

Classification of Water Potability

LIYA LAPIERRE, PRATEEK KAKKAR, NITHARSAN SIVAKANTHAN



Background

Drinkable water is essential to human health

Interested in finding a series of tests to use to predict Potability (safe drinking water)

If possible, find the least number of tests needed to improve accessibility

The Data

Water Quality Dataset (Posted on Kaggle)

Metrics for ~ 3300 bodies of water

Predictors: pH level, Hardness, Solids, Chloramines, Organic Carbon, Trihalomethanes, Turbidity

Response: Potability

■ ~ 2000 bodies of water are not Potable, ~1300 are Potable

Methodology

Data Cleaning

- SVD Imputation

Data Exploration

- Principal Component Analysis

Modeling

- Logistic Regression
- Random Forest Classifier
- K-Nearest Neighbor Classifier

Data Cleaning

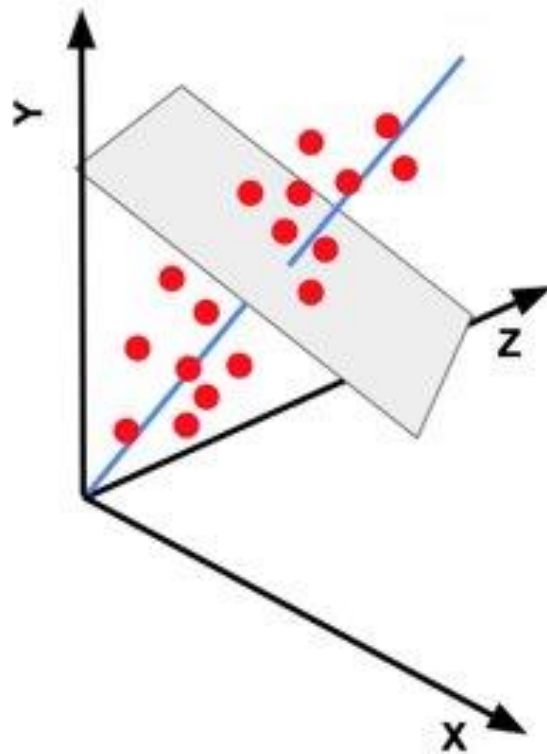
➤ Singular Value Decomposition Imputation

$$\begin{array}{c} \mathbf{A} \\ \left(\begin{array}{ccc} x_{11} & x_{12} & x_{1n} \\ & \ddots & \\ x_{m1} & & x_{mn} \end{array} \right) \\ m \times n \end{array} = \begin{array}{c} \mathbf{U} \\ \left(\begin{array}{ccc} u_{11} & & u_{m1} \\ & \ddots & \\ u_{1m} & & u_{mm} \end{array} \right) \\ m \times m \end{array} \begin{array}{c} \mathbf{S} \\ \left(\begin{array}{ccc} \sigma_1 & & 0 \\ & \ddots & \\ 0 & \sigma_r & \\ & & \ddots & \\ & & 0 & \end{array} \right) \\ m \times n \end{array} \begin{array}{c} \mathbf{V}^T \\ \left(\begin{array}{ccc} v_{11} & & v_{1n} \\ & \ddots & \\ v_{n1} & & v_{nn} \end{array} \right) \\ n \times n \end{array}$$

Source: (Hui, 2019)

