

STATS 101B Final Project

Qianli Wu
Sebastian Rivera-Chepetla
Ovie Soman

1. Introduction

The Random Forest (RF) algorithm is a classification algorithm consisting of decision trees. It begins with one question and two possible answers to that question. Depending on which answer choice is picked, it either asks another question or gives an output or “decision”. The second question may be followed by a third, a fourth, etc.

The decision tree algorithm groups data based on whether the criterion at each node is true or false. The maximum number of nodes (*maxnodes*), the count of observations in terminal nodes (*nodesize*), and the count of predictors (*mtry*) can be tuned. While decision trees are helpful, they are prone to bias and overfitting. To avoid this, Random Forest is used. It is composed of several decision trees, made from several sampling sets taken from the data.

The parameter *ntree* shows the number of trees used, and *replace* signifies whether sampling was done with replacement. *classwt* shows what the weights of the two groups the data is classified into are, and *cutoff* signifies the threshold for data classification. The weight is normally specified for the class denoted by “1” (reference class) - the weight for the “0” class can be obtained by subtracting the reference class weight from 1. *cutoff* signifies the threshold for data classification into the reference class. If an observation does not pass the threshold, it is relegated to the “0” class. RF tells us how likely an observation is to be in the reference class, rather than an absolute class.

The seven tuning parameters mentioned above (*ntree*, *replace*, *mtry*, *nodesize*, *maxnodes*, *classwt*, and *cutoff*) are investigated in this assignment. RF is evaluated by its classification accuracy (the ratio between correctly classified data values and the total data values in a new dataset). Cross-Validation is used, since producing a new dataset might not be feasible, to divide a dataset into two. One part is used to construct the RF and the other to test it. This is repeated many times with different data divisions, and the average accuracy from all these repetitions is reported. Our goal is to select the tuning parameters that most affect CV accuracy. The RF algorithm is built using a dataset containing information that helps predict the occurrence of cardiovascular disease. The resulting dataset with tuning variables and CV accuracy is analyzed.

2. Methodology

a) Experimental Design

Each parameter of concern has 2 levels. Therefore, a complete factorial design required $2^7 = 128$ runs. 7 degrees of freedom are reserved for main effects, 21 for 2-factor interactions,

and 99 for 3 and higher-order interactions. However, due to budget constraints, we could only run up to 35 tests. Taking into consideration the effect sparsity principle, projection property, and sequential experimentation, we end up deciding to go with a $2^{(7-2)}$ design. We used the 'FrF2' function from the 'FrF2' library to recommend a design that would have the largest resolution and smallest aberration possible, as shown in Table 1. Also, we constructed another experimental design using the D-optimal design approach, as shown in Table 2.

##	ntree	mtry	replace	nodesize	classwt	cutoff	maxnodes
## 1	100	2	0	1	0.5	0.2	1000
## 2	1000	2	0	1	0.5	0.8	10
## 3	100	6	0	1	0.5	0.8	10
## 4	1000	6	0	1	0.5	0.2	1000
## 5	100	2	1	1	0.5	0.8	1000
## 6	1000	2	1	1	0.5	0.2	10
## 7	100	6	1	1	0.5	0.2	10
## 8	1000	6	1	1	0.5	0.8	1000
## 9	100	2	0	11	0.5	0.2	10
## 10	1000	2	0	11	0.5	0.8	1000
## 11	100	6	0	11	0.5	0.8	1000
## 12	1000	6	0	11	0.5	0.2	10
## 13	100	2	1	11	0.5	0.8	10
## 14	1000	2	1	11	0.5	0.2	1000
## 15	100	6	1	11	0.5	0.2	1000
## 16	1000	6	1	11	0.5	0.8	10
## 17	100	2	0	1	0.9	0.2	10
## 18	1000	2	0	1	0.9	0.8	1000
## 19	100	6	0	1	0.9	0.8	1000
## 20	1000	6	0	1	0.9	0.2	10
## 21	100	2	1	1	0.9	0.8	10
## 22	1000	2	1	1	0.9	0.2	1000
## 23	100	6	1	1	0.9	0.2	1000
## 24	1000	6	1	1	0.9	0.8	10
## 25	100	2	0	11	0.9	0.2	1000
## 26	1000	2	0	11	0.9	0.8	10
## 27	100	6	0	11	0.9	0.8	10
## 28	1000	6	0	11	0.9	0.2	1000
## 29	100	2	1	11	0.9	0.8	1000
## 30	1000	2	1	11	0.9	0.2	10
## 31	100	6	1	11	0.9	0.2	10
## 32	1000	6	1	11	0.9	0.8	1000

