
CAPSTONE PROJECT

IMPROVED SOURCE OF DRINKING WATER

Presented By:

- 1. Student Name - Amol Govinda Kumbhar**
- 2. College Name - Shram Sadhana Bombay Trust College of Engineering and Technology Jalgaon**
- 3. Department - Computer Engineering**

OUTLINE

- Υ **Problem Statement**
- Υ **Proposed System/Solution**
- Υ **System Development Approach**
- Υ **Algorithm & Deployment**
- Υ **Result (Output Image)**
- Υ **Conclusion**
- Υ **Future Scope**
- Υ **References**

PROBLEM STATEMENT

Example: Access to safe and improved sources of drinking water remains a critical issue in India, especially in rural and underdeveloped regions. Despite ongoing efforts under the Sustainable Development Goals (SDGs), inequalities persist across states and socio-economic groups. The challenge lies in analyzing data to identify patterns and disparities in water accessibility, which hinders evidence-based policymaking to ensure equitable access to clean water.

PROPOSED SOLUTION

- Y The proposed system aims to analyze access to improved drinking water sources using machine learning techniques to identify disparities and inform policy. The solution will consist of the following components:
- Y **Data Collection:**
 - 78th Round Multiple Indicator Survey (MIS) dataset from Al Kosh (e.g., state, rural/urban status, socio-economic group).
 - Additional data: Government reports on water infrastructure and migration trends.
- Y **Data Preprocessing:**
 - Y Clean data to address missing values and inconsistencies.
 - Y Feature engineering: Create indicators like water access percentage and socio-economic indices.
- Y **Machine Learning Algorithm:**
 - Y Use Random Forest for classification to predict water access levels (e.g., improved vs. unimproved).
 - Y Utilize IBM Granite for processing textual data from reports.
- Y **Deployment:**
 - Y Develop a Flask-based web app hosted on IBM Cloud for real-time analysis and visualization.
 - Y Enable policymakers to input regional data and view insights.
- Y **Evaluation:**
 - Y Assess performance using accuracy, precision, and recall.
 - Y Validate findings with statistical tests and continuous data updates.

SYSTEM APPROACH

The "System Approach" section outlines the overall strategy and methodology for developing the water access analysis system. Here's the structure:

- Y **System requirements :**

- Y -Hardware- Laptop/server with 16GB RAM, 2.5GHz CPU.

- Y -Software- IBM Watson Studio, Python 3.9+, Flask, IBM Cloud Pak for Data .

- Y **Library required to build the model :**

- Y - `pandas` : Data manipulation.

- Y - `numpy` : Numerical computations.

- Y - `scikit-learn` : Preprocessing and modeling (Random Forest).

- Y - `ibm-watson` : Integration with IBM Granite.

- Y - `matplotlib/seaborn` : Visualization.

- Y - `flask` : Web app deployment

ALGORITHM & DEPLOYMENT

Y Algorithm Selection

Y - **Algorithm:** Random Forest.

Y - **Justification:** Effective for classification tasks with multiple features, handles imbalanced data, and provides feature importance.

Y Data Input

Y - Features: State, rural/urban status, socio-economic group, household size, migration status, water source type.

Y Training Process

Y - Train Random Forest on the MIS dataset using an 80-20 train-test split.

Y - Apply cross-validation and tune hyperparameters (e.g., number of trees, max depth).

Y Prediction Process

Y - Input regional data via the web app.

Y - Output classification of water access (improved/unimproved) and disparity insights.

Y Deployment

Y - Flask-based web app hosted on IBM Cloud Lite.

