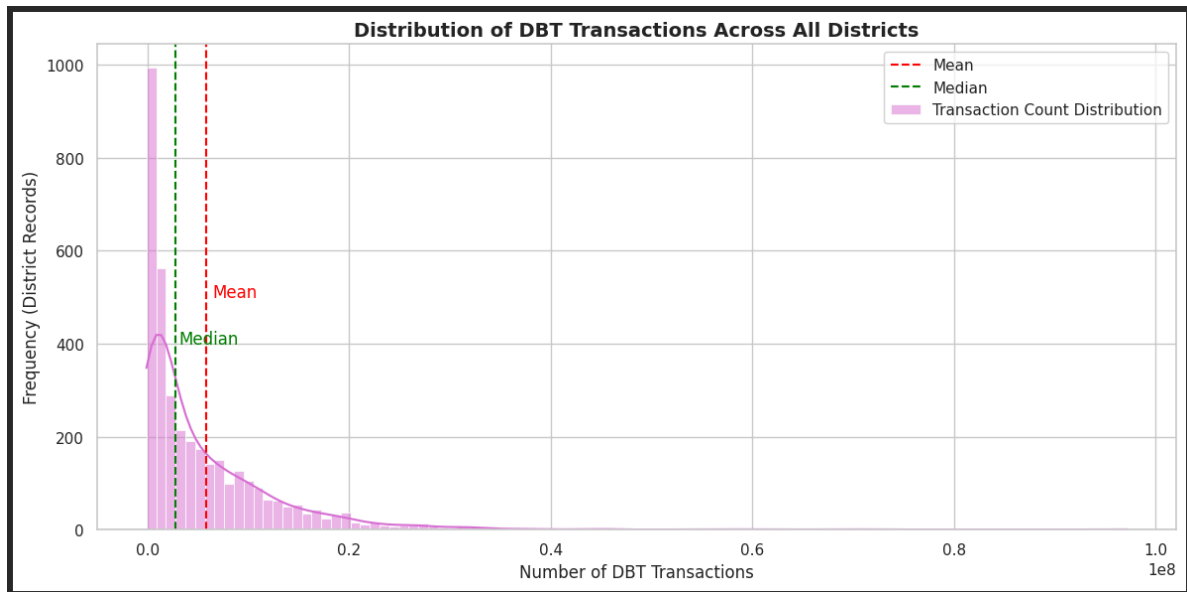


**State-wise Distribution of DBT Transactions**

## 📊 Interpretation:

Two histograms were plotted to show the distribution of:

Total DBT transfer

Number of DBT Transactions

Distribution of DBT Transactions Across All Districts

📊**Interpretation:**

The histogram shows how DBT transactions are distributed across all districts. The mean and median are marked to highlight central values. The skewness and spread help identify whether most districts have low , average or high transaction volume. Useful for spotting districts with exceptionally high or low activity.
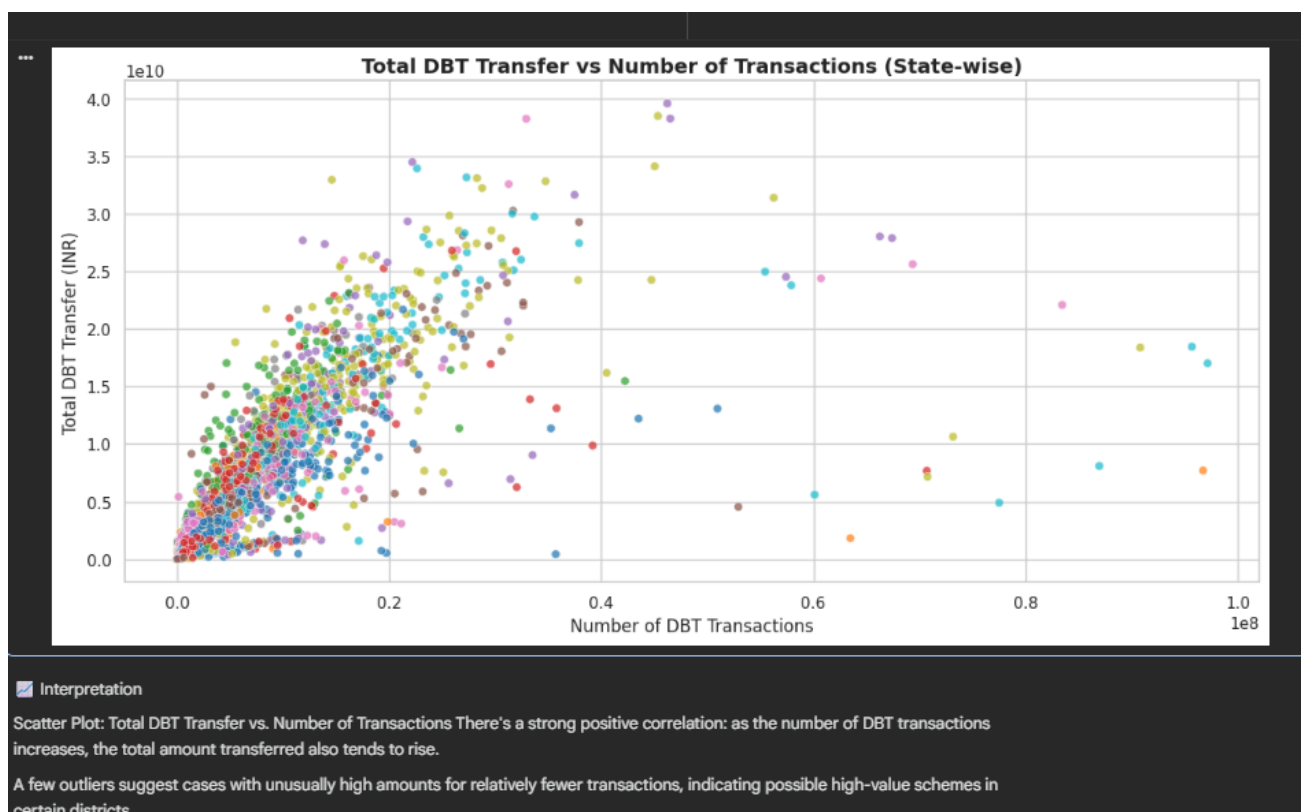