class=design, type= FrF2

Table 1: Fractional Factorial Design

##	ntree	mtry	replace	nodesize	classwt	cutoff	maxnodes
## 1	100	2	0	1	0.5	0.2	10
## 4	1000	6	0	1	0.5	0.2	10
## 6	1000	2	1	1	0.5	0.2	10
## 10	1000	2	0	11	0.5	0.2	10
## 15	100	6	1	11	0.5	0.2	10
## 18	1000	2	0	1	0.9	0.2	10
## 27	100	6	0	11	0.9	0.2	10
## 29	100	2	1	11	0.9	0.2	10
## 34	1000	2	0	1	0.5	0.8	10
## 43	100	6	0	11	0.5	0.8	10
## 45	100	2	1	11	0.5	0.8	10
## 48	1000	6	1	11	0.5	0.8	10
## 49	100	2	0	1	0.9	0.8	10
## 55	100	6	1	1	0.9	0.8	10
## 60	1000	6	0	11	0.9	0.8	10
## 62	1000	2	1	11	0.9	0.8	10
## 66	1000	2	0	1	0.5	0.2	1000
## 71	100	6	1	1	0.5	0.2	1000
## 73	100	2	0	11	0.5	0.2	1000
## 76	1000	6	0	11	0.5	0.2	1000
## 78	1000	2	1	11	0.5	0.2	1000
## 83	100	6	0	1	0.9	0.2	1000
## 85	100	2	1	1	0.9	0.2	1000
## 88	1000	6	1	1	0.9	0.2	1000
## 90	1000	2	0	11	0.9	0.2	1000
## 97	100	2	0	1	0.5	0.8	1000
## 104	1000	6	1	1	0.5	0.8	1000
## 106	1000	2	0	11	0.5	0.8	1000
## 116	1000	6	0	1	0.9	0.8	1000
## 118	1000	2	1	1	0.9	0.8	1000
## 121	100	2	0	11	0.9	0.8	1000
## 127	100	6	1	11	0.9	0.8	1000

Table 2: D-optimal Design

For the Fractional Factorial Design, as shown in Figure 1, there was no aliasing between the main effects. As such, the performance of this design would be acceptable when analyzing the main effects. Therefore, we concluded that the fractional factorial design was better for analyzing main effects and intercept only. However, this was not the case for the aliasing interactions ("replace:mtry", "mtry:classwt", "replace:classwt", "ntree:mtry", "ntree:classwt", "ntree:replace").

Considering a D-optimal Design, as shown in Figure 2, the correlation was weak between the main effects and interactions. Aliasing between main effects, main effects, and specific interactions, and between specific interactions was low. Also, as shown in Figure 4, all VIF values are close to 1 and below 5, which indicates that there was no serious multicollinearity or aliasing issue for the D-optimal design.

