# DS 2010 Final Report

Abhiram Yammanuru, Michael Herchenroder, Ryan Offstein, William Abraham, Ajay Bhagavatula

## About Dataset & Motivation:

Our group decided to determine what factor would have the biggest impact on water potability (drinkability) for humans. Our main motivation behind this was to see what features people would have to look into when determining the drinkability of water. We also looked at this from a humanitarian perspective, as drinkable water is an essential human need, and access to drinking water is a serious issue faced in many areas of the world.

We found a dataset with different variables that each could have a potential impact on the drinkability. The specific variables found in the dataset that impact the potability are Hardness, Solids, Chloramines, Sulfate, Conductivity, Organic Carbon, Trihalomethanes, and Turbidity. With all these variables that could potentially impact the potability, we decided to narrow down on a few for our conjectures. We specifically wanted to compare these variables to each other to see the best predictor of potability. The 3 conjectures/hypotheses that we came up with were that Turbidity is a better predictor of potability than solids, Ph is a better predictor of potability than conductivity, and Hardness is a better predictor of potability than Organic carbon. They can be seen listed below the dataset.

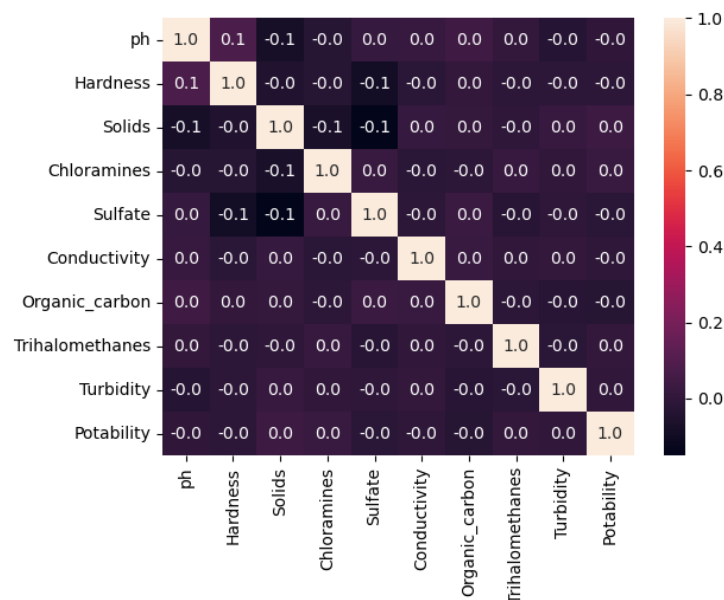| # ph | # Hardness | # Solids | # Chloramines | # Sulfate | # Conductivi... | # Organic_c... | # Trihalomet... | # Turbidity | # Potability |
|---|---|---|---|---|---|---|---|---|---|
| 9.406325752898285 | 216.76215866041707 | 27948.589445431105 | 6.1561107940777742 | 355.47303468825606 | 347.98307827691712 | 16.340715053314568 | 32.108292482849905 | 3.0972779574771165 | 0 |
| 6.32125931522299 | 207.25770955362344 | 8532.1395169595518 | 5.9878765835664684 | 286.48928042140803 | 491.76531319496039 | 10.546886238857326 | 74.502807995095661 | 4.5014572038139855 | 0 |
| | 217.36974600866085 | 17984.327439446446 | 8.5941628848905 5 | | 409.22077647853 76 | 10.213779294934 625 | 18.400012185523 245 | 3.6051541988744 72 | 0 |
| 4.7053564957747 | 103.17358697810718 | 19555.765051566 86 | 6.7672979897863685 | | 370.17816638462 43 | 9.1828336517958 18 | 93.900485873656 55 | 3.4374501340171 87 | 0 |
| 8.896418988301798 | 222.25629288902002 | 8870.9029659422 87 | 6.0113424330021 07 | 332.00263040579 02 | 425.22686170798 15 | 10.847736847274 12 | | 3.7009457327147 453 | 0 |
| 6.581878201548969 | 272.98274466101 01 | 37169.444403538 49 | 8.1147310153849 9 | 416.08348053859 96 | 351.47683939412 354 | 15.129334487820 7 | 79.261026496244 4 | 4.2016628644658 36 | 0 |
| 6.7551459015516 25 | 231.26013129603 174 | 18536.698647477 02 | 8.7571331786021 94 | 342.54801420266 62 | 385.11464773745 956 | 13.888834329605 06 | 79.302435726766 2 | 5.1627297493974 85 | 0 |
| 9.4451298378686 56 | 145.80540244684 383 | 13168.529155675 998 | 9.4444710856229 4 | 310.58337385859 07 | 592.65902097595 07 | 8.6063967469869 45 | 77.577459510356 467 | 3.8751652466165 467 | 1 |
| 9.0248450374175 03 | 128.09669121000 72 | 19859.676475803 88 | 8.0164226491737 37 | 300.15037702033 186 | 451.14348100565 14 | 14.770862942397 969 | 73.778025645975 86 | 3.9852505057435 756 | 1 |
| | 169.97484895701 46 | 23403.637304371 39 | 8.5197298510902 08 | | 475.57356242744 93 | 12.924106818027 212 | 50.861912984034 42 | 2.7473129992921 272 | 1 |
| 6.8001190903158 78 | 242.00808150751 49 | 39143.403328810 09 | 9.5016945877152 7 | 187.17071436243 93 | 376.45659307467 866 | 11.432466347228 74 | 73.777275026262 6 | 3.8549398997210 73 | 1 |
| 7.1741351628079 95 | 203.40893462062 238 | 20401.102461471 397 | 7.6818062872436 73 | 287.08567912256 15 | 315.54990001949 35 | 14.533510036354 244 | 74.405616026957 02 | 3.9398956565986 84 | 1 |
| 7.6579912369982 58 | 236.96088924616 282 | 14245.789121299 43 | 6.2890648594339 3 | 373.16536281008 53 | 416.62418891074 776 | 10.464238582078 508 | 85.852768605027 65 | 2.4372962875595 79 | 1 |
| 8.3229866724022 08 | 207.25246223156 424 | 28049.646283166 227 | 8.2270612831896 18 | 297.81308453289 102 | 358.72586877763 8 | 18.709273368730 52 | 60.911420394398 27 | 4.0521357275526 61 | 1 |

