

# Project Draft

## Introduction

To safeguard agricultural productivity and maintain soil health in Telangana, our research team is conducting an analysis of irrigation water quality across various mandals within the state. Leveraging data obtained from the “Telangana Open Data Portal”, we aim to identify potential correlations between water quality parameters and their implications for irrigation practices. This initiative seeks to provide actionable insights for farmers, agricultural scientists, and policymakers to optimize water resource management strategies. From previous pre-processing of the data, we found that the following features might have a bigger impact on our models.

Data Overview:

- **Mandal:** Administrative division in Telangana where the water sample was collected, serving as the primary geographical unit for our analysis.
- **Village:** Specific villages within each mandal from which samples were collected, offering granular data points for our study.
- **RSC** (Residual Sodium Carbonate): Indicates the excess of carbonate and bicarbonate ions over calcium and magnesium ions in water, crucial for evaluating the risk of soil sodicity.
- **SAR** (Sodium Adsorption Ratio): Measures the sodium hazard to crops and soil, a key factor in assessing water’s suitability for irrigation.
- **Na** (Sodium content): The concentration of sodium in the water, which influences both SAR and the overall salinity hazard.
- **E.C** (Electrical Conductivity): Reflects the water’s salinity level, directly impacting plant water uptake and soil structure.
- **TDS** (Total Dissolved Solids): The overall concentration of dissolved substances in water, indicating potential salinity issues.
- **HCO<sub>3</sub>** (Bicarbonate level): Concentration of bicarbonate in water, affecting soil pH and the precipitation of calcium and magnesium.

- **pH:** The acidity or alkalinity of the water, with significant implications for nutrient availability and microbial activity in the soil.
- **Classification.1:** Safety classification of the water (Safe P.S., Marginal MR, Unsafe U.S.). This variable is our response variable and has been changed to a binary response by changing MR to U.S.

The primary objective of this study is to establish a comprehensive understanding of irrigation water quality within Telangana's agricultural landscapes. By analyzing the relationships between the above water quality indicators and their collective impact on soil and plant health, we aim to categorize water sources into suitability classes for irrigation. This classification will help in formulating guidelines for water use in agriculture, thereby mitigating the risks associated with inappropriate water sources.

### Random Forest Model: (Brayan Gutierrez, Tie Wang, Ulises Ramirez)

To avoid issues such as overfitting, high variance, and poor generalization to unseen data that a single decision tree is prone to, we decided to use an ensemble method for decision trees. Although the bagging ensemble model avoids the issues that a single decision tree presents, a unique issue arises when using the bagging model. Variables in the bagging model tend to be highly correlated as the most important variables will always be chosen in each decision tree the model builds. To avoid the issues that a single decision tree presents, decorrelate the decision trees, and the fact that random forest models are highly powerful, we decided to use this specific ensemble method. However, with deciding to use the random forest model we are sacrificing the interpretability that was present in a single decision tree, computational resources, and we gain the requirement of hyperparameter tuning that is not a major issue in a single decision tree. Like the previous model, our response variable is the **Classification.1** variable with the rest of the variables being our predictors. Below is the formula for our random forest model:

$$\text{Classification.1} \sim \text{Mandal} + \text{Village} + \text{RSC} + \text{SAR} + \text{Na} + \text{E.C} + \text{TDS} + \text{HCO3} + \text{pH}$$

First, we wanted to see how the base random forest with no tuned hyperparameters performs on our data. The data was read, the NA observations were omitted, and the response variable was changed to a binary response. Furthermore, since more than one category was present in the **Mandal** and **Village** variables and for this model to function properly, these variables were one-hot-encoded using the `dummyVars()` function from the `caret` library.

```
library(tree)
```

Warning: package 'tree' was built under R version 4.3.3

```
library(randomForest)
```

Warning: package 'randomForest' was built under R version 4.3.3

randomForest 4.7-1.1

Type rfNews() to see new features/changes/bug fixes.

```
library(caret)
```

Loading required package: ggplot2

Attaching package: 'ggplot2'