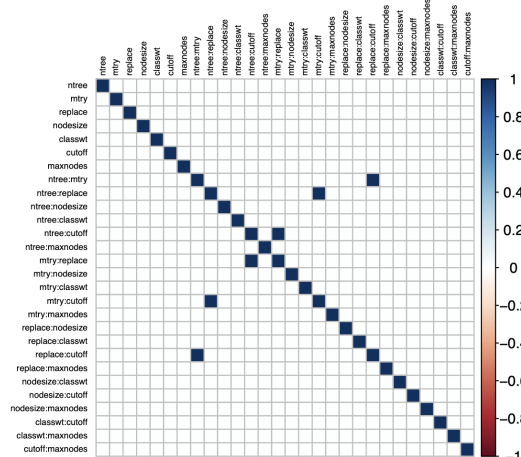


Figure1: Color Map For Factorial Design

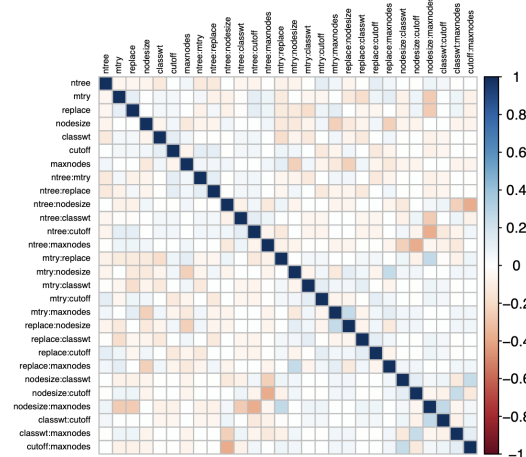


Figure2: Color Map for D-optimal Design

Therefore, the experimental design that we would recommend would be the D-optimal design since there is no serious aliasing issue or a serious violation of multicollinearity for main effects and two-factor interactions. D-optimal designs are desirable since they minimize the variance-covariance matrix of the coefficient's estimates, which facilitates developing precise inferences. Also, In practical terms, since both designs require 32 runs and there's no huge difference between the proportion of two levels, the time it takes to run the experiments are similar for both designs.

VIF Values

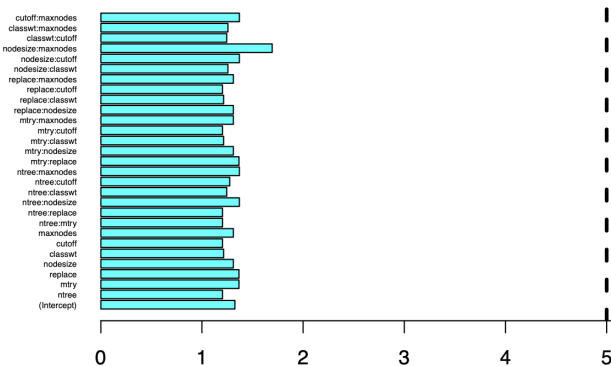


Figure3: VIF of D-optimal Design

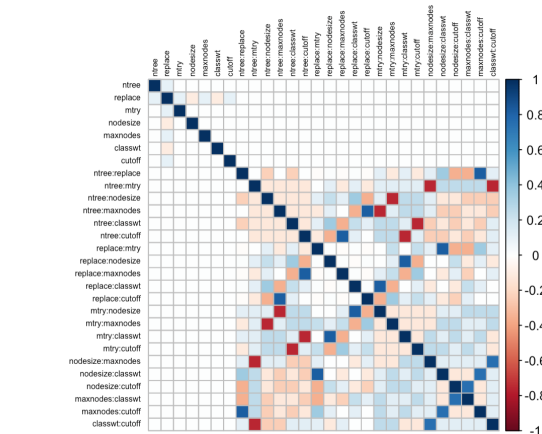


Figure 4: Color Map for Provided Design

For the 3-level design provided on Question 4, as shown in Figure 4, the provided design has lighter cells for aliasing among main effects with other main effects and there's no alias between main effects and two-parameter interactions. However, since the cells at bottom right corner are darker, the degree of aliasing is generally greater than these for our D-optimal design for two-parameter interactions. It is reasonable since the proposed design only has 22 runs, which is less than our D-optimal design with 32 runs.

Also, the proposed 3-level design allows a quadratic model. However, our 2-level D-optimal design doesn't support that. But our D-optimal design has better performance on estimating two-parameter interactions due to less degree of aliasing.

b) Data Analysis

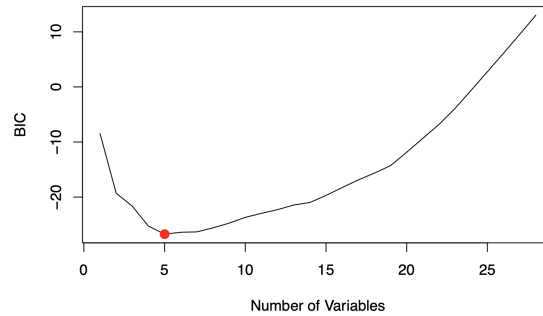
The data was collected using the D-optimal design and subsequently analyzed. The run order of the experiments was randomized so as to maintain independence. The RF algorithm was run on the data collected using the D-optimal design. The tuning parameters in the dataset were converted to factors from numerical variables, and their levels were re-encoded to “-1” and “1”, with the former representing the level with the lower value. A linear model was run on the data using all the tuning parameters as independent variables and the linear model summary was analyzed. There were no significant factors. The difference between the Multiple and Adjusted R^2 values was substantial (about 0.60), which indicated that the high Multiple R^2 value could be due to the inclusion of several independent factors. We would thus need a reduced model.