**Conjectures:**

1. **Ph is a better predictor of potability than conductivity**
2. **Turbidity is a better predictor of potability than solids**
3. **Hardness is a better predictor of potability than Organic carbon**

We hypothesized that Turbidity, Ph, and Hardness were going to be the best predictors of potability from our data set so our conjectures are comparing those features, with other features in the dataset to back up that hypothesis.

The features we are comparing are relatively similar too so we can see which is the better predictor. For example, Turbidity is the opaqueness and thickness of the water, whereas solids is the mass of the suspended solids in the water. Opaqueness and thickness of the water are strong indicators of solids present, so they are similar in that way. Ph and conductivity are related in the sense that they take into consideration the $H^+$ ions in the water. Lower Ph's correlate with lower conductivity and vice versa. The hardness of water is related to Organic Carbon because $CO_2$ induces the dissolution of $CaCO_3$ in the reaction with water, it can contribute to an increase in the hardness of water.

Even though these features seemed to have an impact on each other due to similarities, we found that none of these features have any significant correlation with each other. As seen below, the correlation matrix we ran on each of the datasets shows that none of these features significantly impact each other. This worked to our benefit as we didn't want to have confounding variables, which would have made it a lot more challenging to see which features had the largest impact on potability.

# Data Cleaning methods:

When first taking a look at our dataset, there were numerous missing values. The specific features that were missing values were Sulfate, Ph and Trihalomethanes. It was essential that we dealt with these in order to run different machine learning models on them. To remedy this, we attempted 4 different ways of cleaning our data. These were the deletion of missing values, Min-Max normalization, Simple Imputer, and Classified imputation.

The first method that we used was the deletion of missing values, where if a row in the dataset was missing a value, we would remove it from the training set This was the easiest way of going about cleaning the data. Unfortunately it wasn't effective in providing results, as it removed a lot of the training data that we needed for the model. We did this in pandas by using the df.dropna[] function with Sulfates, Ph and Trihalomethanes.