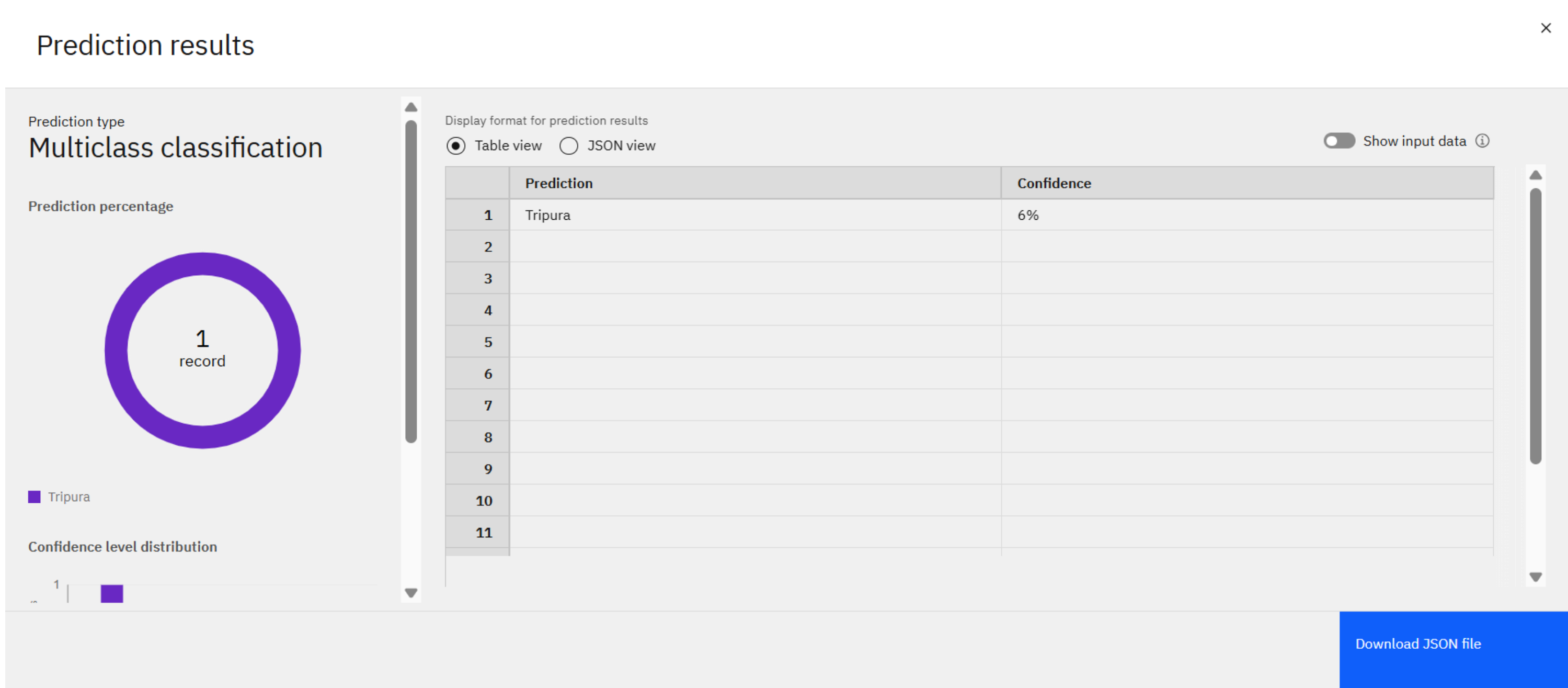
Y - IBM Granite processes textual inputs (e.g., policy documents) for feature extraction.

Y - Provide a scalable dashboard for policymakers.