The Bayes Information Criterion (BIC) was run on the model to trim out the factors that did not account for substantial variability. The formula for BIC is given below:

$$BIC = n \ln\left(\frac{RSS}{n}\right) + \ln(n)k$$

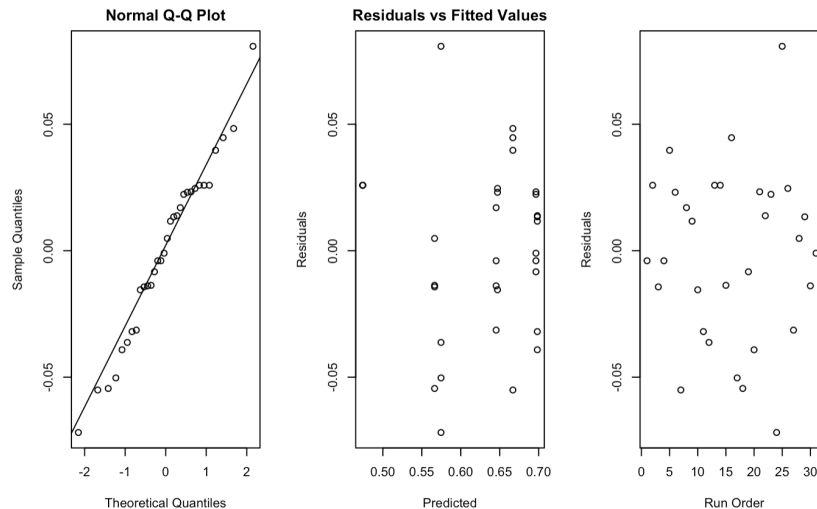
where n is the number of observations, k is the number of predictors, and the RSS is the Residual Sum of Squares for the linear regression. The model with the lowest BIC value (or the least complex model) was picked, and a reduced linear regression was run using the independent variables from this model with the lowest BIC value.

Figure 5: BIC Scores for Sub-Models



As shown in Figure 5, we selected and constructed the best sub-model with lowest BIC value. For the new model, the difference between the Multiple and Adjusted R^2 values was reduced to about 0.05, which was substantially lower. The significant factors were *classwt*, *cutoff*, *maxnodes* and the interaction between *classwt* and *maxnodes* and the interaction between *classwt* and *cutoff*. Next, the residuals were analyzed.

Figure 6: Analysis of Residuals



When analyzing the residual plots in Figure 6, we see that the points on the normal Q-Q Plot follow the qqline fairly well, except at the extremities. There is no evidence of violations of normality. Residuals in the residuals vs. fitted plot distributed evenly across each fitted value. And the distribution of the residuals vs run order plots are fairly homoskedastic. Thus, we say the model adheres to normality assumptions.

Since we have already run 32 tests, we are left with 3 runs for confirmation experiments. Using the final equation shown below, we found the settings that maximize porosity percentage by using the 'optim' function.

Table 3: Tuning Parameters

	Parameter	Low Level	High Level
A	ntree	100	1000
B	replace	0	1
C	mtry	2	6
D	nodesize	1	11
E	maxnodes	10	1000
F	classwt	0.5	0.9
G	cutoff	0.2	0.8

$$\hat{y} = 0.622765 - 0.048193E + 0.028372F + 0.015964G + 0.023461EF + 0.037590EG$$

The optim function suggests that the optimal setting for the model is to have 'ntree', 'replace', 'mtry', 'nodesize', 'maxnodes' and 'classwt' set at the low level and 'cutoff' set at its high level.

```
ntree mtry replace nodesize classwt cutoff maxnodes
-1    -1    -1      -1      -1      1      -1
```

By running the cv.rf() function 3 times with the optimal settings we found earlier, we find that the results of confirmation experiments are within our prediction interval (0.608844, 0.7895945). Furthermore, the mean and the fitted value are close to each other of the results. That is, 0.6992193, and 0.6970846 respectively.

3. Conclusion

Overall, our model was satisfactory. Our success came from the careful methodology that we followed. To start, by using tools such as color maps, we carefully evaluated the strengths of our choice of design, whether that be fractional factorial or D-optimal. For instance, we give up fractional factorial design because it suffers from multicollinearity of 3 pairs of the aliased interactions.

We selected the D-optimal Design since it minimizes multicollinearity and we found no serious aliasing issue for main effects and two-factor interactions.

We are confident our model works well since the residual plots suggest that all assumptions of our model are satisfied, which validate our model. Moreover, in the confirmation experiments, the mean response collected is close to the fitted value of prediction by our model, which proves that our model prediction is accurate. Additionally, we used Bayes Information Criterion to help select a reduced model to avoid overfitting. This procedure allowed us to proceed with a more reliable model.