We also tried simple imputation. Simple imputation essentially replaces the missing values with whatever values desired. The 3 main methods of simple imputation that we used were simple mean imputation, simple median imputation, and simple mode imputation. As stated in the name, simple mean imputation replaces the missing values of the feature with the mean value of that specific feature. Same goes with median and mode except it would make the missing value median and mode respectively. While this did allow us to run the models, and made the accuracy of these models higher, we believed that we could do better. We ran simple imputation using scikitlearns simple imputer package.

```
>>> from sklearn.impute import SimpleImputer
```
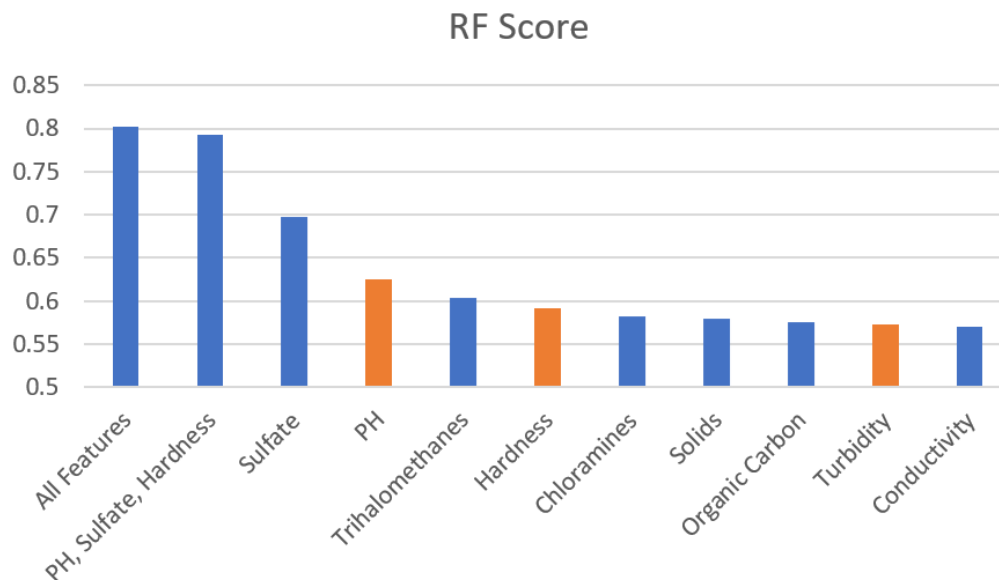
The final and most effective method that we did was classified imputation. It is similar to simple imputation as it also replaces the missing values with the mean value of the feature; however, it does that based on the class the datapoint is assigned to. Potability is a binary value, with 0 meaning that the water is potable and 1 meaning the water isn't potable. What this means is that the values will be imputed based on the classand would replace the empty data based on the 1 or 0 for the potability. This way it is a lot more unified in the way that the values are replaced which allowed us to get the highest accuracy for our model this way. The code that we used can be seen below.

```
#Nan to Mean
df['ph']=df['ph'].fillna(df.groupby(['Potability'])['ph'].transform('mean'))
df['Sulfate']=df['Sulfate'].fillna(df.groupby(['Potability'])['Sulfate'].transform('mean'))
df['Trihalomethanes']=df['Trihalomethanes'].fillna(df.groupby(['Potability'])['Trihalomethanes'].transform('mean'))
```

Once our data wasn't missing values, the last cleansing procedure we applied to our data was min-max normalization. This is a technique in which the lowest value in a certain feature becomes a 0 and the highest value becomes a 1 and all the other values become scaled relatively. While at first this seemed like a good idea as the data would get normalized, it became challenging as a lot of the missing values were assigned 0, which inhibited and skewed the models quite a bit.
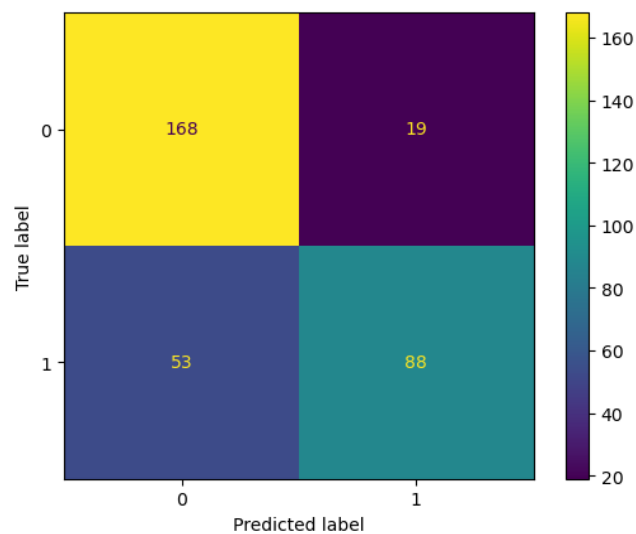
## Analysis:

We trained several different classification models on our data test, including K-Nearest Neighbors, Random Forest, Logistic Regression, Gradient Descent, Support Vector Machine, and a Neural Network. Factors such as the method of data cleaning and imputation significantly impacted the accuracy of the models. We used 90% of the dataset for training and 10% for testing, as this struck a balance between giving the models enough data to make accurate predictions while preventing overfitting of the data. The high accuracy we achieved was 80.2% with the Random Forest Model that used all the features as predictors. The three features that best predicted potability were pH, Sulphates, and Hardness. Training a Random Forest Model using these three features as predictors resulted in an accuracy of 79.3%. Below is the accuracy of the Random Forest Model on the different individual and combinations of features in the dataset.

### RF Score



**Confusion Matrix:**

Overall, our confusion matrix displayed the successes of our model. Obviously the ideal confusion matrix correctly guesses that water is drinkable when it actually is and vice versa, but our confusion matrix showed that our model was fairly accurate. For water potability, the most important thing we kept in mind was minimizing the false positives, as water that isn't drinkable when ingested is critically harmful to the human body. Our confusion matrix only had approximately 5 percent of the values to be false positive, which is a good statistic. Although minimizing false negatives is important, not drinking drinkable water is a good example of being safe rather than sorry, so we weren't as concerned with this. The model with the least percentage of false positives was Random Forest, and this was the main reason why we stuck with that model. One thing to keep in mind however, is the fact that there were more data points of undrinkable water vs drinkable which would skew the model towards making more negative predictions.



## Results of Conjectures:

For each of the features in the conjectures, we decided to run the Random Forest Model as it was the most accurate. We then compared them to each other to determine if our conjectures were correct or not. Below the summary of each conjecture a bar plot can be seen with all the accuracies of the model for the specific conjectures.
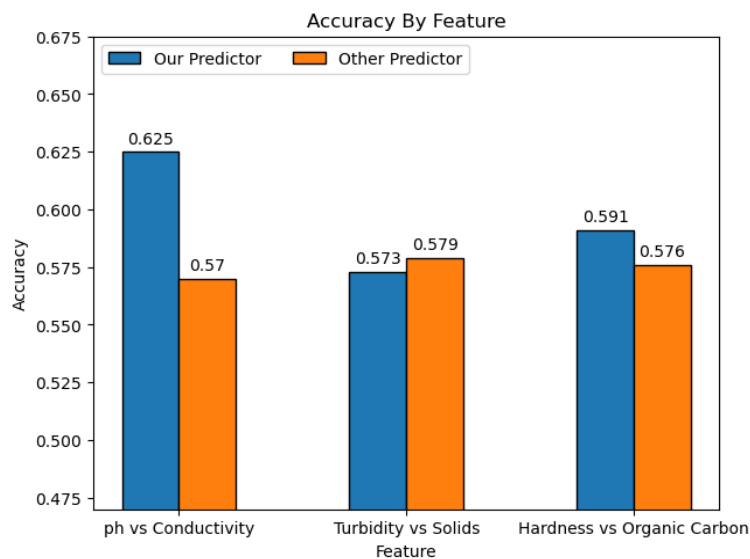
**Conjecture 1**

Our first conjecture was that **pH is a better predictor of potability than conductivity.** This conjecture was proven true due to our analysis, and was the only conjecture out of the 3 with a significant enough margin in between the accuracies. The accuracy of the pH model was 62.5 percent, when compared to the conductivity predictor at 57 percent. This conjecture was therefore proven correct

## Conjecture 2

Our second conjecture was that **Turbidity is a better predictor of potability than solids.** Although solids were a better predictor of potability than turbidity, the difference in accuracy was insignificant. Individually, both features were not good predictors of potability as they only correctly predicted potability around 60% of the time. We originally hypothesized that turbidity would be a better feature to base potability on, but solids are a marginally better predictor of drinkability. This conjecture was proven marginally wrong through our model.

## Conjecture 3

Our third conjecture was that **Hardness is a better predictor of potability than Organic carbon.** This conjecture was proven true when we analyzed the accuracy by feature for the dataset, as Hardness was given a higher accuracy rating than Organic Carbon. However, since both of these values were close to 50% accuracy, they were both ignored as single-column predictors because they were not accurate enough. We originally predicted that Hardness is a better predictor of potability when compared with It can be considered that this conjecture was marginally true.