Data Exploration

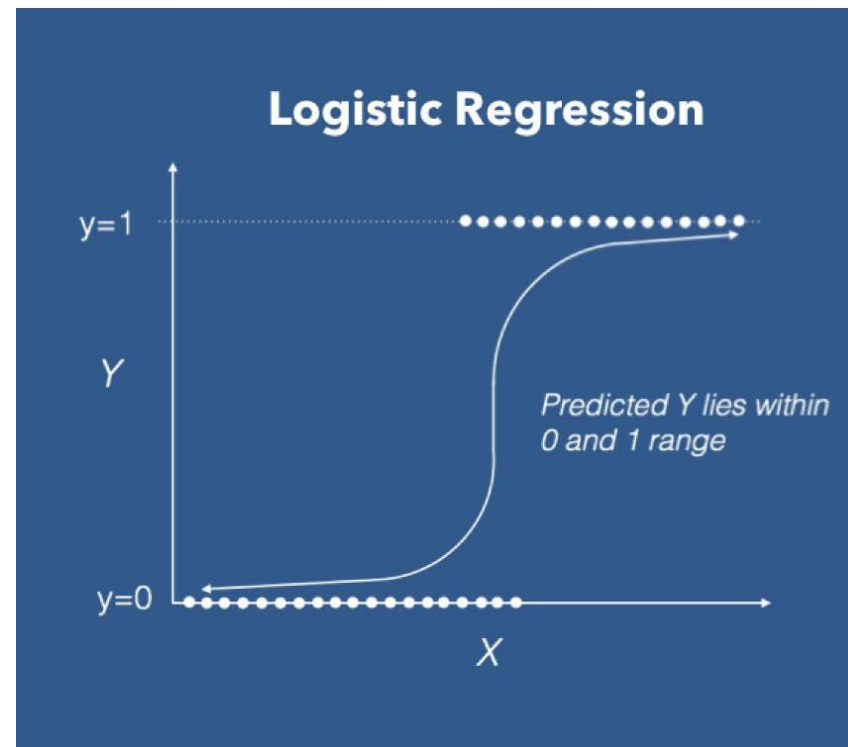
➤ Principal Component Analysis



Source: (LearnOpenCV)

Classification Methods

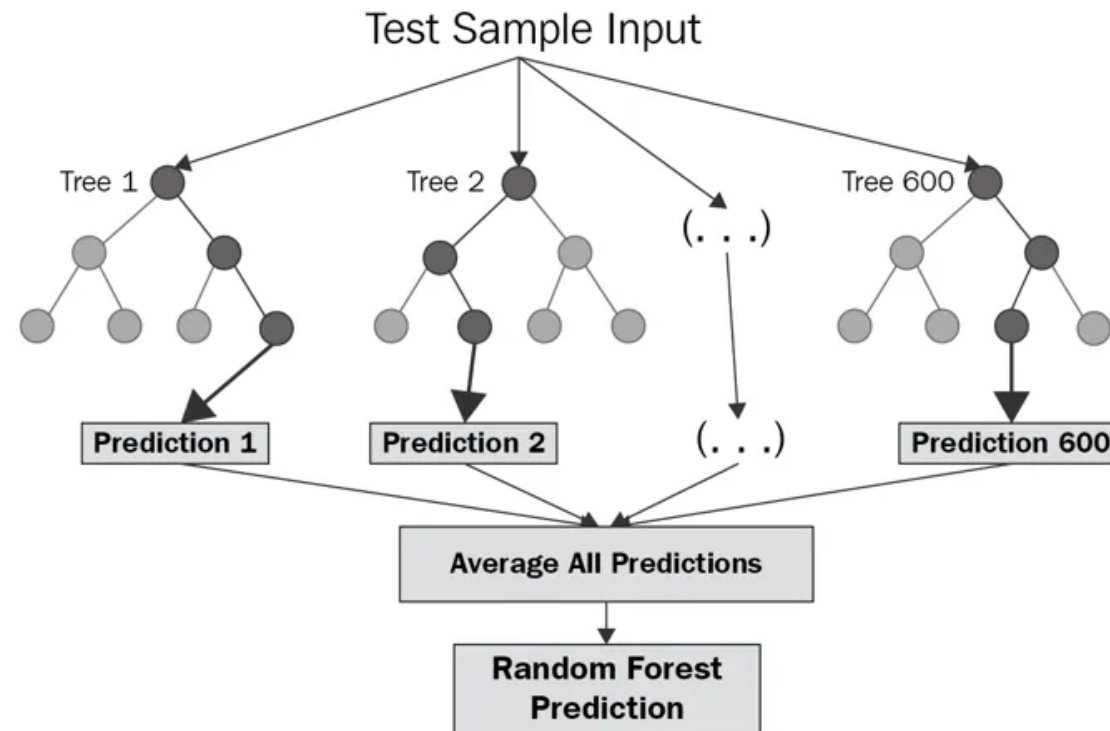
➤ Logistic Regression



Source: (Rajput, 2018)

Classification Methods

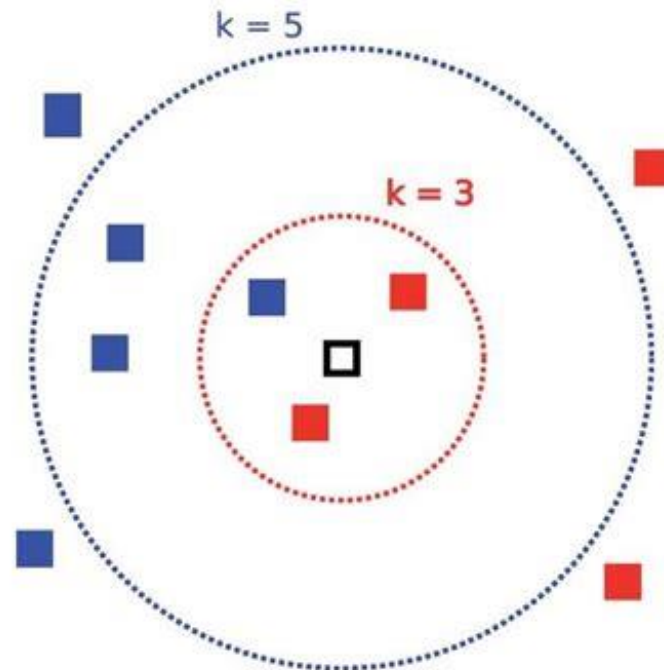
➤ Random Forest Classifier



Source: (corporatefinanceinstitute)