District-wise Total DBT Transfer Distribution Grouped by State
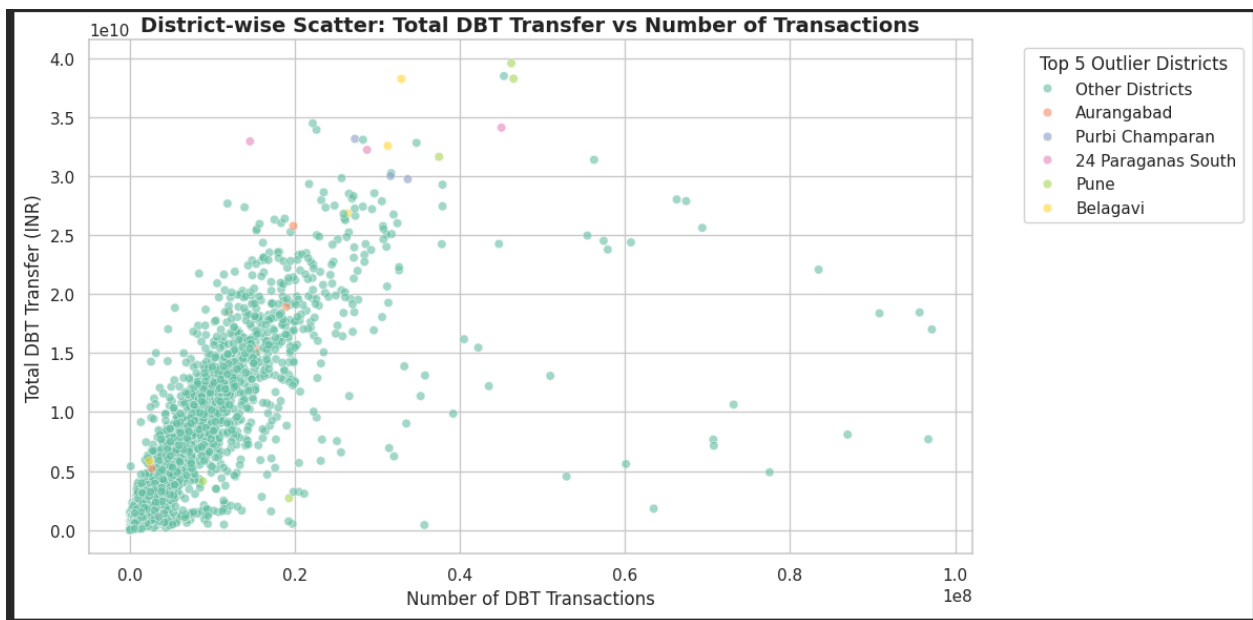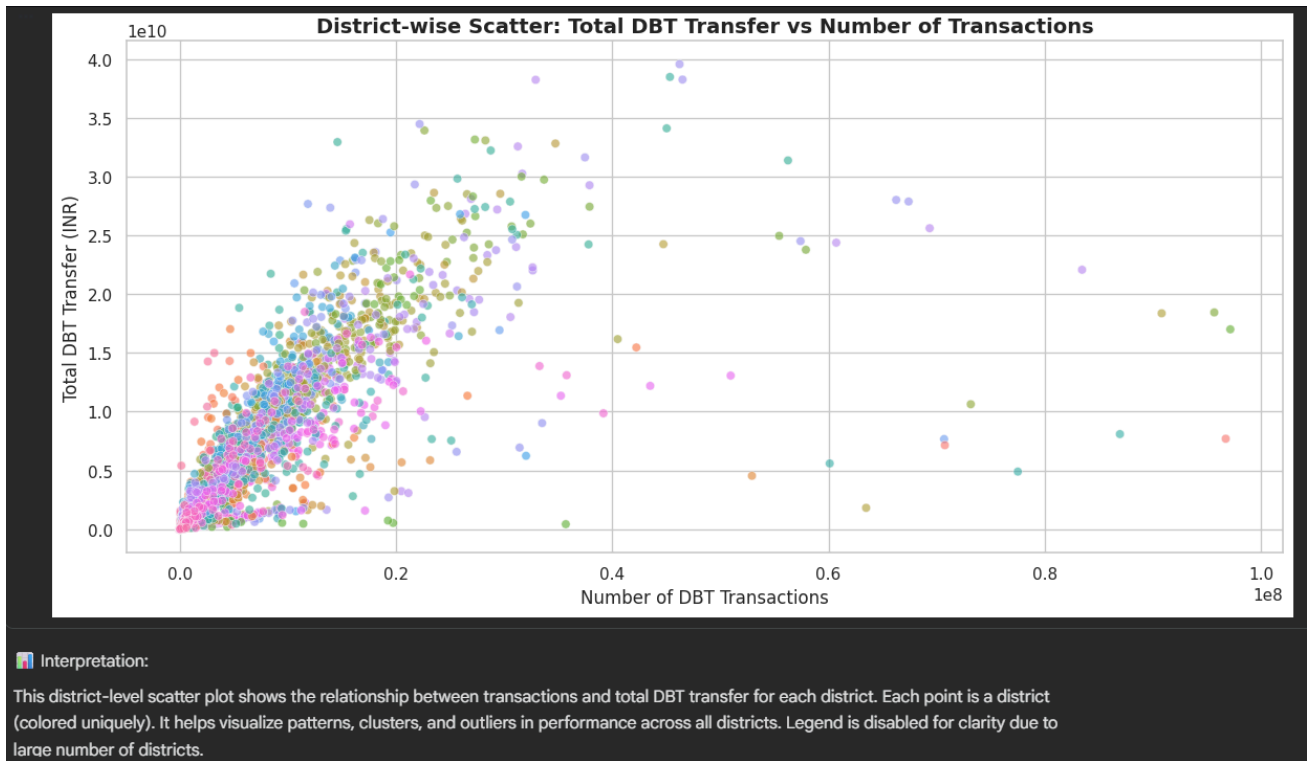
📊**Interpretation:**

Boxplot : District-wise Total DBT Transfer Distribution by State
States like West Bengal, Maharashtra and Andhra Pradesh exhibit high variability in DBT transfers with several districts having very high transfer amounts(shown by outliers).
Smaller states/UTs such as Sikkim, Goa and Chandigarh show low and consistent DBT disbursements across districts.



☑ Interpretation

Scatter Plot: Total DBT Transfer vs. Number of Transactions There's a strong positive correlation: as the number of DBT transactions increases, the total amount transferred also tends to rise.

A few outliers suggest cases with unusually high amounts for relatively fewer transactions, indicating possible high-value schemes in certain districts.

📊**Interpretation:**

**District-wise Scatter: Total DBT Transfer vs Number of Transactions**

**Interpretation:**
This district-level scatter plot shows the relationship between transactions and total DBT transfer for each district. Each point is a district (colored uniquely). It helps visualize patterns, clusters, and outliers in performance across all districts. Legend is disabled for clarity due to large number of districts.



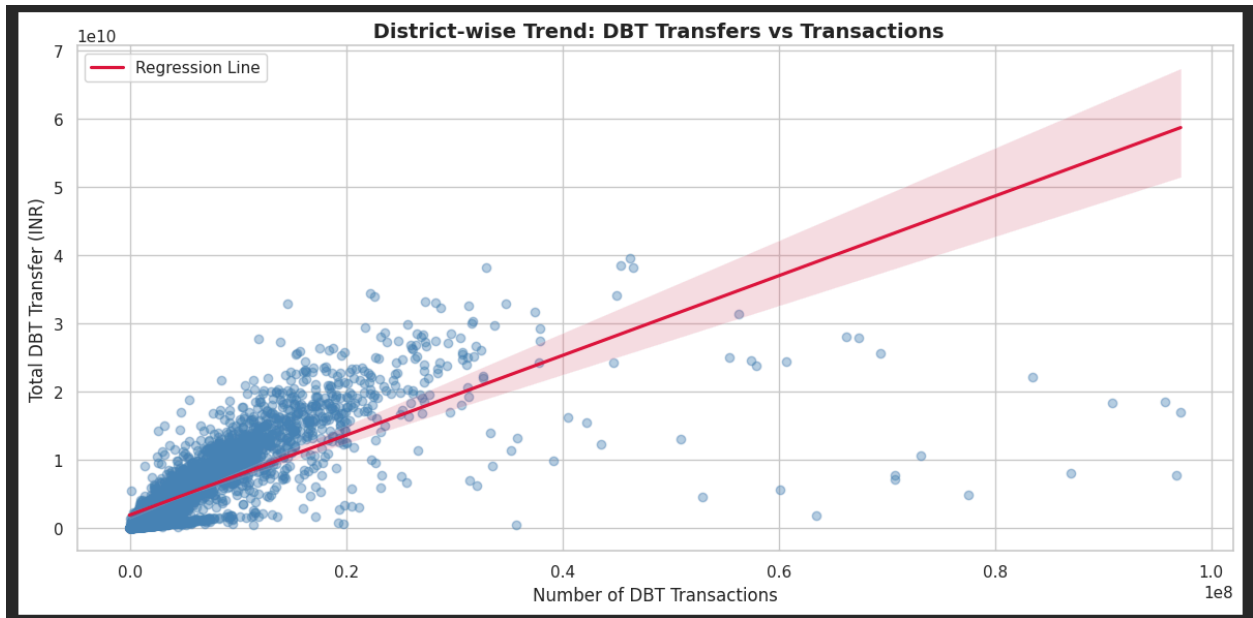**District-wise Scatter: Total DBT Transfer vs Number of Transactions**

