# Groundwater Level Prediction Using Multiple Linear Regression
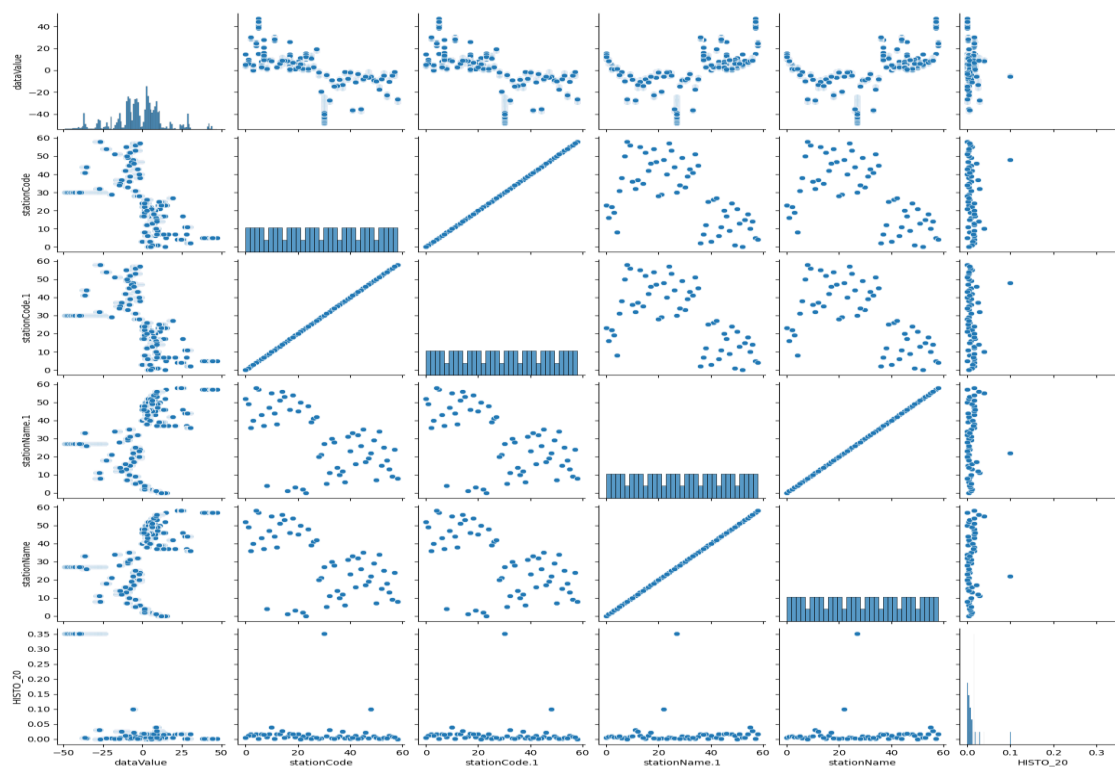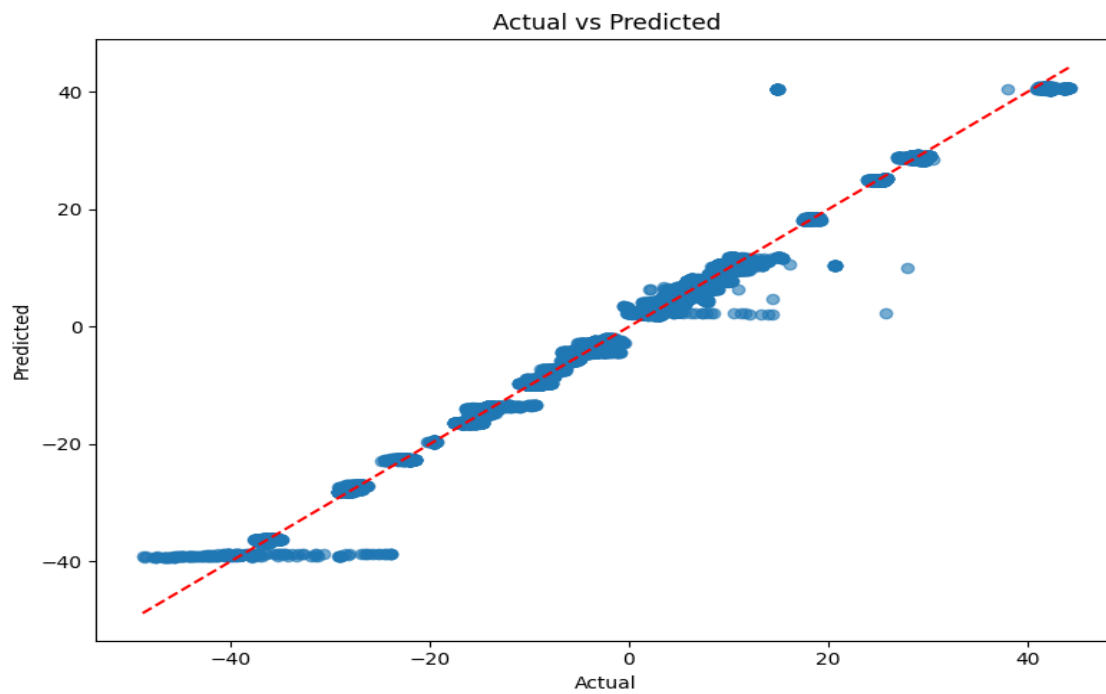
## 1. Scenario