Classification Methods

➤ K Nearest Neighbors



Source: (ResearchGate)

Results

Logistic Regression

- Accuracy: 61%
- Recall: 0%
- Precision: 0%
- Specificity: 99%

Potable?	Yes	No
P- Yes	0	1
P- No	255	400

Random Forest

- Accuracy: 68%
- Recall: 35%
- Precision: 66%
- Specificity: 89%

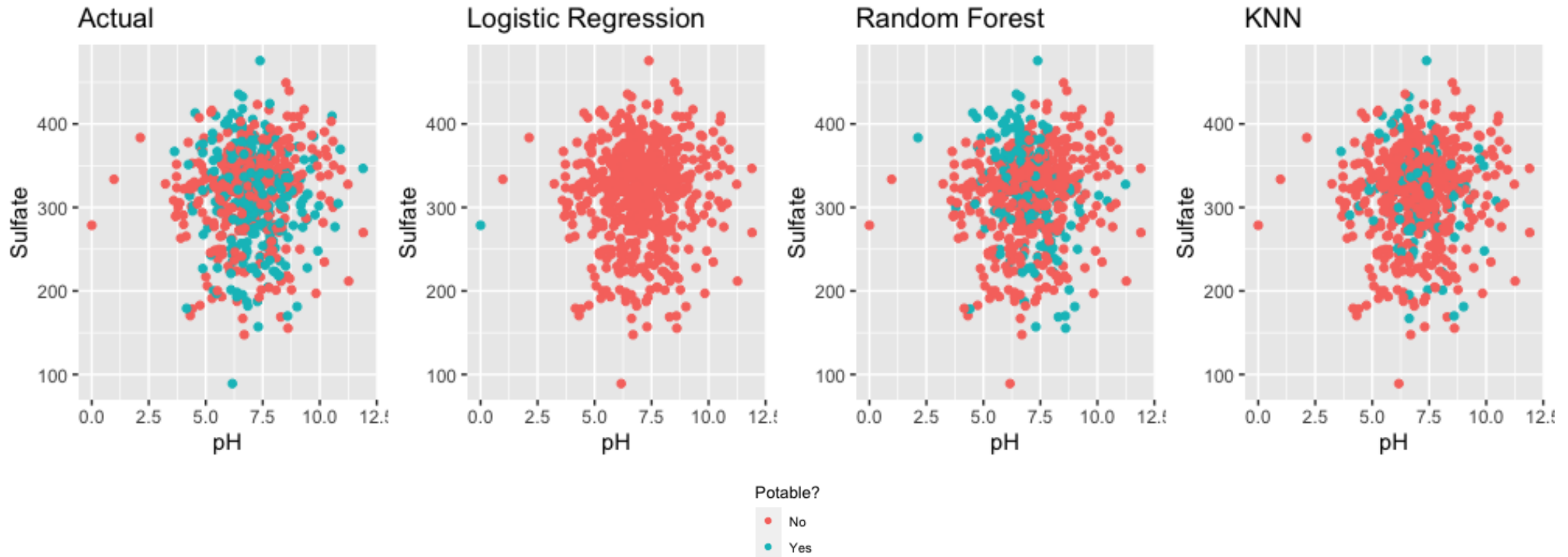
Potable?	Yes	No
P- Yes	89	46
P- No	166	355

KNN

- Accuracy: 65%
- Recall: 25%
- Precision: 62%
- Specificity: 91%

Potable?	Yes	No
P- Yes	63	38
P- No	192	363

Results Plotted



Conclusions

- We can predict, with moderate accuracy, the potability of water using several machine learning algorithms
- KNN models were very accurate at identifying the samples that were not potable
- The models with only three predictor variables (pH, sulfate, chloramines) were still moderately successful

Future Considerations

- Bacterial levels are an important indicator of potability
- Consider using undersampling or oversampling due to unbalanced data
- Tune model based on other metrics