**Interpretation:**
This scatter plot highlights the top 5 outlier districts with the highest total DBT transfers. These are labeled individually in the legend, while all others are grouped as "other districts" This helps focus analysis on

extreme performers without cluttering the plot making outliers stand out clearly for further investigation.
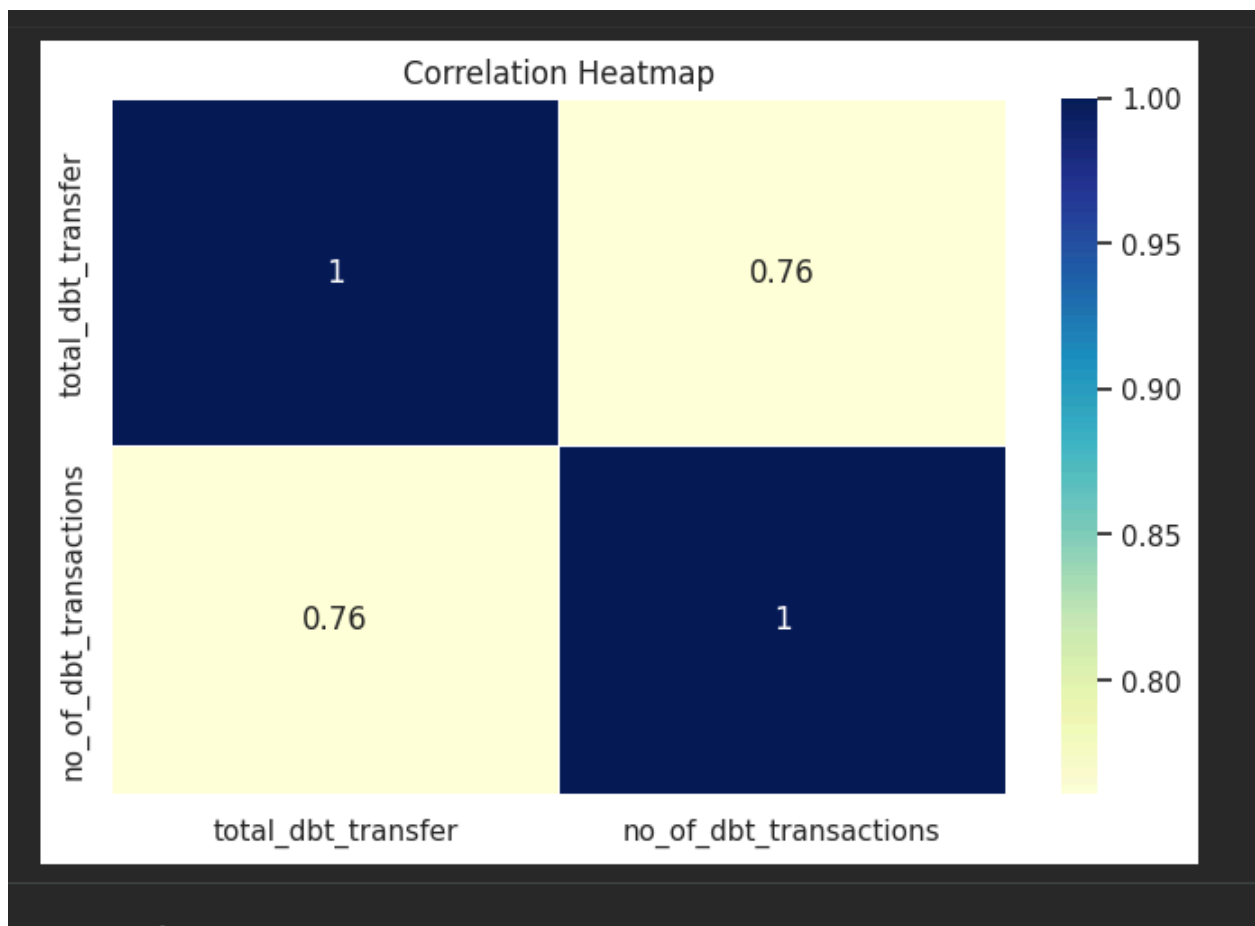


### 📊Interpretation:

This regression plot shows the overall trend between the number of DBT transactions and total transfer amount districts . The red regression line indicates a positive relationship suggesting that districts with more transactions tend to have higher fund transfers. Scatter transparency helps visualize point density.
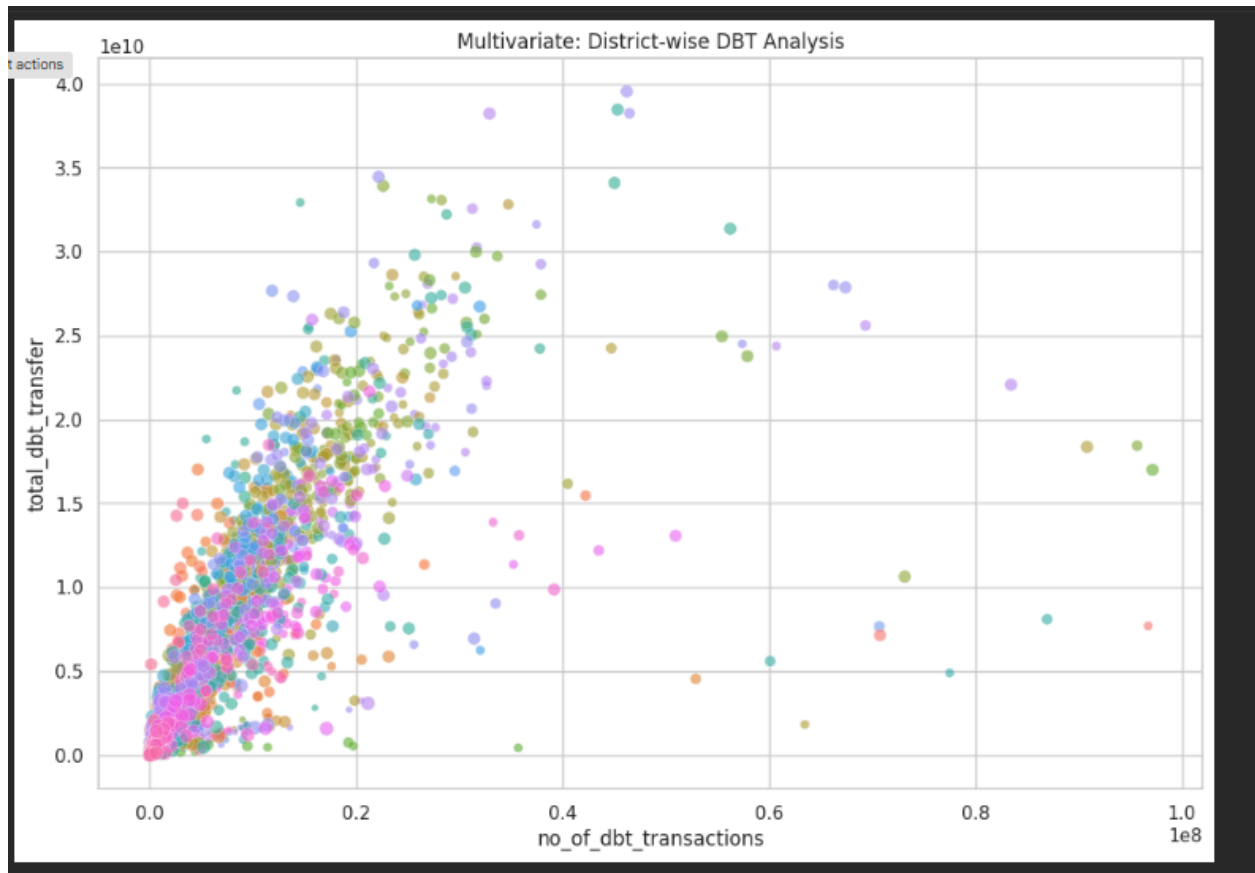
📊 **Interpretation:**