We recommend three-level experiments with more runs for future experimentation. It is because three-level experiments allow for quadratic models. Thus, there are several advantages that promote the use of three-level experiments.

4. Appendix



Final Project

Ovie Soman, Qianli Wu, Sebastian Rivera-Chepetla

05 June, 2022

Contents

Question 1. Propose a fractional factorial design for the problem. In addition, propose an experimental design constructed using the optimal design approach.	3
The goal of this project is to identify the tuning parameters of random forest that affect the cross-validation accuracy.	3
Fractional Factorial Design	3
An Experimental Design Constructed using the Optimal Design Approach.	6

Question 2. Compare the optimal design with the fractional factorial design in practical and statistical terms. For instance, what is the performance of the designs for studying the main effects of the tuning parameters only? Can they estimate all two-parameter interactions? Why or why not? How do they compare in terms of multicollinearity?	7
2.1: what is the performance of the designs for studying the main effects of the tuning parameters only?	7
2.2 Can they estimate all two-parameter interactions?	9
2.1 Color Map for D-optimal Design	10
3. Compare	12
Question 3. Recommend one experimental design between the two options in Question 1. Motivate your decision.	13
Question 4. Using a commercial software, the TAs and I came up with the experimental design shown in Table 2. How does your recommended design in the previous question compare with this one?	14
4.1 Color Map for Table 2 Design	15
Question 5. Collect data using your recommended design in Question 3.	16
Question 6. Conduct a detailed data analysis. What are the influential tuning parameters? What is the final model that links the tuning parameters to the cross-validation accuracy? Does the final model provide a good fit to the data?	17
A model with interaction terms with optimal design	18
Bayes Information Criterion (BIC)	18
Analysis of Residuals	19
Running Confirmation Tests	19

Question 1. Propose a fractional factorial design for the problem. In addition, propose an experimental design constructed using the optimal design approach.

The goal of this project is to identify the tuning parameters of random forest that affect the cross-validation accuracy.

- We have seven tuning parameters:

	Parameter	Low Level	High Level
A	ntree	100	1000
B	replace	0	1
C	mtry	2	6
D	nodesize	1	11
E	maxnodes	10	1000
F	classwt	0.5	0.9
G	cutoff	0.2	0.8

Fractional Factorial Design

Since there are 7 factors each with 2 levels, a complete factorial design 2^7 requires 128 runs. Only 7 degrees of freedom correspond to main effects, 21 correspond to two-factor interactions, and the rest of 99 are associated with three-factor and higher interactions.

Due to limitation of resources, we can only run at most 35 tests. We assume that certain high-order interactions are negligible.

Since

- The effect sparsity principle
- The projection property
- Sequential experimentation

We choose the One-Quarter Fraction of the 2^7 Design, i.e., 2^{7-2} Design.

Letting the ‘FrF2’ function recommend one design. Generally, the recommended designs are those in Table 8.14 in Chapter 8. The designs recommended by the function have the largest possible resolution. They also have the smallest aberration as defined in Section 8.4.1. Fractional factorial designs with minimum aberration are preferred because they minimize the aliasing among the factors’ effects.

```
factors <- list(ntree = c(100, 1000),
               mtry = c(2, 6),
               replace = c(0, 1),
               nodesize = c(1, 11),
               classwt = c(0.5, 0.9),
               cutoff = c(0.2, 0.8),
               maxnodes = c(10, 1000))
my.design <- FrF2(nruns = 32, nfactors = 7, randomize = F, factor.names = factors)
print(my.design)
```

```
design.info(my.design)$catlg.entry  
design.info(my.design)$aliased
```

An Experimental Design Constructed using the Optimal Design Approach.

Suggestion is to start with 32 runs

3 runs are saved for confirmation experiments or follow-up experiment.

```
# Create a full design for 7 factors each with 2 levels
full.design <- FrF2(nruns = 128, nfactors = 7, randomize = F, factor.names = factors)
# Consider a model with 35 runs
alternative.design <- optFederov(~.^2,
                                full.design, nTrials = 32, nRepeats = 1000)
print.data.frame(alternative.design$design)
```

Question 2. Compare the optimal design with the fractional factorial design in practical and statistical terms. For instance, what is the performance of the designs for studying the main effects of the tuning parameters only? Can they estimate all two-parameter interactions? Why or why not? How do they compare in terms of multicollinearity?