The following object is masked from 'package:randomForest':

margin

Loading required package: lattice

```
ground_water_quality_2022_post <- read.csv("ground_water_quality_2022_post.csv", header =  
ground_water_quality_2022_post = na.omit(ground_water_quality_2022_post)  
  
# Omitting CO3 and Season since they're the same  
ground_water_quality_2022_post$CO3 = NULL  
ground_water_quality_2022_post$season = NULL  
  
# Selecting the most useful features  
gw_df = ground_water_quality_2022_post[,c(23, 21, 16, 9, 10, 11, 8, 3, 4)]  
  
# One-hot-encoding Mandal and Village  
dummy = dummyVars(" ~ .", data = gw_df)  
dum_df = data.frame(predict(dummy, newdata = gw_df))  
  
# Changing Marginal Safe (MR) to Unsafe (U.S.)
```

```

Classification.1 = ground_water_quality_2022_post[,24]
final_gw_df = data.frame(dum_df, Classification.1)
final_gw_df$Classification.1 = gsub("MR", "U.S.", final_gw_df$Classification.1)
final_gw_df$Classification.1 = as.factor(final_gw_df$Classification.1)

```

To obtain the testing error for a classification random forest model, we would need a confusion matrix. Since we are testing this model 10 times with different 80% training/20% testing sets, we decided to store all of the confusion matrices for these different splits in a list called `conf_mat`. Furthermore, each model's out-of-bag score was calculated and stored in the OOB list that is called before the loop. The model was trained on the 10 different training sets and then tested on 10 different testing sets in a for loop as the `cv.glm()` function would only work on a regression random forest model. The seed used in our model was set to 1. Additionally, we plotted the variable importance for each training/testing split. The code for this can be seen below:

```

conf_mat = list()
OOB = rep(0,10)

set.seed(1)

for (i in 1:10) {

  # Splitting data
  sample = sample(1:nrow(final_gw_df), 0.8 * nrow(final_gw_df))
  training = final_gw_df[sample,]
  testing = final_gw_df[-sample,]

  # Random Forest Model
  rf_gw = randomForest(Classification.1 ~., data = training, proximity = T, importance = T)

  # OOB of Models
  conf = rf_gw$confusion[, -ncol(rf_gw$confusion)]
  OOB[i] = 1 - (sum(diag(conf)) / sum(conf))

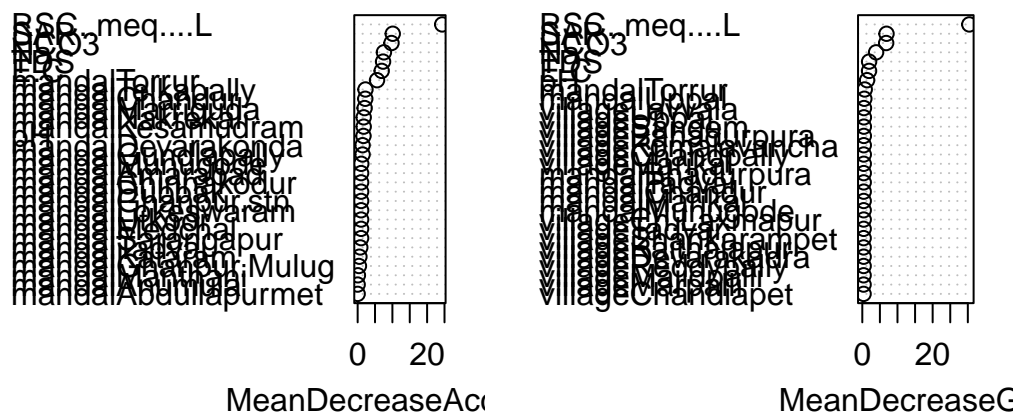
  # Importance Plots
  importance(rf_gw)
  varImpPlot(rf_gw, main = paste('Ground Water Random Forest Split', i, sep = ' '))

  test.pred = predict(rf_gw, newdata = testing)

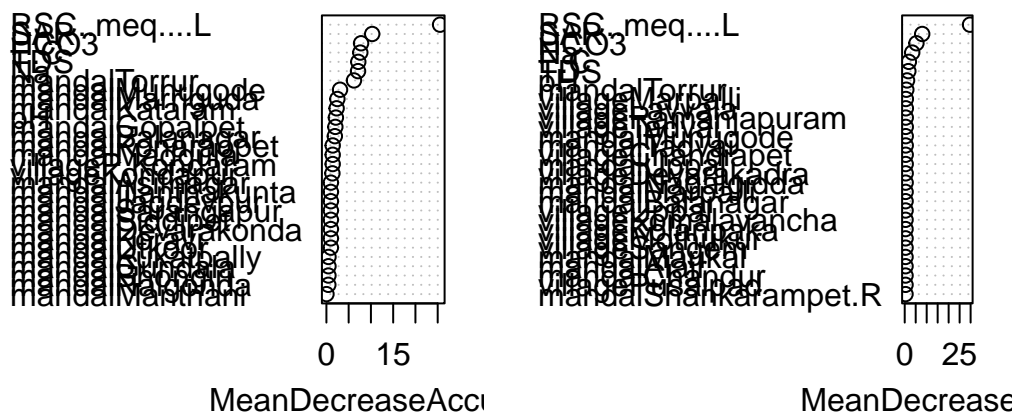
  # Confusion Matrix