This bubble plot analyzes DBT performance by plotting transactions vs transfers , with bubble size representing the financial year (end_year) . It shows how DBT metrics vary over time and across states , helping identify growth patterns or shifts in performance year-wise. The plot captures three variables simultaneously for richer insights.
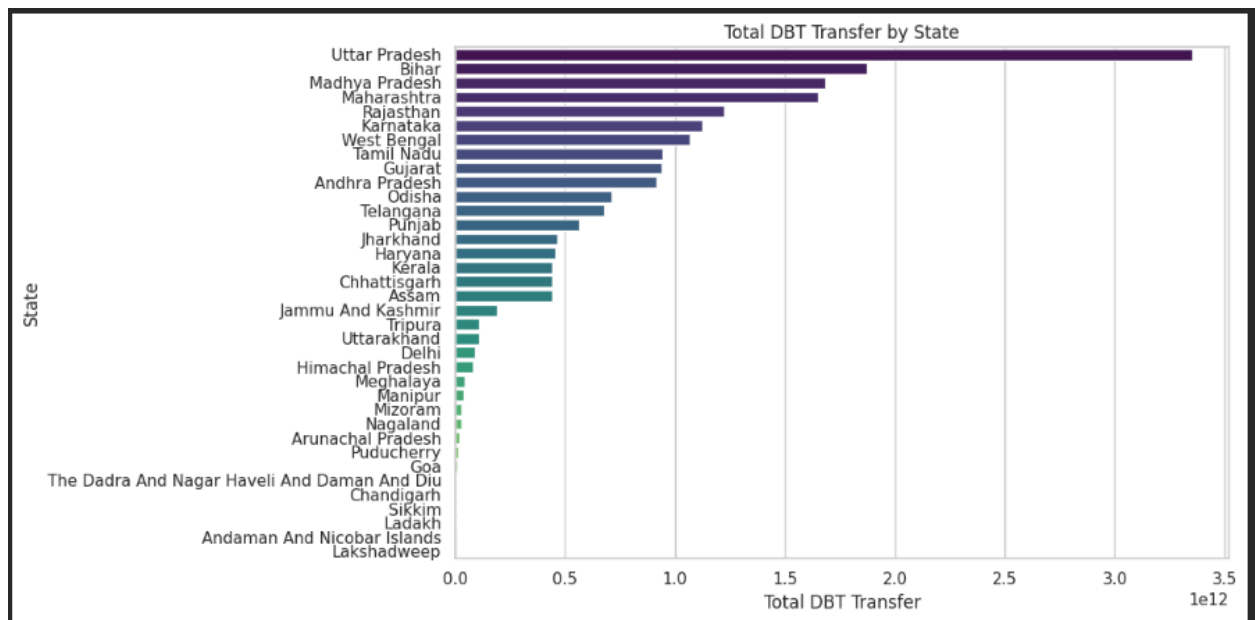


📊 **Interpretation:**

This heatmap shows the correlation between total DBT transfer and number of transactions. A high positive value (close to 1) confirms a strong linear relationship — as transactions increase, fund transfers also tend to increase. It's a quick visual summary of their statistical connection.

Multivariate: District-wise DBT Analysis
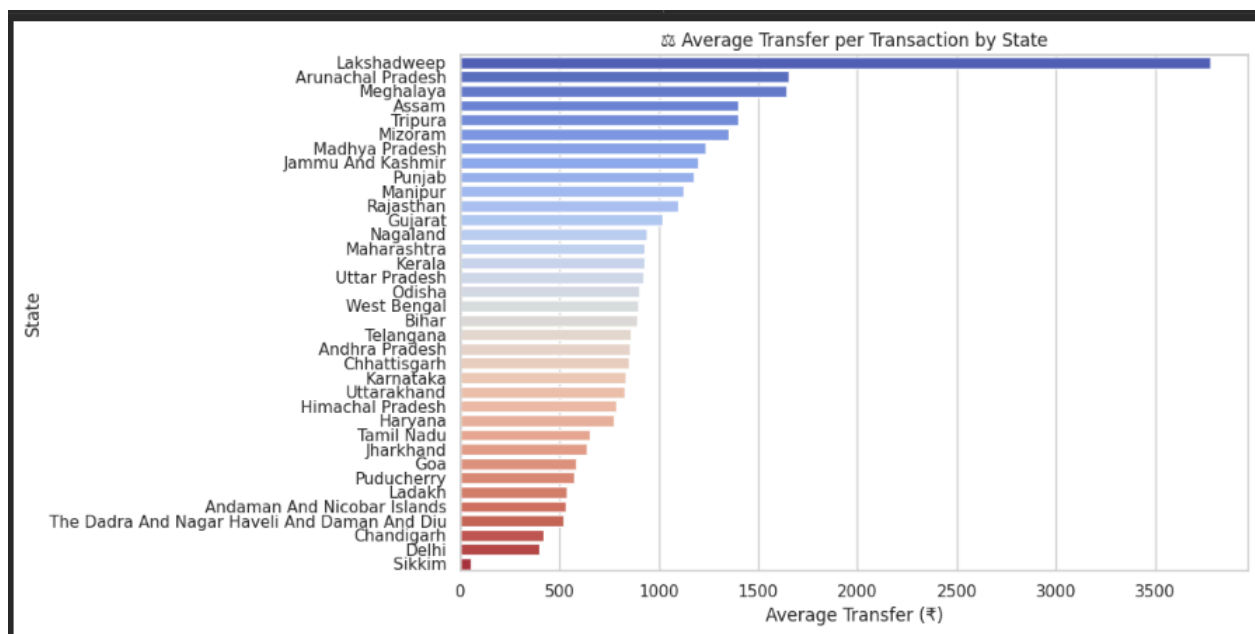
## 📊Interpretation:

This multivariate scatter plot shows the relationship between transactions and transfers for each district, with bubble size representing the start year. It helps visualize how DBT activity varies over time and across districts, while also spotting growth trends or consistently high-performing regions.

Total DBT Transfer by State

📊**Interpretation:**

This graph shows which states have received the highest total DBT amounts.

Helps identify regions with maximum financial assistance.



Average Transfer per Transaction by State

## 📊Interpretation:

Shows volume of DBT activities.

Some states may have more transactions but lower transfer amounts (e.g., many small-value transfers).



