# Classifying Safe Water Sources

Potential Potables

By: William Grennan

# Earth's Most Valuable Resource

- Fresh Water accounts for only ~2.5% of all water on earth

- But not all fresh water is safe to drink

- Water Pollution takes many forms
  - pH imbalance
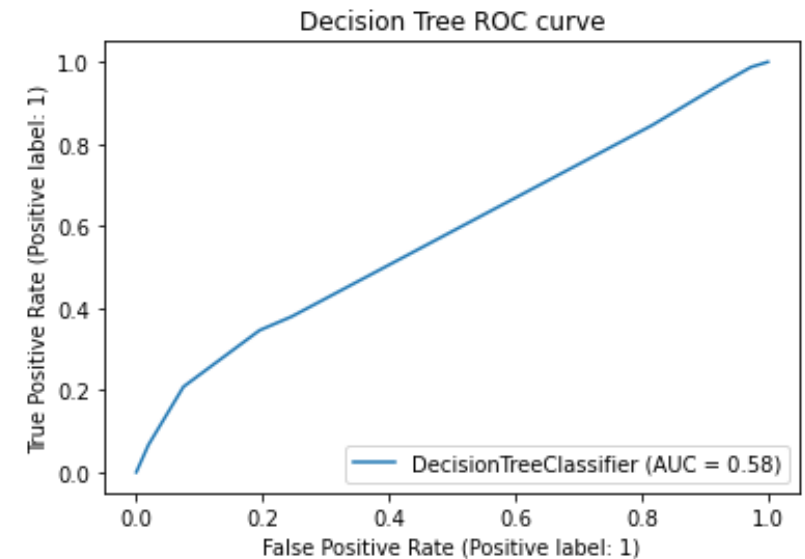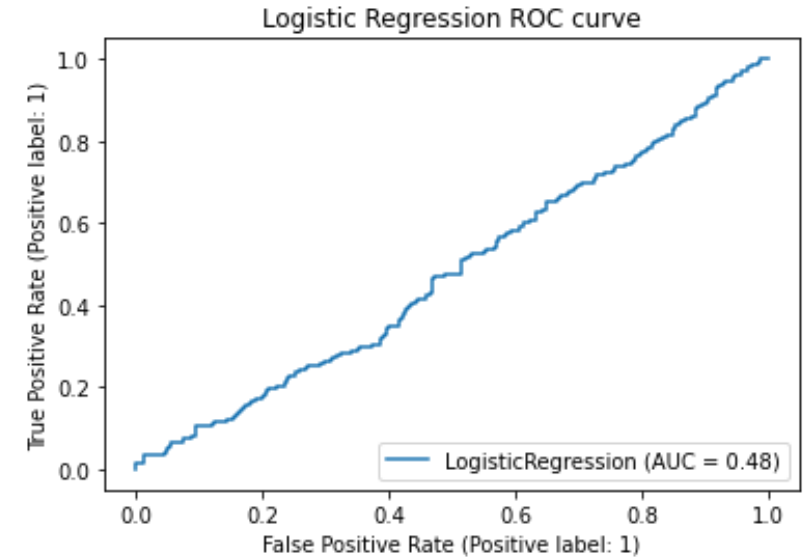  - Clarity
  - Toxic substances
  - Etc.

# How Can Data Science Help?

- 9 key measurements can be taken from the water sources quickly
  - pH, Hardness, Dissolved Solids, Chloramines, Dissolved Sulfates, Conductivity, Organic Carbon, Trihalomethanes, and Turbidity
- Using these measurements, we can narrow the search
- 2 Objectives depending on need
  - **Precision – Quickly find water most likely to be safe**
  - Recall – Identify sources that were previously overlooked