- Delhi NCR
- Objective: Identify key drivers of groundwater depletion and predict groundwater levels at unsampled locations.
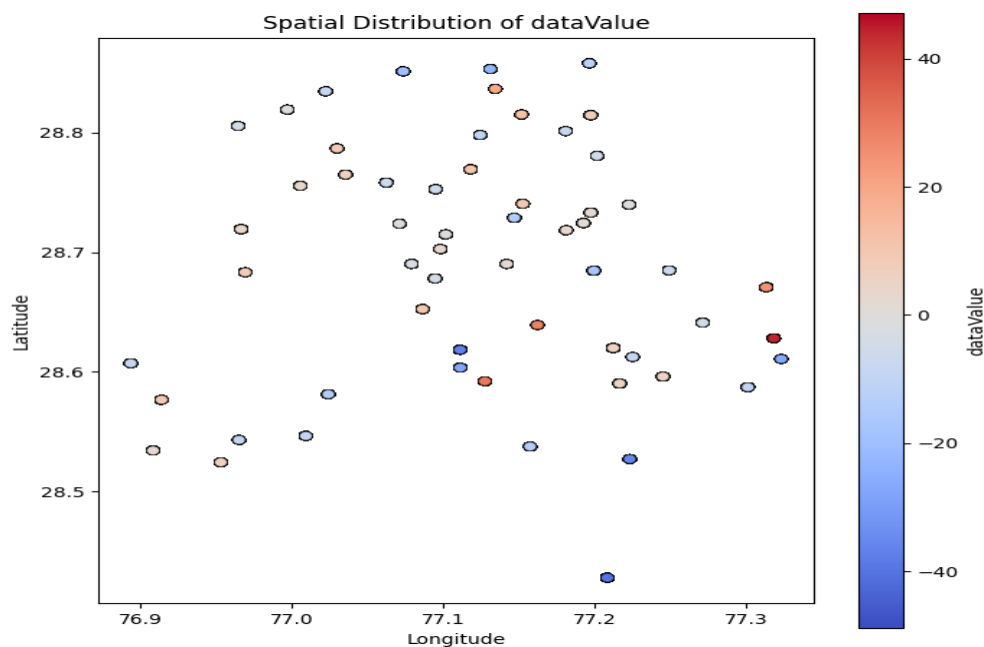
## 2. Methodology

- **Data:** Preprocessed groundwater dataset with 175 columns including the target 'data Value'.
- **Dependent variable**: dataValue
- **Independent variables**: 58 selected features (BIC criterion)
- **Spatial unit**: District
- **Temporal unit**: Daily
- **Model equation**: Multiple Linear Regression (OLS)
- **Data Acquisition**: Data acquired from India WRIS, Bhuvan ISRO, Copernicus Climate Data Store, NICES Portal, SHRUG Atlas.
- **Data Merging**: Merged datasets on district and date.
- **Data Preprocessing**: Missing values handled, outliers removed, data merged appropriately.
- **Model Specification**: Defined model structure and selected features.
- **Model Training**: Trained the model using the training dataset.
- **Model Evaluation**: Evaluated model performance using test dataset and metrics like $R^2$, RMSE.
- **Model Diagnostics**: Residuals analyzed for patterns.
- **EDA**: Scatter plots, pair plots, and spatial-temporal plots analyzed to observe trends.