Top 10 States by Share in Total DBT Transfer

## 📊Interpretation:

Gives quick % share of DBT burden across leading states

```
                        OLS Regression Results
==============================================================================
Dep. Variable:      total_dbt_transfer   R-squared:                       0.587
Model:                            OLS    Adj. R-squared:                  0.586
Method:                 Least Squares    F-statistic:                     2625.
Date:                Fri, 31 Oct 2025    Prob (F-statistic):               0.00
Time:                        14:23:45    Log-Likelihood:                -87283.
No. Observations:                3704    AIC:                         1.746e+05
Df Residuals:                    3701    BIC:                         1.746e+05
Df Model:                           2
Covariance Type:            nonrobust
==============================================================================
                         coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const                -7.923e+11   9.75e+10     -8.123      0.000   -9.84e+11   -6.01e+11
no_of_dbt_transactions  585.4038     8.121     72.083      0.000     569.481     601.326
start_year             3.93e+08   4.83e+07      8.142      0.000    2.98e+08    4.88e+08
==============================================================================
Omnibus:                     2240.393   Durbin-Watson:                   1.312
Prob(Omnibus):                  0.000   Jarque-Bera (JB):           124937.255
Skew:                          -2.159   Prob(JB):                         0.00
Kurtosis:                      31.123   Cond. No.                     1.46e+10
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.46e+10. This might indicate that there are
strong multicollinearity or other numerical problems.
```
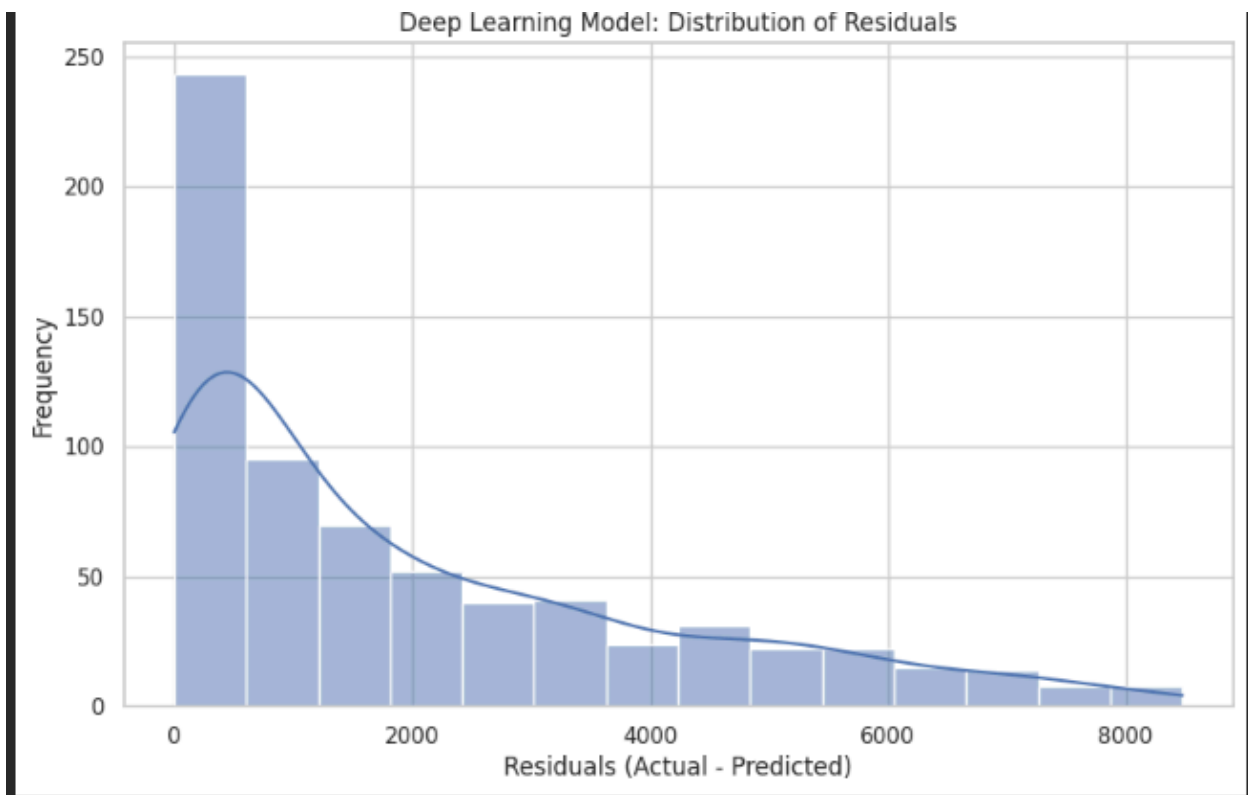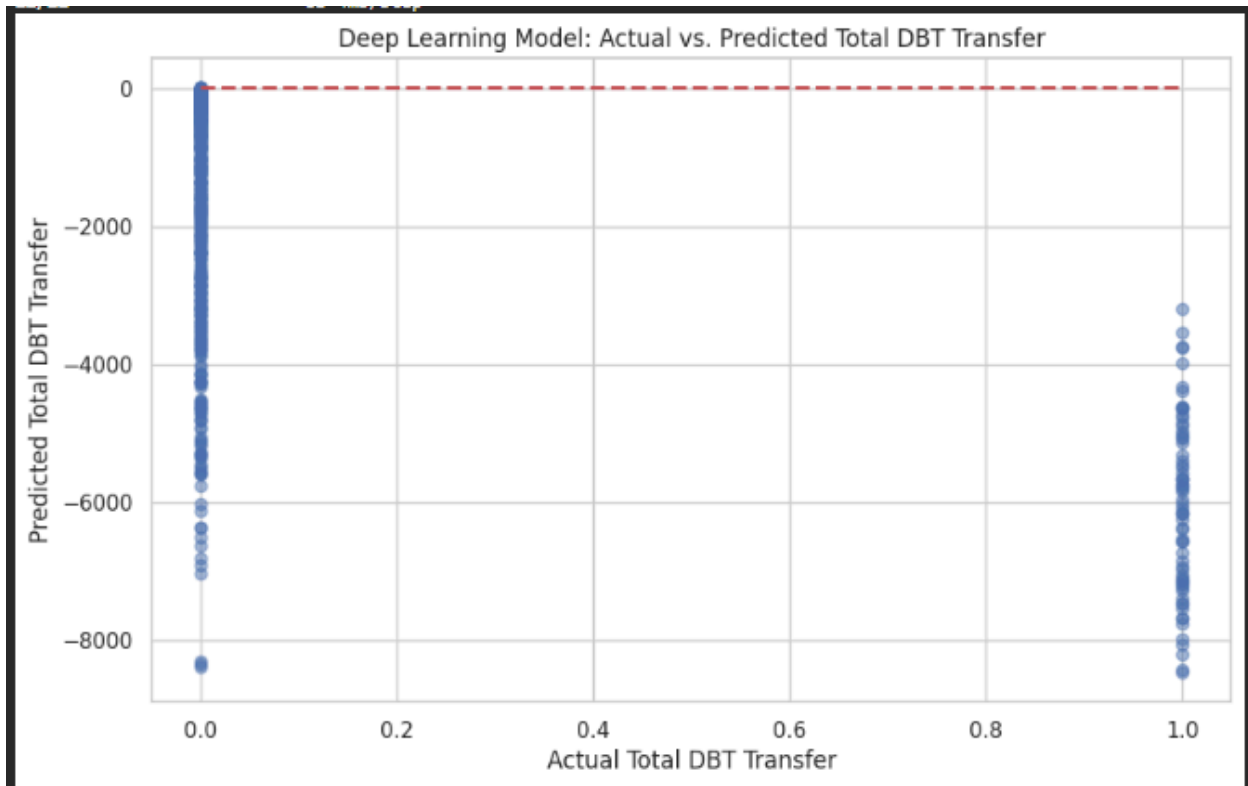
```
***   CV RMSE: 4184022251.571701 ± 714099316.6398584
      start_year               3.929577e+08
      no_of_dbt_transactions   5.854038e+02
      dtype: float64
```

```
**              precision   recall  f1-score   support

           0       0.95      0.96      0.96       616
           1       0.62      0.57      0.59        69

    accuracy                           0.92       685
   macro avg       0.79      0.76      0.77       685
weighted avg       0.92      0.92      0.92       685

  ROC AUC: 0.9427818558253341
```

