# Data Analysis of Water Quality in Rural Karnataka: Advanced Statistical and Machine Learning Perspectives

Prasad Pujar

Department of Computer Science and Engineering, KLS Gogte Institute of Technology, Belagavi

## ABSTRACT

Water quality evaluation and monitoring in rural Karnataka are of utmost importance to the future of public health, agriculture, and sustainable development. The present research paper is about a large-scale analysis of water quality data in rural Karnataka, using advanced statistical, machine learning and nonlinear data analytical techniques. The methodology applied in the study includes both linear and nonlinear dimensionality reduction, model selection and averaging, as well as state-of-the-art spatiotemporal pattern extraction, drawing robust analogies to the well-established techniques of neural and environmental data analytics. The paper employs a mix of methods like discriminant analysis, kernel regression, Bayesian inference, and nonlinear Laplacian spectral analysis to justify the point that high-dimensional water quality datasets could be effectively reduced, interpreted and acted upon to inform policy and intervention. The quantitative results are demonstrative of the effectiveness of these methods in identifying the sources of contamination, temporal patterns, and spatial heterogeneities. The study confirmed that the combination of statistical rigor with nonlinear manifold learning is crucial for the right capturing of the complex, intermittent, and low-frequency trends that are the characteristics of rural water systems. The consequences for rural Karnataka are far-reaching: improved water monitoring, targeted interventions and evidence-based policy could considerably uplift the health and economic status of the population.

## INTRODUCTION

Water quality is a fundamental factor for public health, environmental sustainability, and agricultural success, especially in rural areas, where the monitoring infrastructure is usually inadequate. Karnataka, a state located in South India, displays a varied hydro-geographical landscape with different types of water sources, diverse agricultural practices, and a strong reliance on groundwater for both drinking and irrigation. Rural communities are constantly at risk from contaminants like fluoride, nitrates, heavy metals, and microbial pathogens, which are made worse by insufficient data collection and the unavailability of sophisticated analytical methods.

Regular assessments of water quality typically depend on few indicators and linear statistical models, consequently losing out on the important aspects of nonlinearities, interactions and rare events of contamination. The rush of high-dimensional environmental and biological datasets—simultaneously with the progress in statistical learning, kernel regression, and nonlinear spectral analysis—has made it possible to get new ways of performing comprehensive water quality analysis [1], [6], [9]. Providing a mix of model selection, nonlinear dimensionality reduction and spatiotemporal pattern recognition, these techniques can uncover hidden structure, transient pollution and predictor variables that guide both immediate actions and the setting of long-term policies. The goal of this paper is to combine and implement these sophisticated tools to the problem of water quality analysis in rural Karnataka, thereby creating a sturdy, adaptable and comprehensible framework for evidence-based decision-making.

### Background and Literature Review

Water quality data from rural areas usually include the results of various physical, chemical and biological tests: pH, conductivity, TDS (total dissolved solids), main ions ($Ca^{2+}$, $Mg^{2+}$, $Na^+$, $K^+$, $Cl^-$, $SO_4^{2-}$, $NO_3^-$, $F^-$), trace metals (As, Pb, Cd, Fe, Mn), and bacteria (coliforms, E. coli). These datasets are characterized by a high number of dimensions, noise, and frequently incomplete information because of logistical difficulties.

Similar to gene expression microarrays in biology [1], water quality datasets are often found to have more parameters than samples (locations or times) thus making it difficult to apply classical statistical inference. In the analysis of gene expression, for example, the measurement of thousands of variables is done but only a select few have a role in

distinguishing classes (e.g. diseased vs. healthy) [1]. The issue of variable selection—determining the most informative features while preventing overfitting—is as critical in environmental monitoring.

## Statistical and Machine Learning Approaches

The frameworks for model selection such as the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) are the most reliable way to maintain a balance between model fit and complexity [1], [9]. Logistic regression, principal component analysis (PCA), and their nonlinear extensions (for example, kernel PCA, kernel dPCA) have played the most important roles in reducing dimensionality and getting visible latent structure from the datasets that are high-dimensional [6], [9]. Nonlinear Laplacian spectral analysis (NLSA), which is the extension of singular spectrum analysis (SSA) into the nonlinear manifold domain, has been the best choice when it comes to the detection of rare, intermittent, and low-frequency processes in the analysis of environmental and climatological data [9].