**Challenges:**

We ran into a lot of challenges when working with this dataset. One of the main issues that we ran into was missing values. The columns with the missing values were Ph, Sulfate, and Trihalomethanes. With these missing values, it's impossible to run machine learning algorithms on the data to predict the potability.

While these allowed us to run the machine learning algorithms that we desired on the dataset, the accuracies of these models weren't ideal as they were around 50 percent. We soon realized that the closer to 50 percent that they got, the models were doing nothing better than guessing because our dataset condition is binary with 0 meaning the water is not potable and 1 means that the water is potable.

To combat this, we decided to do a more complex version of the imputation that we did earlier. This type of imputation is called classified imputation. Classified Imputation input values that provided the water was potable separate from those values that stated that the water wasn't potable. The values that provided 1 as the result would have a different imputation than the values that provided 0. This would fill in the empty values more purposefully and provided us with a Random Forest model closer to 80 percent accuracy.
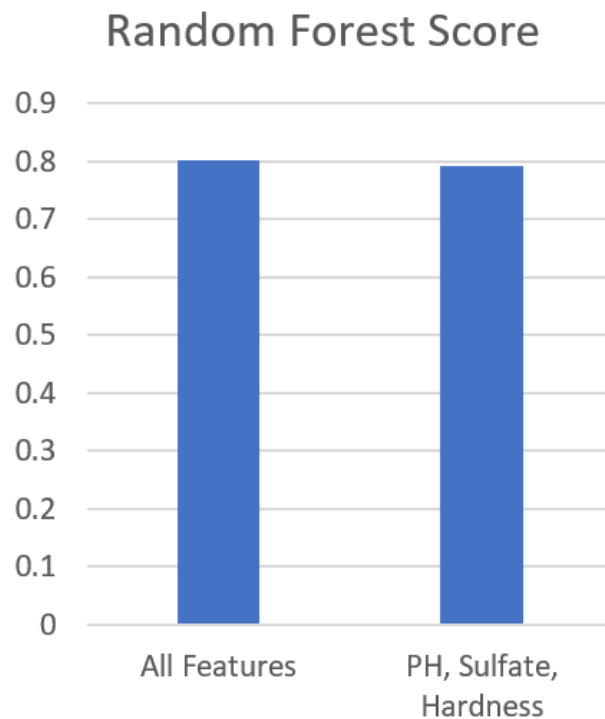
**Conclusion:**

To summarize, our main motivation behind this was to see what features people would have to look into when determining the drinkability of water. We looked at this from a humanitarian perspective, as drinkable water is an essential part of survival.

We cleaned the data, and ran multiple machine learning models on the dataset to determine the feature(s) with the biggest impact of potability. To ensure that our findings were correct, we made a confusion matrix to determine how correctly the model is outputting potability.

After thorough investigation with numerous models with all the different features, we have come to the recommendation to look out for specifically pH, Sulfate and Hardness as the biggest contributors to if the water is potable or not. In our bar chart below, when comparing the Random Forest Model accuracy with potability, it is nearly identical. The accuracy for all of the features is 80%, when compared to the 79.3% that is provided with the 3 features listed above.

From a business standpoint with minimal resources, our recommendation for any humanitarian effort to improve potability would be testing of Hardness, pH, and Sulfate in the water. Therefore they could use the resources gained from making less measurements to collect more samples for those three features.

## Random Forest Score



**Member Contributions:**

Abhi Yammanuru - Organized meeting times, performed simple imputation on the dataset, worked on slides, worked on introduction, conjectures, data cleaning methods and challenges sections of the essay

Michael Herchenroder- Assembled and organized the slides of the presentation, and helped figure out methods of analysis and data cleaning.

Ryan Offstein- Wrote code for imputing the data, trained several of the machine learning models, specifically the support vector machine model, created the confusion matrix, worked on several of the slides.

William Abraham- Wrote speaker notes for slides, added graphs, cleaned up slide visibility, thought of different ideas for report structure, wrote, connected and compiled code in jupyter notebook

Ajay Bhagavatula- Wrote code to clean and impute data, trained several machine learning models, created visualizations, worked on the analysis section of the essay, worked on the presentation