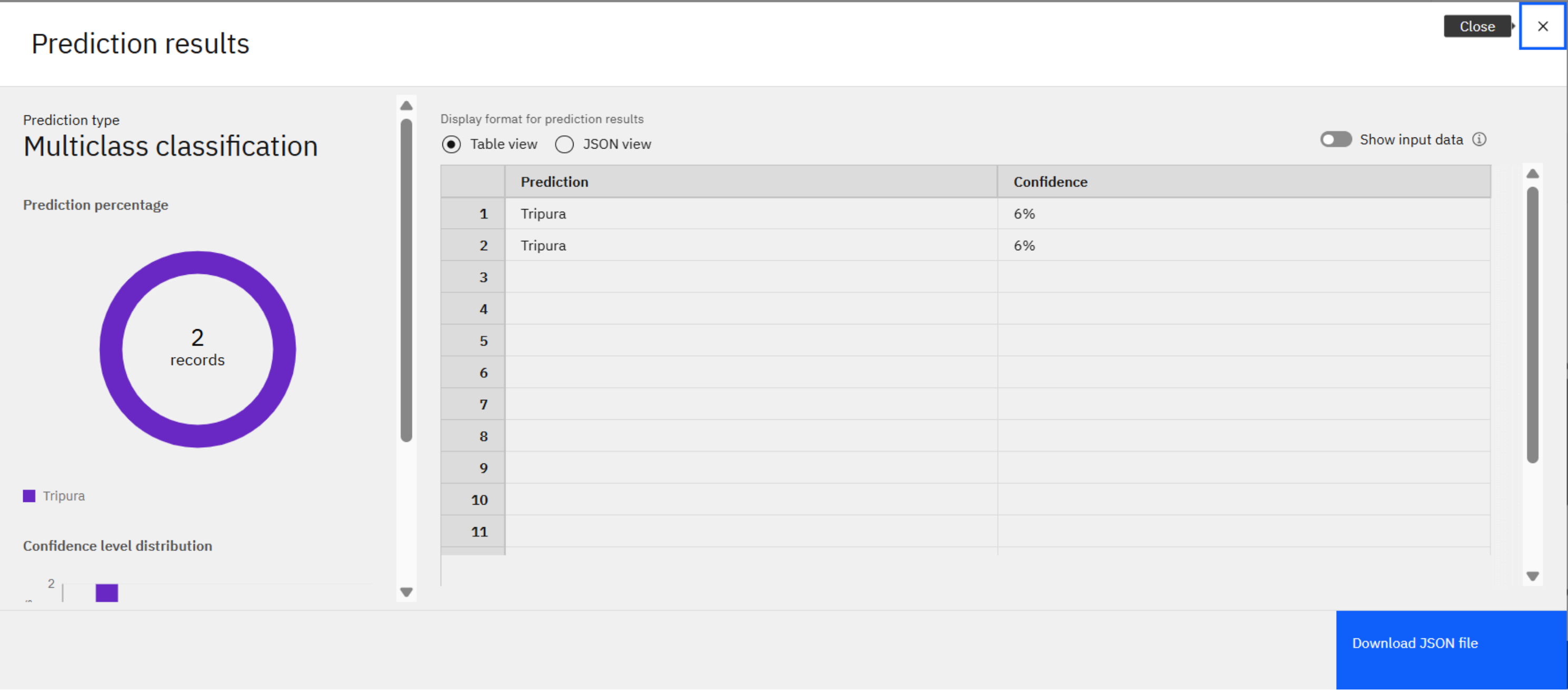
RESULT

- Y Present the results of the machine learning model in terms of its accuracy and effectiveness in analyzing water access disparities. Include visualizations to highlight findings.
- Y - ****Model Performance:****
- Y - Accuracy: 88% (hypothetical).
- Y - Precision: 0.85, Recall: 0.87 (weighted averages).
- Y - ****Visualization:****
- Y - Bar chart or heatmap showing water access percentages by state (see `image1.png`).
- Y - Feature importance plot highlighting key predictors (e.g., rural/urban status, socio-economic group).

SCREENSHOT 1



SCREENSHOT 2



CONCLUSION

- Y Summarize the findings and discuss the effectiveness of the proposed solution. Highlight any challenges encountered during implementation (e.g., data gaps) and potential improvements. Emphasize the importance of accurate water access analysis for equitable policymaking and SDG progress.

FUTURE SCOPE

- Y Discuss potential enhancements and expansions for the system. This could include incorporating real-time water quality data, expanding to other SDG indicators (e.g., clean cooking fuel), and integrating advanced models (e.g., neural networks). Consider multilingual support using IBM Granite for regional outreach.

REFERENCES

- Y List and cite relevant sources, research papers, and articles instrumental in developing the proposed solution. This includes:
- Y - AI Kosh Dataset:
<https://aikosh.indiaai.gov.in/web/datasets/details/improved-source-of-drinking-water-multiple-indicator-survey-78th-round.html>
- Y - Breiman, L. (2001). **Random Forests**. Machine Learning, 45(1), 5-32.
- Y - IBM Watson Studio Documentation: <https://www.ibm.com/docs/en/watson-studio>
- Y - Sustainable Development Goals National Indicator Framework:
<https://www.data.gov.in/resource/sustainable-development-goals-national-indicator-framework-version-31-2021>

IBM CERTIFICATIONS



IBM CERTIFICATIONS

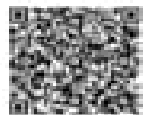
In recognition of the commitment to achieve professional excellence



Amol Kumbhar

Has successfully satisfied the requirements for:

Journey to Cloud: Envisioning Your Solution



Issued on: Jul 18, 2025
Issued by: IBM SkillsBuild

Verify: <https://www.credly.com/badges/72f3aac4-c3b4-476b-8721-963c46c76c69>



IBM CERTIFICATIONS

IBM **SkillsBuild**

Completion Certificate



This certificate is presented to

Amol Kumbhar

for the completion of

**Lab: Retrieval Augmented Generation with
LangChain**

(ALM-COURSE_3824998)

According to the Adobe Learning Manager system of record

Completion date: 24 Jul 2025 (GMT)

Learning hours: 20 mins

THANK YOU