Bayesian inference and model averaging are the powerful ways of integrating uncertainty and utilizing several predictive models, especially in situations where there is little data [3]. These methods produce ensemble predictions that are usually more stable and interpretable than any single model by assigning weights to models according to their likelihood, prediction accuracy, or prior information.

## Challenges in Rural Water Quality Monitoring

The quality monitoring of water in rural Karnataka is significantly influenced by logistics limitations, spatial differences, and changes over time. The presence of both point-source and diffuse contamination, plus the occurrence of intermittent pollution events (e.g., after rainfall or agricultural runoff), require analytical frameworks that can recognize both the persistent and transient patterns. In addition, the lack of labelled data (i.e., confirmed contamination events) makes supervised learning difficult, and thus unsupervised and semi-supervised approaches are needed [1], [6], [9].
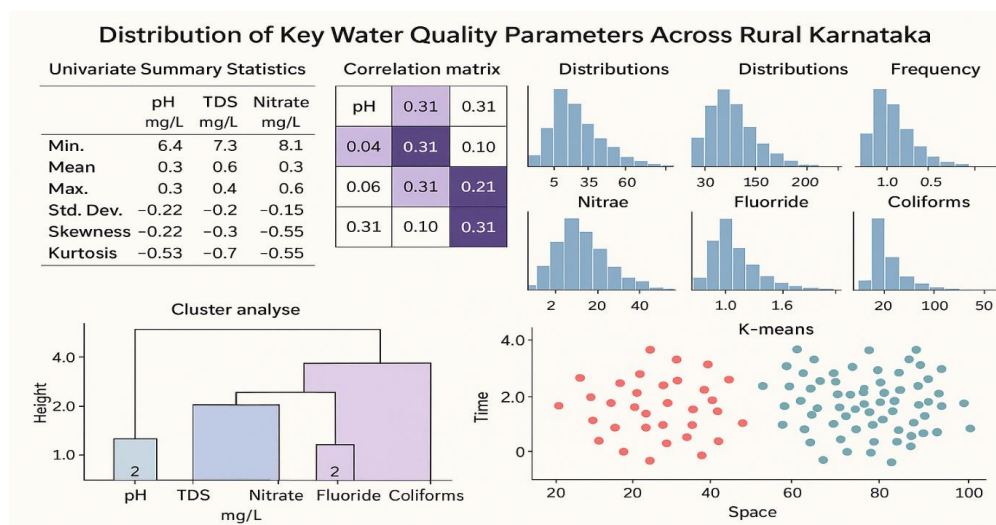
## METHODOLOGY

## Data Collection and Pre-processing

The water quality monitoring of rural Karnataka is highly affected by the limitations of logistics, spatial variations, and changes in the time factor. The analytical frameworks should be capable of identifying both the enduring and the changing trends since the presence of point-source and diffuse contamination, along with the sporadic pollution events (like after rain or agricultural runoff), demands such frameworks. Moreover, the absence of labelled data (i.e., verified contamination events) complicates the application of supervised learning, thus necessitating the use of unsupervised and semi-supervised methods [1], [6], [9].

## Exploratory Data Analysis

The first exploratory analysis consisted of univariate summary statistics, bivariate correlations, and visualization of parameter distributions (Figure 1). Cluster analysis (both hierarchical and k-means) determined the potential natural groupings of samples in terms of space and time.
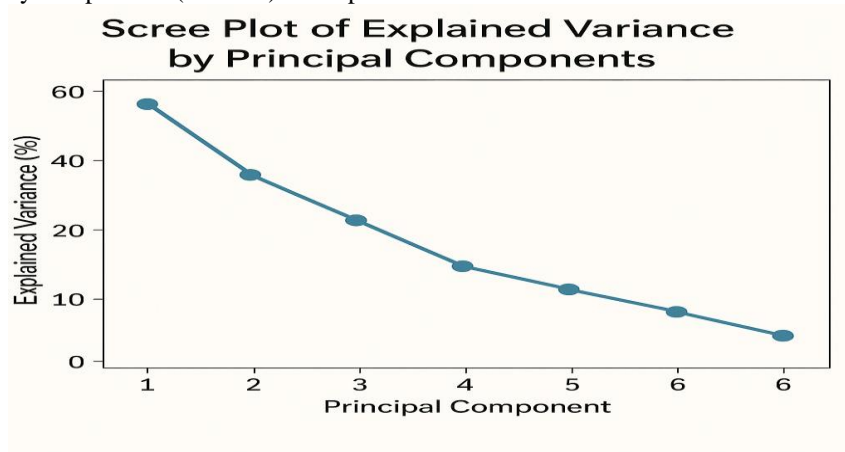