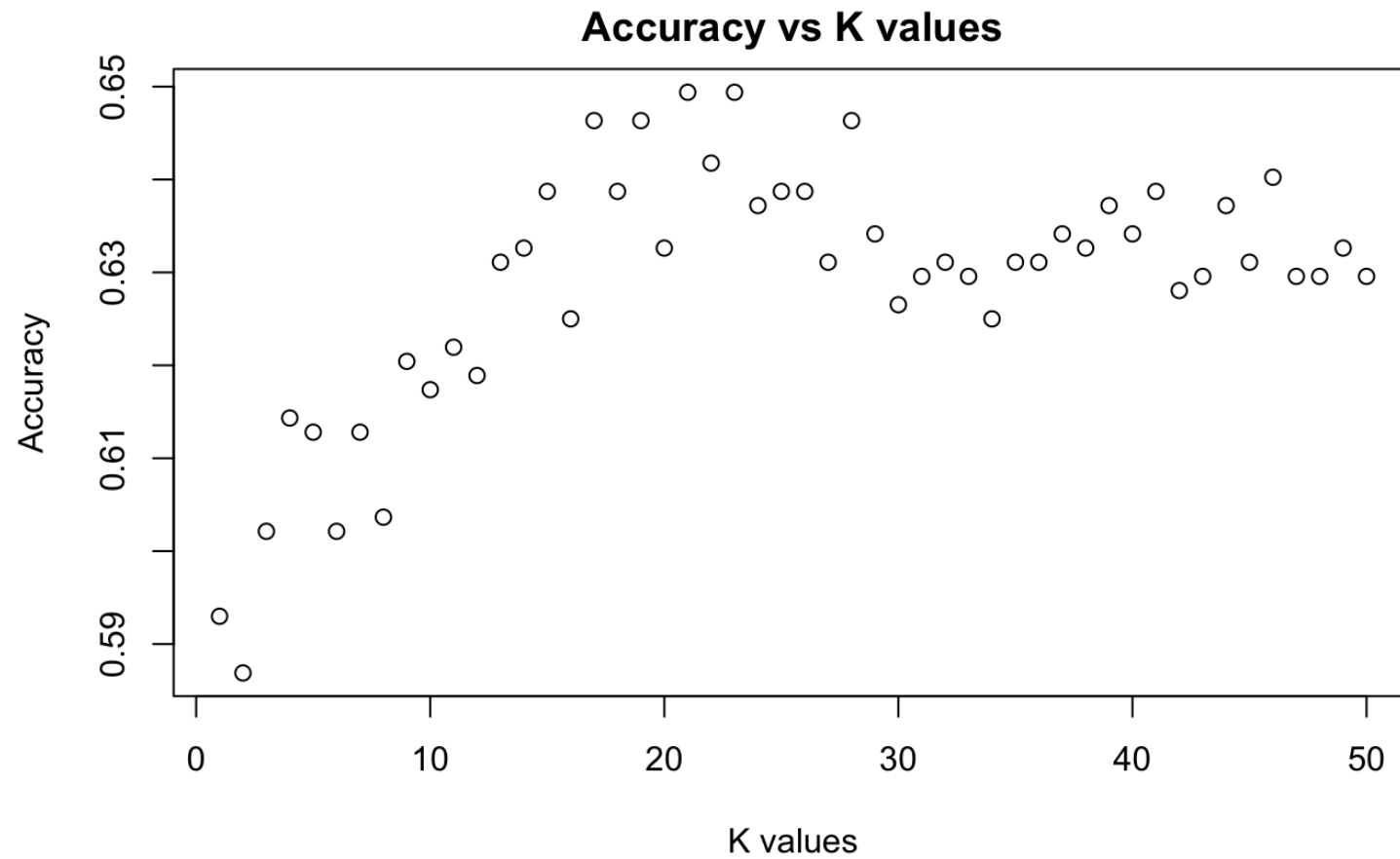
References

Environmental Protection Agency. (n.d.) National Primary Drinking Water Regulations. *EPA.gov*.
https://www.epa.gov/sites/default/files/2016-06/documents/npwdr_complete_table.pdf

Kadiwal, Aditya. (n.d.) Water Quality. *Kaggle*. <https://www.kaggle.com/datasets/adityakadiwal/water-potability?resource=download>

Appendix

Tuning KNN



Random Forest w/ Less predictors

CONFUSION MATRIX

		Actual	
		No	Yes
Predicted	No	288	164
	Yes	113	91

DETAILS

Sensitivity 0.357	Specificity 0.718	Precision 0.446	Recall 0.357	F1 0.397
	Accuracy 0.578		Kappa 0.078	

KNN w/ Less predictors

CONFUSION MATRIX

		Actual	
		No	Yes
Predicted	No	360	185
	Yes	41	70

DETAILS

Sensitivity 0.275	Specificity 0.898	Precision 0.631	Recall 0.275	F1 0.383
Accuracy 0.655		Kappa 0.192		