```

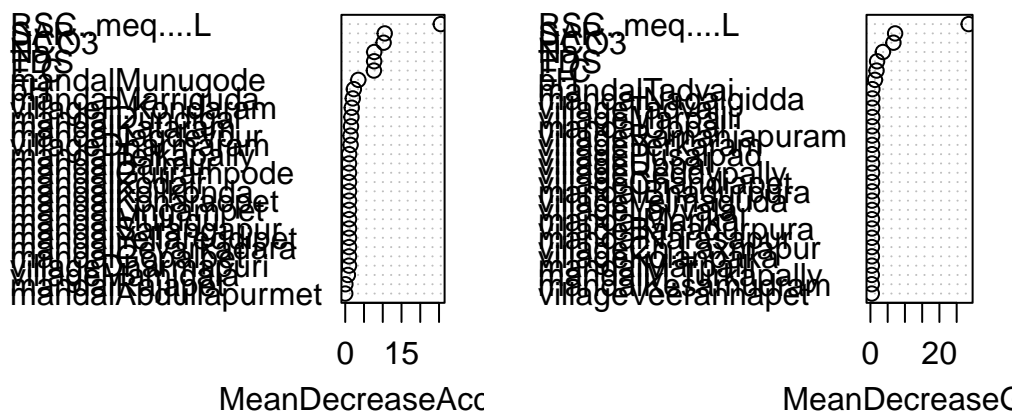
## Ground Water Random Forest Split 1



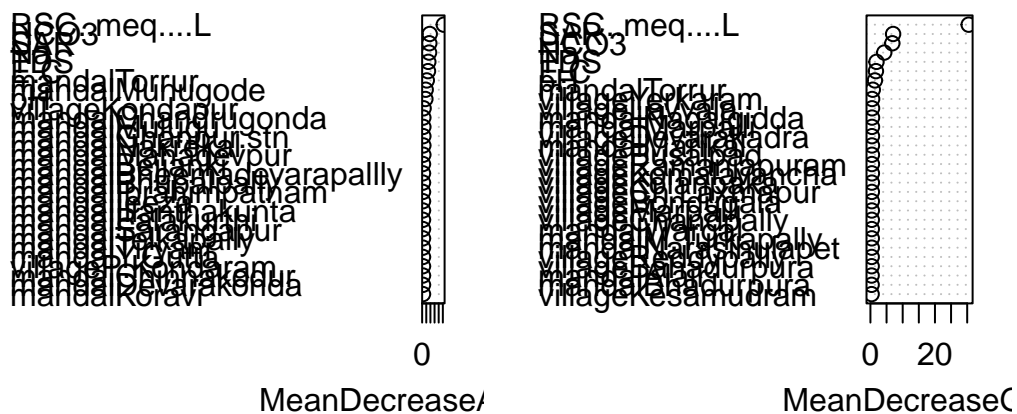
## Ground Water Random Forest Split 2



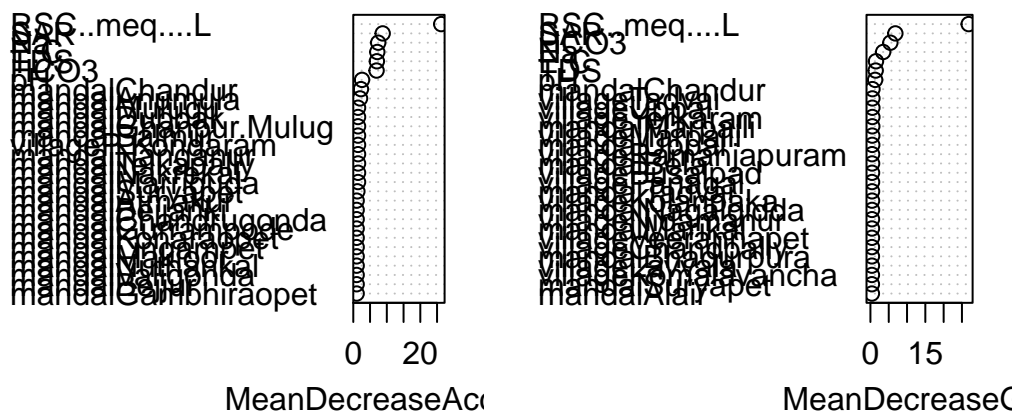
## Ground Water Random Forest Split 3



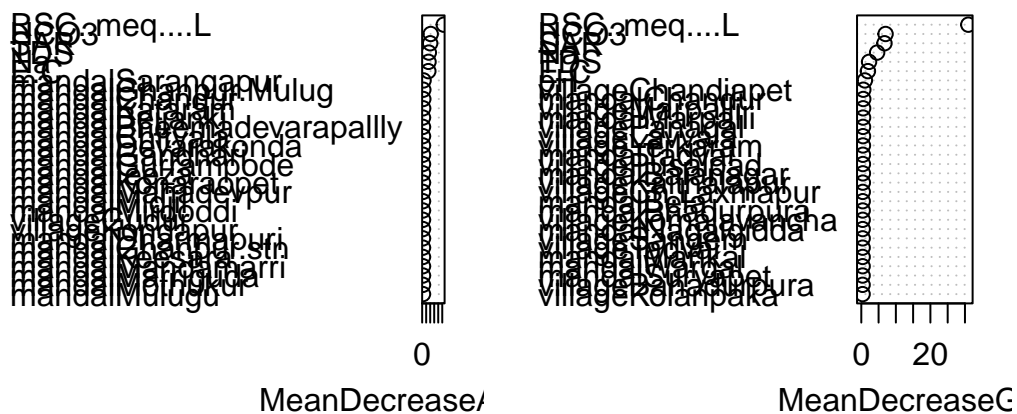
## Ground Water Random Forest Split 4



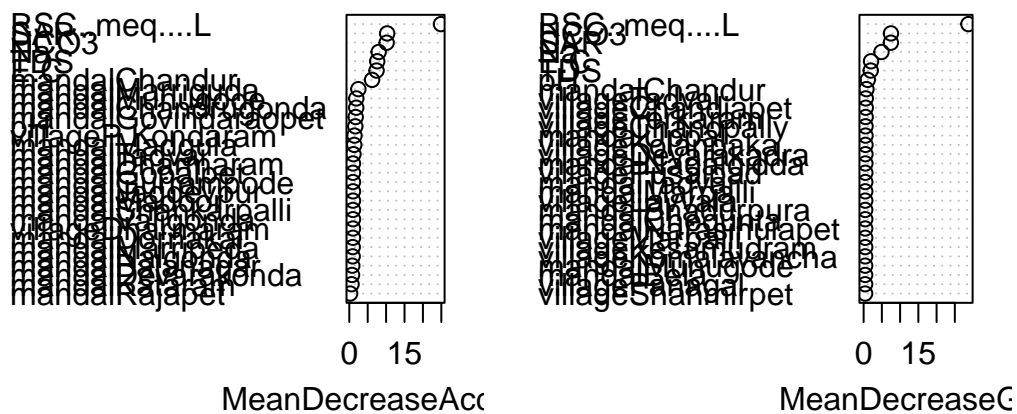
## Ground Water Random Forest Split 5



## Ground Water Random Forest Split 6

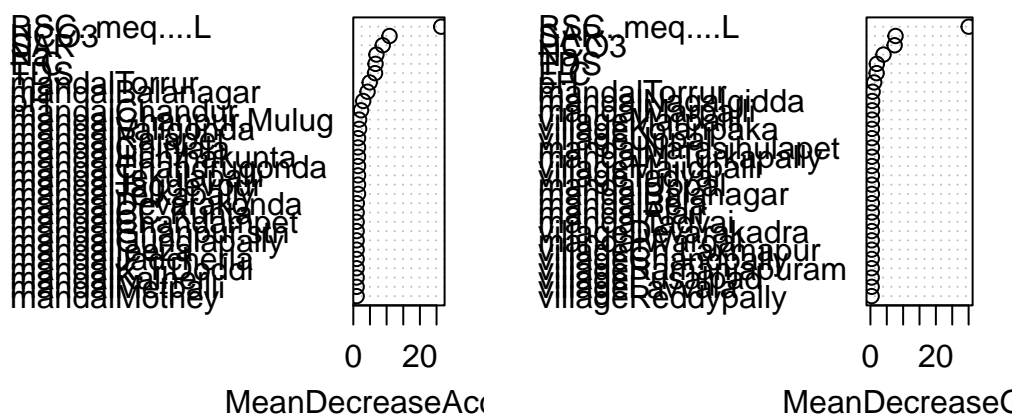


## Ground Water Random Forest Split 7

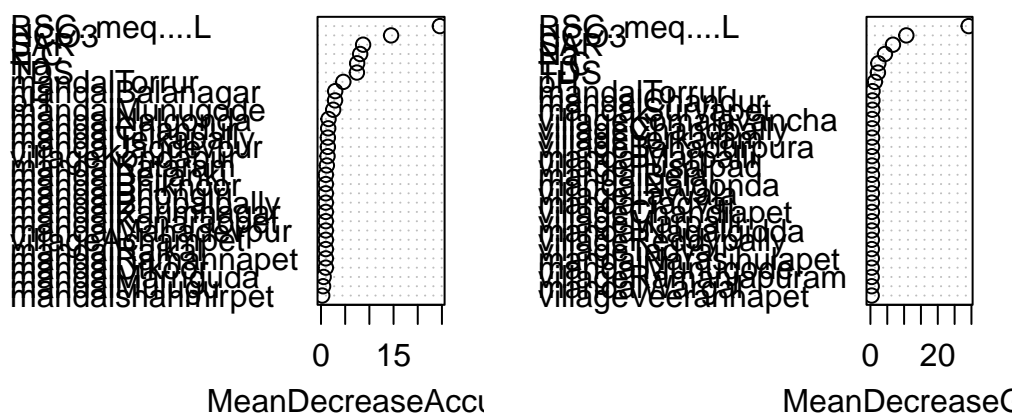




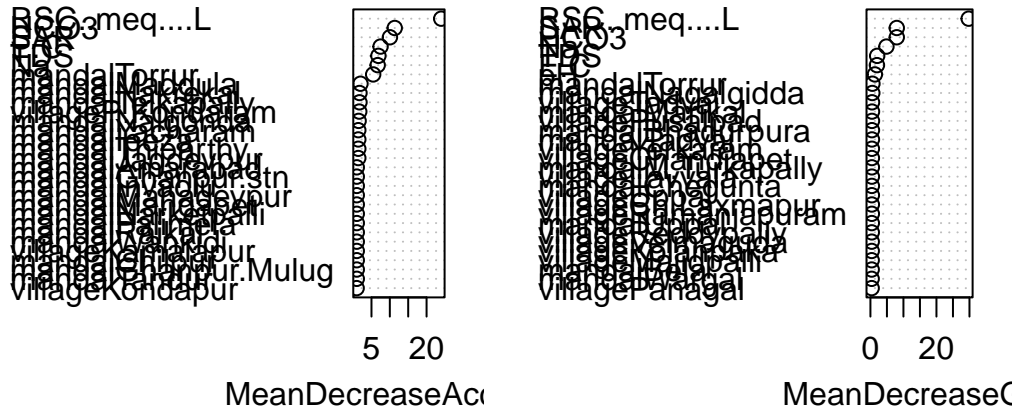
## Ground Water Random Forest Split 8



## Ground Water Random Forest Split 9



## Ground Water Random Forest Split 10



We can see that the 3 most important variables in each split were RSC, SAR, and