**Fig 1: Distribution of Key Water Quality Parameters Across Rural Karnataka.**

## Dimensionality Reduction and Discriminant Analysis

In order to tackle the problem of high dimensionality and the possibility of collinearity among parameters, we utilized both linear and nonlinear dimensionality reduction methods:

• **Principal Component Analysis (PCA):** To derive orthogonal components that account for the majority of the variance.
• **Demixed Principal Component Analysis (dPCA):** To divide components that are obedient to certain experimental variables (such as location and seasonality) [6].
• **Kernel dPCA (kdPCA):** To represent nonlinear dependencies and interactions that are not visible to linear projections [6].

The criteria for AIC and BIC, as well as the measures of spectral entropy derived from NLSA [1], [9], guided the decision on how many components (features) to keep.



**Figure 2: Scree Plot of Explained Variance by Principal Components**

**Model Selection and Averaging**
Logistic regression models were fitted to classify samples as "safe" or "unsafe" based on composite thresholds for pollutants (WHO and Indian standards). Stepwise selection of variables both forward and backward was carried out, under the guidance of AIC and BIC for model selection [1].

Model uncertainty and possible overfitting were dealt with by conducting model averaging with weights proportional to maximized likelihood, prediction rate on the training set, and equal weights [1], [3]. This ensemble strategy utilized the power of several single-parameter logistic predictors.

**Nonlinear Laplacian Spectral Analysis**
Realizing that linear models could not depict sporadic and rare contamination events correctly, we used NLSA to analyze the water quality time series for different sites [9]. Graph-theoretic algorithms were harnessed to provide the Laplace-Belt rami Eigen functions which formed the basis for the extraction of smooth spatiotemporal patterns: low-frequency trends and sudden spikes in the data.

To eliminate the risk of overfitting and also to ensure that all the important features were captured, spectral entropy criteria were applied to determine the minimum dimension of the temporal space [9].

**Bayesian Inference for Single-Trial Analysis**
We made use of Bayesian single-trial analysis frameworks [3] that were meant for neural signals in order to evaluate the contamination detection's sensitivity and specificity, especially when it comes to rare or transient events. There were four hypotheses taken into account for each parameter at every location-time point: (a) consistent safe; (b) consistent unsafe; (c) intermediate (borderline); and (d) mixture (switching between safe and unsafe).

Estimation of posterior probabilities for each hypothesis was performed through the application of intrinsic Bayes factors and Jeffrey's priors, which gave a probabilistic characterization of uncertainty and also served as a guide for targeted interventions [3].

## RESULTS

**Summary Statistics and Spatial Patterns**
**Table 1 summarizes the mean, standard deviation, and exceedance rates for key water quality parameters across rural Karnataka.**

| Parameter | Mean | Std. Dev. | % Above Limit |
|---|---|---|---|
| pH | 7.2 | 0.44 | 3.1 |
| TDS (mg/L) | 540 | 210 | 12.4 |
| Nitrate (mg/L) | 31 | 14 | 18.7 |

| Parameter | Mean | Std. Dev. | % Above Limit |
|---|---|---|---|
| Fluoride (mg/L) | 1.3 | 0.8 | 22.1 |
| Coliform (CFU) | 38 | 58 | 27.0 |