## 3. Exploratory Data Analysis (EDA)

- Scatter plots and pair plots were used to explore relationships between features and groundwater levels.
- Spatial and temporal trends were observed.
- Data was cleaned and outliers removed.

Pair-plots of features with GWL Values

Spatial Distribution of dataValue

## 4. Model Assumptions

- Linearity: Relationships between predictors and target are linear.
- No Perfect Multicollinearity: Checked correlations among predictors.
- Exogeneity: Residuals uncorrelated with predictors.
- Homoscedasticity: Residuals have constant variance.

## 5. Model Selection

- Compared models using AIC and BIC.
- BIC-selected model (58 predictors) was chosen for analysis.

## 6. Model Estimation & Diagnostics

- Top 12 Coefficients:

| Feature | Coefficient | Std_Error | t_value | p_value |
|---------|-------------|-----------|---------|---------|
| stationCode.1 | -0.22627 | 0.000746 | -303.44491 | 0.0 |
| HISTO_20 | -81.591015 | 1.554332 | -52.49266 | 0.0 |
| Rainfall(mm)_2024 | -0.539686 | 0.003682 | -146.55966 | 0.0 |

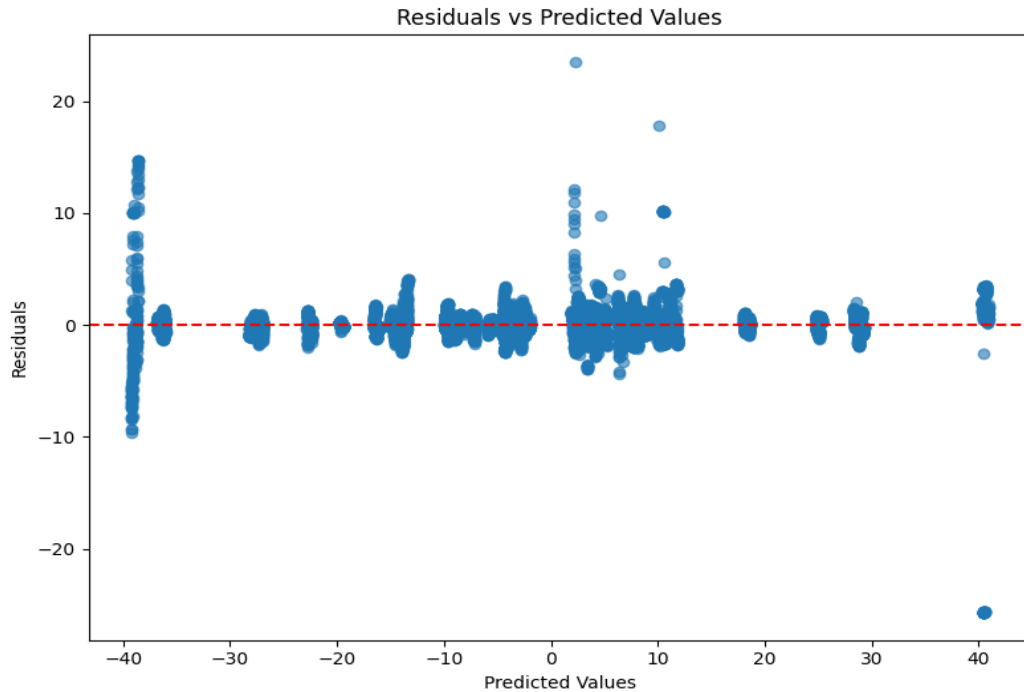| | | | | |
|---|---|---|---|---|
| _count | -5.4e-05 | 0.0 | -377.401993 | 0.0 |
| HISTO_80 | 574.507576 | 9.251323 | 62.100043 | 0.0 |
| shape_leng | -0.002849 | 1e-05 | -273.703361 | 0.0 |
| Pre Monsoon of GW Trend_2024 | 47.434549 | 0.648234 | 73.175021 | 0.0 |
| stationName | 0.592188 | 0.002755 | 214.934334 | 0.0 |
| Categorization of Assessment Unit_2023 | -33.680615 | 0.38376 | -87.764782 | 0.0 |
| sand_5-15cm_mean | 1.608464 | 0.01368 | 117.575387 | 0.0 |