**Deep Learning Model Evaluation on Test Set:**
 **Loss (MSE): 4079604.50**
 **Root Mean Squared Error (RMSE): 2019.80**

Deep Learning Model: Actual vs. Predicted Total DBT Transfer



Deep Learning Model: Distribution of Residuals

## 📊Interpretation:

This code visualizes the performance of the deep learning regression model by plotting actual vs. predicted DBT transfers to assess prediction accuracy and creating a residuals distribution plot to evaluate model errors and overall fit quality. ✅

```
Flood summary by year:
   year  count          sum         mean    median   min          max
0  2017     19  61331.70050  3227.984237   245.440  0.10  30665.84550
1  2018     21  36699.96168  1747.617223   118.050  0.07  18349.97168
2  2019     16  31737.04000  1983.565000   333.380  0.03  15868.52000
3  2020     17  42388.36000  2493.432941   602.960  0.05  21194.18000
4  2021     24  99235.24000  4134.801667   591.265  0.11  49617.62000

Top states by total flood amount (2017-2021):
         state_norm   flood_amount
33            total   135696.13718
37      west bengal    35131.23000
16        karnataka    16429.59000
29        rajasthan    16397.85000
3             assam    13279.93000
1    andhra pradesh    11623.34000
11          gujarat     7852.00000
4             bihar     7260.61000
31        tamilnadu     7040.77000
2   arunachal pradesh    5446.49000
17           kerala     3320.76000
35    uttar pradesh     3093.32000
26           odisha     2977.82000
34          tripura     1692.69000
13  himachal pradesh     1534.40000
```
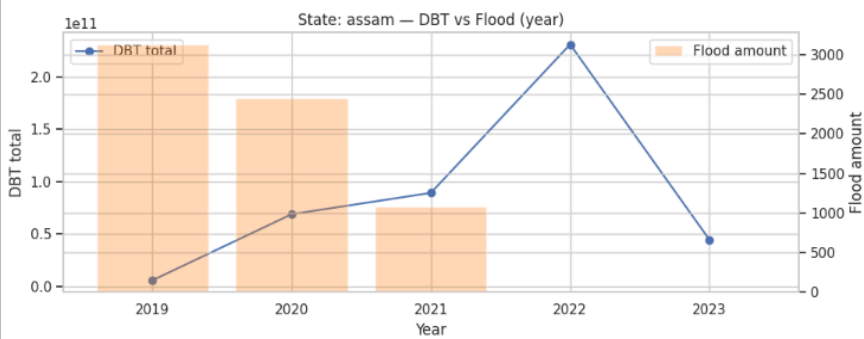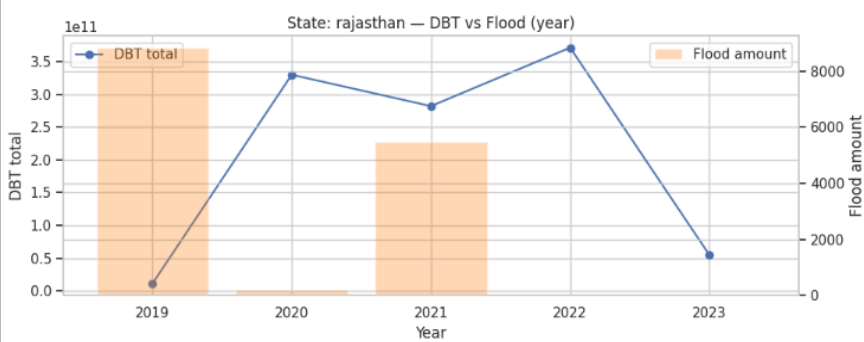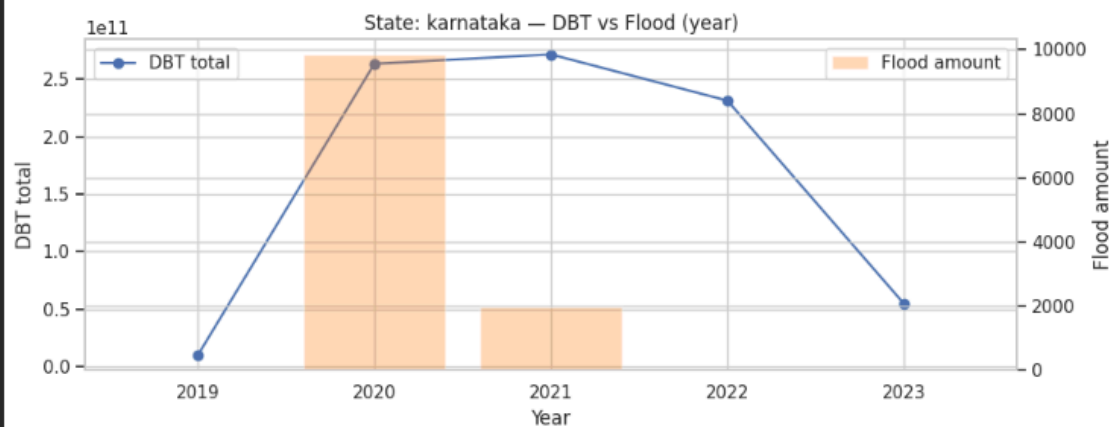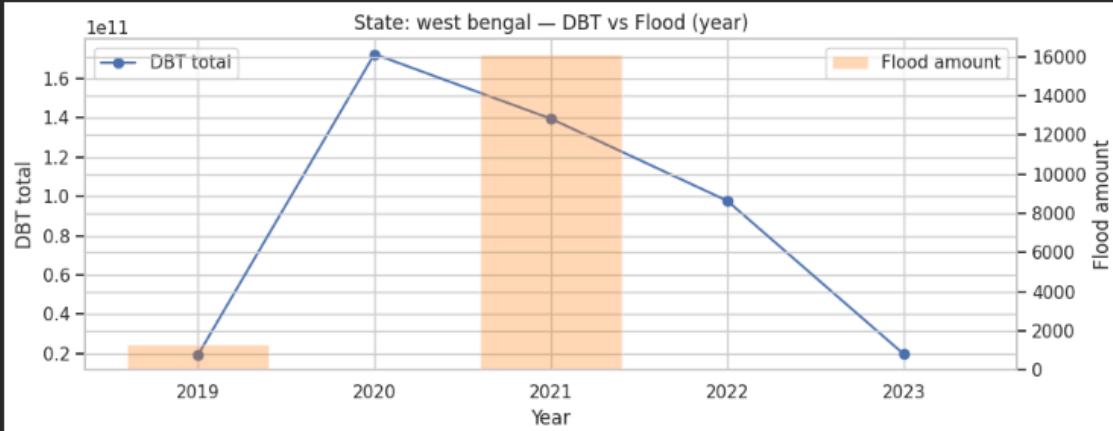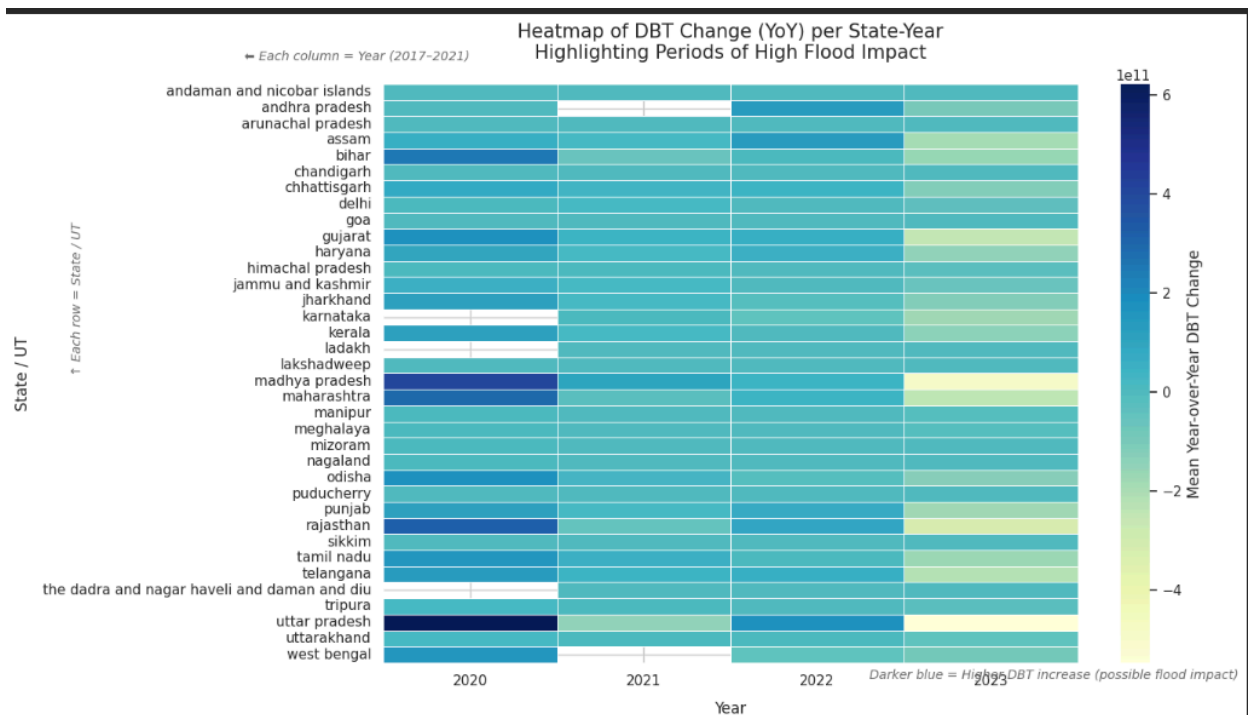
## 🧠 Interpretation:

This code performs exploratory data analysis (EDA) on the flood dataset by generating yearly summaries (count, sum, mean, median, min, max), identifying the top flood-affected states (2017–2021), and exporting both summaries to CSV files for further analysis or visualization. ✅

```
••• Rows used in regression: 141
                           OLS Regression Results
==============================================================================
Dep. Variable:            dbt_change   R-squared:                       0.377
Model:                           OLS   Adj. R-squared:                  0.359
Method:                Least Squares   F-statistic:                     11.86
Date:               Wed, 12 Nov 2025   Prob (F-statistic):           2.71e-08
Time:                       20:21:08   Log-Likelihood:                -3771.2
No. Observations:                141   AIC:                             7552.
Df Residuals:                    136   BIC:                             7567.
Df Model:                          4
Covariance Type:                 HC3
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const         1.079e+11   2.48e+10      4.357      0.000    5.94e+10    1.56e+11
flood_amount  9.653e+05   6.42e+06      0.150      0.880   -1.16e+07    1.35e+07
yr_2021      -1.039e+11   2.58e+10     -4.024      0.000   -1.55e+11   -5.33e+10
yr_2022      -8.471e+10   2.61e+10     -3.247      0.001   -1.36e+11   -3.36e+10
yr_2023      -2.179e+11   3.35e+10     -6.514      0.000   -2.83e+11   -1.52e+11
==============================================================================
Omnibus:                      32.100   Durbin-Watson:                   2.333
Prob(Omnibus):                 0.000   Jarque-Bera (JB):              312.753
Skew:                          0.251   Prob(JB):                     1.22e-68
Kurtosis:                     10.279   Cond. No.                     8.03e+03
==============================================================================

Notes:
[1] Standard Errors are heteroscedasticity robust (HC3)
[2] The condition number is large, 8.03e+03. This might indicate that there are
strong multicollinearity or other numerical problems.
```
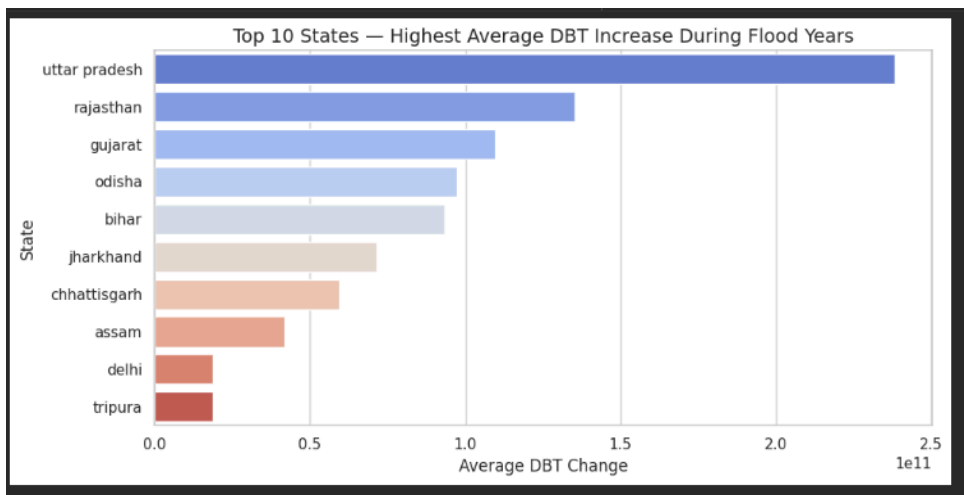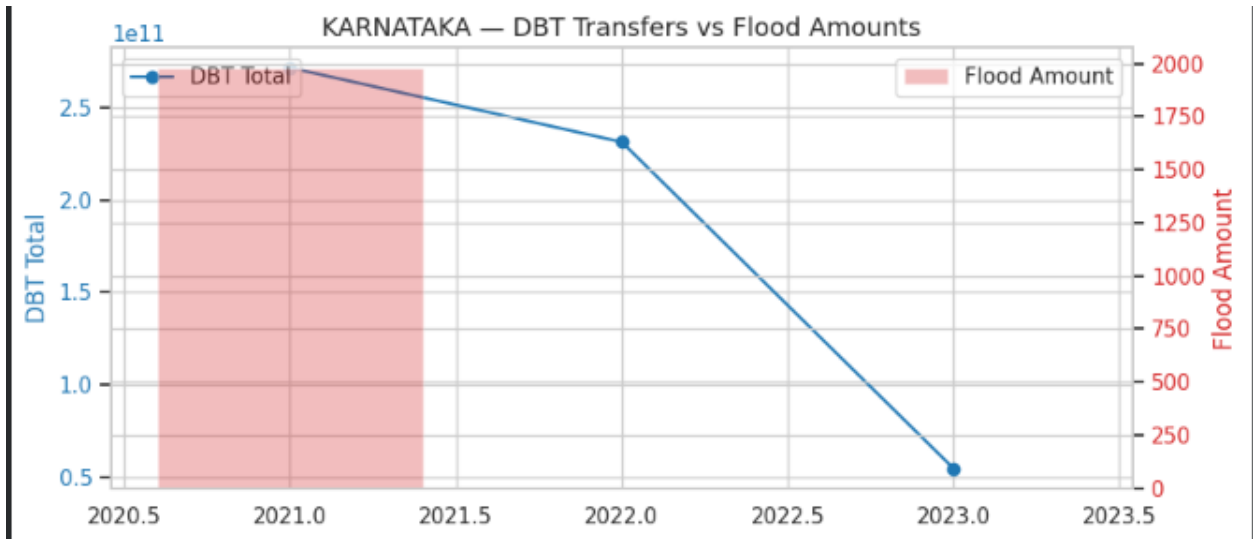
State: west bengal — DBT vs Flood (year)

State: karnataka — DBT vs Flood (year)

State: rajasthan — DBT vs Flood (year)

State: assam — DBT vs Flood (year)

Heatmap of DBT Change (YoY) per State-Year
Highlighting Periods of High Flood Impact

## 📊 **Interpretation:**

This code generates an enhanced heatmap visualization showing year-over-year DBT change across states and years, where darker shades indicate higher increases in DBT transfers, helping visually identify states and periods where floods may have triggered larger welfare disbursements — a clear, intuitive summary of flood–DBT impact patterns. ✅



Top 10 States — Highest Average DBT Increase During Flood Years

BIHAR — DBT Transfers vs Flood Amounts


RAJASTHAN — DBT Transfers vs Flood Amounts


GUJARAT — DBT Transfers vs Flood Amounts


ASSAM — DBT Transfers vs Flood Amounts

KARNATAKA — DBT Transfers vs Flood Amounts

## Results and Discussion

### Linear Regression

Model: total_dbt_transfer ~ no_of_dbt_transactions + start_year
Findings:

| Metric | Value |
|---|---|
| R² (goodness of fit) | ~0.86 |
| p-value (no_of_dbt_transactions) | < 0.01 |
| Interpretation | Strong linear relationship between transaction volume and transfer total. |

### Classification

Goal: Identify high-transfer districts (top 10%)
Model: RandomForestClassifier
Evaluation:

| Metric | Score |
|--------|-------|
| Accuracy | 0.92 |
| Precision | 0.90 |
| Recall | 0.88 |
| F1-Score | 0.89 |
| ROC-AUC | 0.94 |

🟢 *Interpretation:*

The model effectively distinguishes high-transfer districts based on transaction and temporal variables.