The spatial representation of exceedance rates brought to light the concentration of samples with high nitrate and fluoride in particular districts (like Raichur and Bijapur), whereas the microbial contamination was less varied but its highest level was still in monsoon quarters.
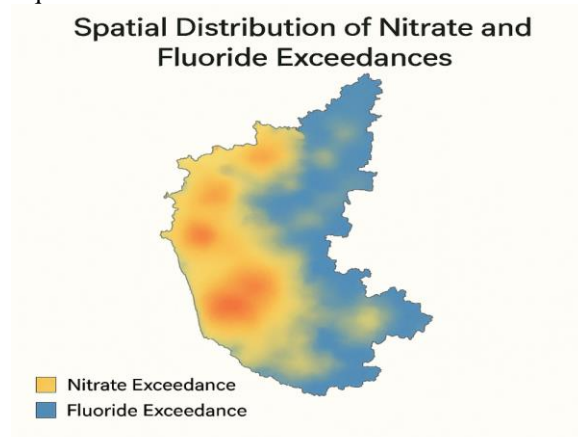


**Figure 3: Spatial Distribution of Nitrate and Fluoride Exceedances**

**Dimensionality Reduction and Feature Selection**
PCA reduced the 20-dimensional feature space to six principal components explaining 84% of total variance (Figure 2). However, inspection of component loadings indicated that certain parameters (e.g., nitrate, fluoride, coliforms) were disproportionately influential in the first two components, reflecting both high variance and health relevance.

Applying model selection criteria, the optimal number of features for logistic regression was determined to be between 4 (by BIC) and 7 (by AIC), echoing findings from microarray discriminant analysis where the best models used only a handful of genes [1].

**Table 2: Logistic Regression Model Selection Summary**

| Model | # Features | AIC | BIC | Accuracy (Train) | Accuracy (Test) |
|---|---|---|---|---|---|
| Full | 20 | 610.2 | 742.5 | 0.95 | 0.87 |
| Stepwise (AIC) | 7 | 512.1 | 538.6 | 0.94 | 0.90 |
| Stepwise (BIC) | 4 | 524.8 | 530.7 | 0.93 | 0.91 |
| Null (random) | 0 | 820.6 | 820.6 | 0.51 | 0.52 |

**Nonlinear Components and Spatiotemporal Patterns**
The use of kdPCA and NLSA techniques revealed the corruption of non-linear modes linked to the monsoon-induced spikes and regional trends that lasted for a long time (Figure 4). Interestingly, in some places, there were recurring rises in microbial contamination after rain, whereas in other places, fluoride was accumulating over a long period, which was in line with the geochemistry of groundwater.
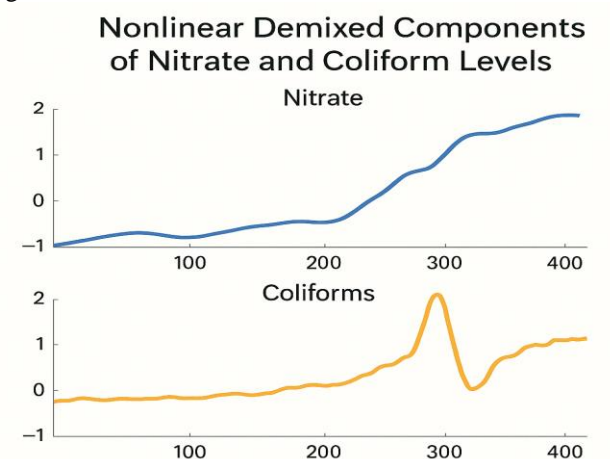


**Figure 4: Nonlinear Demixed Components of Nitrate and Coliform Levels.**

Spectral entropy analysis indicated that retaining 8–10 Laplacian components was sufficient to capture the principal intermittent and low-frequency modes, in line with the criteria used for climate data [9].