# What Successes Can We Find?

- Data modeling can be tricky

- Lots of measurement overlap between the classes

- Decision Trees and Forests tend to perform better than Logistic Regression

# Directly Compare the Models in Streamlit

- Selectable Model
- View Model Hyperparameters
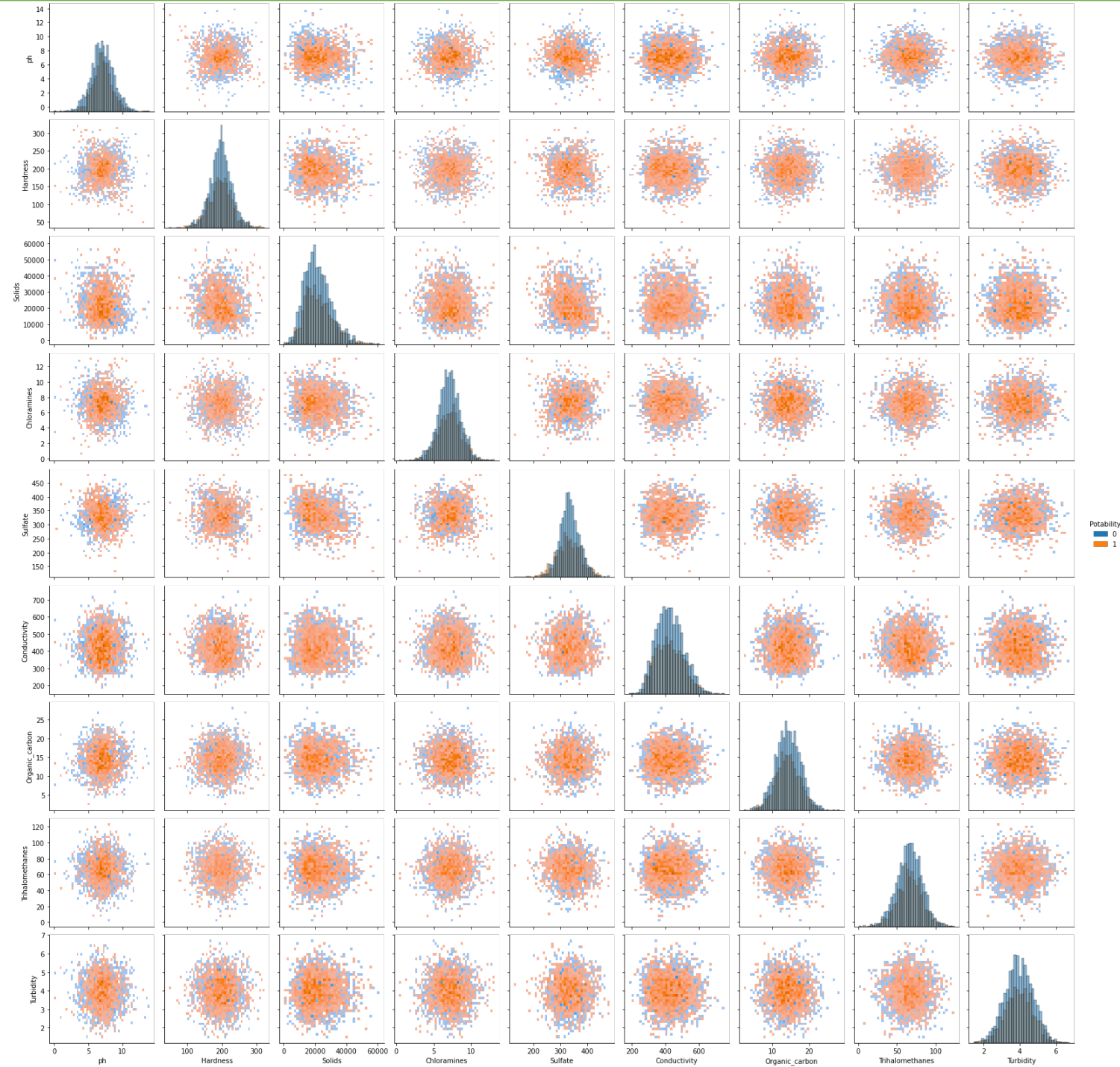- Viewable ROC Chart
- Adjustable Confusion Matrix

# Conclusions

- XGBoost and Random Forests perform the best on test data
- Current modeling is too risky for a definitive answer
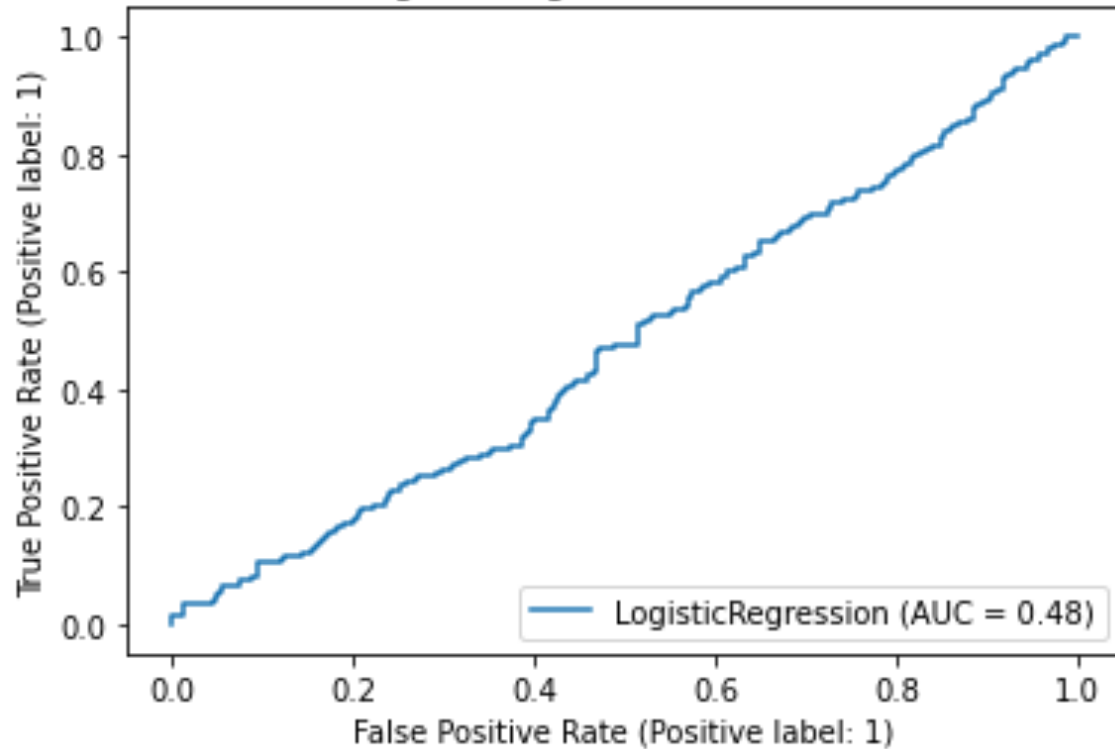- Further scientific study may find more reliable measurements

| Testing Data | Accuracy | ROC AUC |
|---|---|---|
| XGBoost | 61.28 % | 57.25 % |
| KN Neighbor | 60.82 % | 56.26 % |
| Random Forest | 64.02 % | 60.82 % |

# Appendix (Pair Plot)

# Appendix (ROC Curves)

# Appendix (ROC Curves)

# Appendix (ROC Curves)

# Appendix (ROC Curves)