HCO3, with SAR and HCO3 competing for second most important across the different splits. With the confusion matrices for each training/testing split now saved, we used these matrices to now calculate the testing error of each split.

```
# Test Error Calculations
test_errors = rep(0, 10)

for (i in 1:10) {
  TP = conf_mat[[i]][2, 2] # True Positives
  FP = conf_mat[[i]][1, 2] # False Positives
  FN = conf_mat[[i]][2, 1] # False Negatives
  TN = conf_mat[[i]][1, 1] # True Negatives

  # Calculating test error rate
  test_errors[i] <- (FP + FN) / sum(conf_mat[[i]])
}

scores = list(test_error = test_errors, OOB_score = OOB)

scores
```

```
$test_error
[1] 0.024390244 0.006097561 0.006097561 0.006097561 0.006097561 0.000000000
[7] 0.012195122 0.000000000 0.012195122 0.012195122
```

```
$OOB_score
[1] 0.015313936 0.006125574 0.015313936 0.007656968 0.009188361 0.007656968
[7] 0.015313936 0.006125574 0.016845329 0.021439510
```

```
(avg_test_error = sum(test_errors) / 10)
```

```
[1] 0.008536585
```

We found that the base random forest model with 10 different training/testing splits had a test error between 0%-2.4%. With an average testing error of 0.8%. Comparing the testing errors and the OOB scores, we see that they are not far off from each other, meaning that the base random forest model performs well with unseen data and data that it did not see during training. This suggests that the base model learned the underlying patterns in the data well. This result is very promising, however, (*write something about less U.S. observations*). Next, we then moved on to tuning the hyperparameters of the random forest model on our data...