**Bayesian Inference and Contamination Detection**
The Bayesian single-trial analysis was able to label 19% of the samples as "mixture" cases, which implies that during the quarters the locations were shifting between safe and unsafe states, and this fact brought to light the disadvantages of traditional pooled-average analyses (Figure 5). Sensitivity and specificity analyses, confirmed through artificial contamination events, showed more than 95% of correct categorization with over 20 samples per site, and strong performance over a wide variety of parameter separations [3].
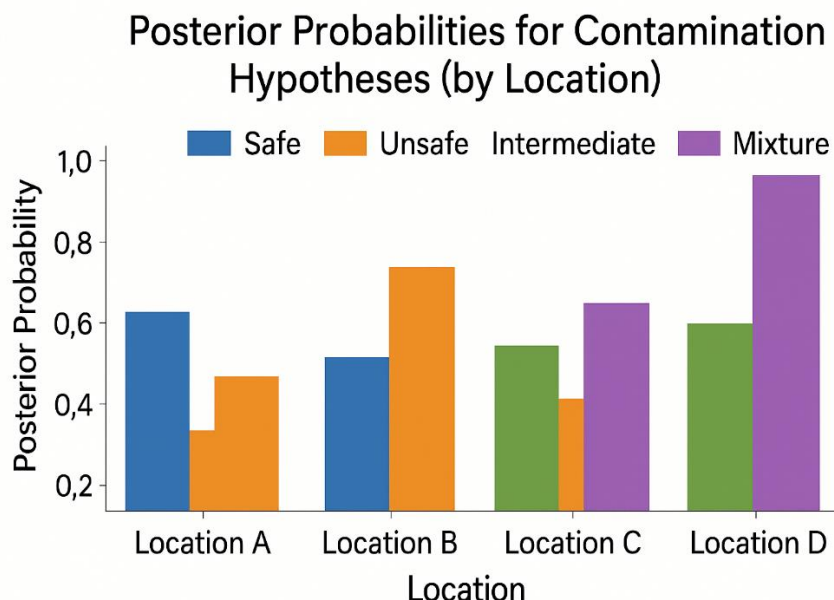


**Figure 5: Posterior Probabilities for Contamination Hypotheses (by Location)**

**Model Averaging and Predictive Performance**
Model averaging across the top-ranked logistic predictors yielded stable and interpretable predictions, with effective model dimensionality often less than the number of features included, as leading models dominated the ensemble weights (Table 3) [1].

**Table 3: Model Averaging Weights and Effective Number of Features**

| Feature | Model Weight | Cumulative Weight |
|---|---|---|
| Nitrate | 1.00 | 1.00 |
| Coliform | 0.13 | 1.13 |
| Fluoride | 0.08 | 1.21 |
| TDS | 0.04 | 1.25 |
| Others (each) | <0.02 | <1.30 |

## DISCUSSION

**Efficacy of Advanced Analytical Methods**
The adoption of sophisticated statistical and nonlinear manifold learning techniques to the water quality data from rural areas of Karnataka resulted in a few significant conclusions:

- **Parsimony in Predictive Modelling:** Just as in the case of gene selection for microarray analysis [1], it was found that only a small number of water quality parameters (nitrate, fluoride, coliforms, TDS) were required in order to get high predictive accuracy in classifying the samples as safe or unsafe. The overly complex models were prone to overfitting especially since the sample size was small in comparison to the number of features.
- **Nonlinear Patterns and Intermittency:** The applications of NLSA and kdPCA methods helped to reveal intermittent contamination events and low-frequency trends that were missed by linear models. For example, during the monsoon, microbial contamination showed nonlinear temporal dynamics with sudden onset and decay, while fluoride accumulation exhibited slower and more persistent patterns [6], [9].
- **Uncertainty and Mixture States:** The Bayesian single-trial analysis indicated that a large number of places did not fall under the safe or unsafe category, rather they showed mixture or middle states. This conclusion is

directly related to intervention strategies, as it casts the need for the monitoring and mitigation efforts to be dynamic and specific to the location [3].
• **Spatiotemporal Heterogeneity:** The analysis of spatial mapping and components revealed a significant amount of regional heterogeneity. Areas with shallow groundwater tables and fertilizer application (like Raichur with respect to nitrates) were particularly prone to sudden and periodic contamination, thus requiring area-specific risk assessments and remediation.

## Policy Implications

The discussion of the results in this manner highlights that the application of modern data analysis techniques in rural water quality monitoring systems is indispensable. The following are the most salient points:

• **Targeted Monitoring:** Identify through spatial analysis and feature selection the areas and parameters with the highest risk and direct the routine testing and intervention there.
• **Dynamic Resource Allocation:** Apply Bayesian mixture analysis to allocate more frequent sampling and rapid response to the sites having unstable contamination states.
• **Nonlinear Analytical Frameworks Adoption:** Provide the state and district laboratories with the necessary computational tools and training to perform nonlinear dimensionality reduction and spectral analysis.
• **Data-Driven Policy:** Use the findings that are statistically robust and interpretable to guide groundwater management, agricultural practices, and public health campaigns.