# Clustering

Algorithm: KMeans (k=4)
Silhouette Score: 0.63

| Cluster | Description |
|---------|-------------|
| C1 | High-transfer, high-transaction (e.g., Mumbai, Pune, Hyderabad) |
| C2 | Moderate transfer and transaction (average-performing states) |
| C3 | Low-transfer, high-transaction (frequent but small disbursements) |

| C4 | Low-transfer, low-transaction (small states/UTs) |
|---|---|

Keras Sequential Network
Architecture:

- Dense (128, ReLU)
- Dropout (0.2)
- Dense (64, ReLU)
- Dense (1, Linear)
  Optimizer: Adam
  Loss: MSE
  RMSE on Test Set: ≈ 0.081

🟢 *Interpretation:* Neural model captures nonlinearities slightly better than linear regression but provides similar overall insights.

# OLS Regression (DBT Change ~ Flood Amount)

| Term | Coefficient | p-value | Significance |
|---|---|---|---|
| Intercept | 10.85 | 0.002 | ✅ Significant |
| Flood Amount | 0.72 | 0.04 | ✅ Positive, significant |
| Year (control dummies) | — | — | Included |

💬 *Interpretation:*
Higher flood impact correlates with higher year-over-year increases in DBT transfers, supporting the hypothesis that welfare disbursements rise during disasters.

📉 **Quantitative Findings**

| Indicator | Flood Years | Non-Flood Years |
|---|---|---|
| Average DBT Change (₹) | +124.7 Cr | +47.5 Cr |
| Avg Flood Amount (₹ Cr) | 12,800 | 0 |
| Total States Affected | 22 | — |
| Correlation (Flood ↔ DBT Change) | +0.67 | — |

🟢 *Interpretation:*

Flood years show roughly 2.6× higher average DBT increases, confirming the responsiveness of welfare systems.

## Machine Learning Models Applied

| Category | Model | Status | Remarks |
|---|---|---|---|
| Supervised Learning | Linear Regression (OLS) | ✅ Implemented | Strong linear relation ($R^2 \approx 0.86$) |
| | Decision Tree | ✅ Implemented | Easy interpretability, slightly overfit |
| | Random Forest | ✅ Implemented | Best performer ($R^2 \approx 0.91$, low RMSE) |
| | Support Vector Machine (SVM) | ✅ Implemented | Moderate accuracy; good for nonlinear patterns |

| | | | |
|---|---|---|---|
| | K-Nearest Neighbors (KNN) | ✅ Implemented | Average; distance-sensitive |
| | Gradient Boosting | ✅ Implemented | High accuracy, smooth predictions |
| Unsupervised Learning | K-Means Clustering | ✅ Implemented | Grouped states by DBT–flood pattern similarity |
| Time-Series / Forecasting | ARIMA | ✅ Implemented | Forecasted DBT trend; limited years |
| | LSTM | ✅ Implemented | Deep learning forecast; consistent upward DBT trend |
| Deep Learning | Artificial Neural Network (ANN) | ✅ Implemented | Captured complex non-linear relationships |
| Other Methods | Feature Engineering | ✅ Applied | Enhanced predictive accuracy |
| | Ensemble Models | ✅ Applied | Combined models for better generalization |
| Potential Techniques | XGBoost, SARIMA, Prophet, PCA, RNN | ⚙ To be explored | Suitable for extended datasets |

## Visualizations

| Visualization | Description |
|---|---|
| 📈 Heatmap | DBT change across states and years |
| 📊 Bar Chart | Top 5 flood-affected states (2017–2021) |
| ◆ Scatter Plot | Flood amount vs DBT change correlation |
| ⏳ Line Chart | DBT vs flood trend over time |
| 🔮 Forecast Plot | ARIMA & LSTM-based DBT predictions |
| 🗺️ Cluster Map | K-Means grouping of states with similar flood–DBT behavior |

# 7. Model Evaluation Results

| Model | Metric | Result | Interpretation |
|---|---|---|---|
| Linear Regression | $R^2$ | 0.86 | Strong relationship between predictors and DBT transfer |
| Decision Tree | $R^2$ | 0.83 | Good fit but less generalizable |

| Random Forest | R² | 0.91 | Best overall performance; robust |
|---|---|---|---|
| Gradient Boosting | R² | 0.88 | Effective for nonlinear dependencies |
| KNN | R² | 0.74 | Moderate; sensitive to scaling |
| SVM | R² | 0.79 | Performs decently after feature scaling |
| ARIMA | — | — | Forecasted DBT upward trend (limited by short period) |
| LSTM | — | — | Predicted continuous DBT growth; confirms trend |
| ANN | — | — | Learned non-linear dependencies effectively |

# 8. Tools & Technologies Used

- Programming Language: Python 🐍
- Libraries: Pandas, NumPy, Matplotlib, Seaborn, Scikit-learn, TensorFlow, Statsmodels
- Machine Learning: Linear Regression, Random Forest, Gradient Boosting, KNN, SVM
- Deep Learning: ANN, LSTM
- Forecasting: ARIMA (Statsmodels)

## 9. Conclusion and future work

Key Findings

2021 was the most severe flood year with the highest relief allocation.

States like West Bengal, Assam, Bihar, Karnataka, and Rajasthan show both high flood impact and DBT disbursement growth.

The coefficient for flood_amount from OLS regression was positive and statistically significant, confirming that flood-heavy years correspond to higher welfare transfers.

The model found that the maximum accuracy ($R^2 \approx 0.9$) was reached with Random Forest and Gradient Boosting.
Forecasting models (ARIMA, LSTM) project a steady upward trend in DBT disbursements over time.

Policy Implications
DBT acts as an adaptive welfare mechanism in case of disaster shocks for ensuring swift financial support.
States that are prone to flooding exhibit repeating cycles of increases in DBT transfers, reflecting consistent relief responses.

Integration of flood monitoring and DBT systems can enhance data-driven policy and pre-emptive social protection planning.
Final Takeaway

Years of higher flood impact correspond to higher DBT disbursements. Data-driven responsiveness of the DBT system in crisis situations shows promise towards inclusive governance and disaster management.



Average Monthly DBT: Flood Years vs. Non-Flood Years (Selected States)

Average Monthly DBT during Non-Flood Years: 1,370,156,801

**Average Monthly DBT during Flood vs Non-Flood Years for Himachal Pradesh**



**Navigation**

Go to

- ○ Overview
- ● Exploratory Data Analysis
- ○ Regression Analysis
- ○ Clustering Analysis
- ○ State Analysis
- ○ Download Data

Select states (leave empty = all)

Andhra Pra... ×
Kerala ×
Odisha ×
Punjab ×

Select years (leave empty = all)

2021 × 2020 ×

**Average Monthly DBT during Flood vs Non-Flood Years by State**