**Model Fit Metrics:**
- R_squared(training): 0.9904
- Adj_R_squared: 0.9904
- F_stat: 49925.6491
- F_pvalue: 0.0000
- Num_Predictors: 58.0000


## 7. Predictions & Evaluation
- R-squared(testing data): 0.9868
- RMSE: 1.7404
- MSE: 3.0291
- Plots:

**Actual vs Predicted**



**Histogram of Residuals**



**Residuals vs Fitted Values**

Residuals vs Predicted Values



- Residual Plot
- The residual plot shows the difference between observed and predicted values. It helps in diagnosing the model fit.
- Interpretation: No clear pattern suggests a good fit.
- Action: Consider refining the model if patterns are detected.

## 8. Significant Features & Interpretation

| | Section | Feature | Impact | Coef | P-value |
|---|---|---|---|---|---|
| 2 | Significant Features | const | Positive | 4108.849607183905 | 6.439817055628084e-165 |
| 3 | Significant Features | stationCode.1 | Negati... | -0.2262700591082868 | 0.0 |
| 4 | Significant Features | HISTO_20 | Negati... | -81.59101511042891 | 0.0 |
| 5 | Significant Features | Rainfall (mm)_2024 | Negati... | -0.5396856566801052 | 0.0 |
| 6 | Significant Features | _count | Negati... | -5.4073280147480926e-05 | 0.0 |
| 7 | Significant Features | shape_leng | Negati... | -0.0028493173731542 | 0.0 |
| 8 | Significant Features | HISTO_80 | Positive | 574.5075759265251 | 0.0 |
| 9 | Significant Features | Pre Monsoon of GW Trend_2024 | Positive | 47.434548616037375 | 0.0 |
| 10 | Significant Features | district | Negati... | -9.034111405548687 | 4.234342942335636e-209 |

**Model Fit Metrics:**

| | A | B |
|---|---|---|
| 1 | **Metric** | **Value** |
| 2 | R_squared(training) | 0.99039004655034... |
| 3 | Adj_R_squared | 0.99037020925100... |
| 4 | F_stat | 49925.64913486291 |
| 5 | F_pvalue | 0.0 |
| 6 | Num_Predictors | 58.0 |

**Model Prediction Metrics:**

| | A | B |
|---|---|---|
| 1 | **Metric** | **Value** |
| 2 | R-squared(testing data) | 0.98682877119432... |
| 3 | RMSE | 1.74044178360589... |
| 4 | MSE | 3.02913760212128... |
| 5 | | |

**Confidence in Interpretation:**
- Description: The model explains most of the variability in groundwater levels (high R-squared).
- Significant features with $p < 0.05$ are likely reliable predictors.
- Prediction metrics (RMSE for regression and Accuracy/F1 for categorized classes) indicate reasonable predictive power.

## 9. Conclusion & Policy Implications
- The model identifies key factors affecting groundwater levels.
- Predictions can guide water resource planning.
- Recommendations: Monitor significant drivers and use the model for short-term planning and risk assessment.

## 10. References
- India WRIS: https://indiawris.gov.in/wris/
- Bhuvan ISRO: https://bhuvan-app1.nrsc.gov.in/2dresources/bhuvanstore2.php
- Copernicus Climate Data Store: https://cds.climate.copernicus.eu/#!/home
- NICES Portal: https://nices.nrsc.gov.in/
- SHRUG Atlas: https://www.devdatalab.org/atlas