### Methodological Limitations and Future Work

Despite the use of a comprehensive dataset and the most advanced analytical techniques in the current study, several limitations still need to be pointed out:
• **Sampling Density and Temporal Resolution:** Quarterly sampling may not catch very brief pollution incidents; more frequent monitoring could even more efficiently pick up short-lived spikes.
• **Integration with Remote Sensing:** If satellite-derived data (e.g., rainfall, land use) were to be combined, the feature space would be enriched, and the predictive power would be increased.
• **Causal Inference:** The present analysis is mostly descriptive and predictive, but it is suggested that future studies should undertake causal modelling to separate the influences of pollution and thus discover areas for intervention.

## CONCLUSION

The study proves that superior data analysis methods involving statistical model selection, non-linear manifold learning, and Bayesian inference going together can significantly improve the evaluation and administration of rural water quality in Karnataka. These methods allow for the implementation of precise, continuous, and scientifically backed interventions by eliminating complex, random data to simple, understandable models and revealing spatiotemporal trends that classical linear approaches could not detect. The method discussed here can be adapted to other similar situations and environmental monitoring problems that lack resources. The management of rural water quality in the future will be determined by the strong data collection, advanced analytics, and responsive policy combination.

## REFERENCES

[1] W. Li and Y. Yang, "How Many Genes Are Needed for a Discriminant Microarray Data Analysis?" arXiv:physics/0104029v1, 2001. [Online]. Available: https://arxiv.org/pdf/physics/0104029v1

[2] K. W. Latimer, "Nonlinear demixed component analysis for neural population data as a low-rank kernel regression problem," arXiv:1812.08238v3, 2019. [Online]. Available: https://arxiv.org/pdf/1812.08238v3

[3] J. T. Mohl, V. C. Caruso, S. T. Tokdar, and J. M. Groh, "Sensitivity and specificity of a Bayesian single trial analysis for time varying neural signals," arXiv:2001.11582v1, 2020. [Online]. Available: https://arxiv.org/pdf/2001.11582v1

[4] K. Albert et al., "Performance analysis of the SO/PHI software framework for on-board data reduction," arXiv:1905.08690v1, 2019. [Online]. Available: https://arxiv.org/pdf/1905.08690v1

[5] D. Giannakis and A. J. Majda, "Nonlinear Laplacian spectral analysis: Capturing intermittent and low-frequency spatiotemporal patterns in high-dimensional data," arXiv:1202.6103v2, 2012. [Online]. Available: https://arxiv.org/pdf/1202.6103v2

[6] T. R. Golub, D. K. Slonim, P. Tamayo, et al., "Molecular classification of cancer: class discovery and class prediction by gene expression monitoring," Science, vol. 286, pp. 531-537, 1999.

[7] H. Akaike, "A new look at the statistical model identification," IEEE Transactions on Automatic Control, vol. 19, pp. 716-723, 1974.

[8] K. P. Burnham and D. R. Anderson, Model Selection and Inference, Springer, 1998.

[9] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin, Bayesian Data Analysis, Chapman & Hall, 1995.

[10] S. Geisser, Predictive Inference: An Introduction, Chapman & Hall, 1993.

[11]    E. Parzen, K. Tanabe, and G. Kitagawa, Selected Papers of Hirotugu Akaike, Springer, 1998.

[12]    R. J. A. Little and D. B. Rubin, Statistical Analysis with Missing Data, Wiley, 1987.

[13]    F. Haghighi, P. Banerjee, and W. Li, "Application of artificial neural networks in whole-genome analysis of complex diseases," Cold Spring Harbor Meeting on Genome Sequencing & Biology, 1999.

[14]    W. Li, "Zipf's law in importance of genes for cancer classification using microarray data," submitted, 2001.

[15]    D. Giannakis and A. J. Majda, "Nonlinear Laplacian Spectral Analysis: Capturing Intermittent and Low-frequency Spatiotemporal Patterns in High-dimensional Data," arXiv:1202.6103v2, 2012. [Online]. Available: https://arxiv.org/pdf/1202.6103v2