2.1: what is the performance of the designs for studying the main effects of the tuning parameters only?

```
# Visualize the aliasing in the design.
FD.alt <- (desnum(my.design)) # Extract the design.
# Create the model matrix including main effects and two-factor interactions.
FX.alt <- model.matrix(~., data.frame(FD.alt))

# Create color map on pairwise correlations.
contrast.vectors.correlations.alt <- cor(FX.alt)
corrplot(contrast.vectors.correlations.alt, type = "full", addgrid.col = "gray",
          tl.col = "black", tl.srt = 90, method = "color", tl.cex=0.5)
```

2.1.1 Color Map for Main Effects of Fractional Factorial Design Since all cells other than the diagonal are white, there's no alias among all main effects and intercept for fractional factorial design. Therefore, we have the performance of fractional factorial design is good for main effects and intercept.

```
# Visualize the aliasing in the design.
OD.alt <- alternative.design$design # Extract the design.

# Convert factors to -1 and 1
OD.alt <- sapply(OD.alt,function(x)(as.integer(x) -1.5)*2 )
OX.alt <- model.matrix(~., data.frame(OD.alt))

# Create color map on pairwise correlations.
contrast.vectors.correlations.alt <- cor(OX.alt)
corrplot(contrast.vectors.correlations.alt, type = "full", addgrid.col = "gray",
          tl.col = "black", tl.srt = 90, method = "color", tl.cex=0.5)
```

2.1.2 Color Map for Main Effects of D-optimal Design

2.1.3 VIF for Main Effects of D-Optimal Design

```
# We need to include the intercept when computing the VIF.
X.alt <- model.matrix(~., data.frame(OD.alt))
# Variance-covariance matrix of 46-run design. Assuming sigma^2 = 1
var.eff.one <- diag(solve(t(X.alt)%*%X.alt))
# Set the left margin of plot be 1
par(oma=c(0,1,0,0))
#create horizontal bar chart to display each VIF value
barplot(nrow(X.alt)*var.eff.one, main = "VIF Values", horiz = TRUE, col = "cyan", las=1, cex.names=0.4,
#add vertical line at 5
abline(v = 5, lwd = 3, lty = 2)
```

From the VIF plot above, we have there's no factor has VIF greater than 5. Thus, we don't have a severe multi-collinearity or aliasing issue.

2.2 Can they estimate all two-parameter interactions?

2.2.1 Color Map for Fractional Factorial Design

```
# Create the model matrix including main effects and two-factor interactions.
FX.alt <- model.matrix(~.^2-1, data.frame(FD.alt))

# Create color map on pairwise correlations.
contrast.vectors.correlations.alt <- cor(FX.alt)
corrplot(contrast.vectors.correlations.alt, type = "full", addgrid.col = "gray",
         tl.col = "black", tl.srt = 90, method = "color", tl.cex=0.5)
```

From the color map above, we get there's alias among some pairs of two-parameter interactions.

From the results of `design.info` in Question 1, we have $AB=CF$, $AC=BF$, and $AF=BC$, where $A=ntree$, $B=mtry$, $C=replace$, $D=nodesize$, $E=classwt$, $F=cutoff$, and $G=maxnodes$. The aliasing is also shown on the map above, as the corresponding dark cells for $mtry:classwt$ with $ntree:replace$, $replace:classwt$ with $ntree:mtry$, and $replace:mtry$ with $ntree:classwt$.

Therefore, we have the fractional factorial design's performance is best for all main effects and two-factor interactions except these aliasing factors (AB , CF , AC , BF , AF , and BC).

2.1 Color Map for D-optimal Design

```
OX.alt <- model.matrix(~.^2-1, data.frame(OD.alt))

# Create color map on pairwise correlations.
contrast.vectors.correlations.alt <- cor(OX.alt)
corrplot(contrast.vectors.correlations.alt, type = "full", addgrid.col = "gray",
          tl.col = "black", tl.srt = 90, method = "color", tl.cex=0.5)
```

From the D-optimal Design above, we have most cells with light colors, which means the strength of correlation between the main effects and two-parameter interactions are weak. Therefore, we have the degree of aliasing among the main effects, between the main effects and the two-parameter interactions, and among the two-parameter interactions exist but are fairly low.

2.2.3 VIF for D-Optimal Design

```
# We need to include the intercept when computing the VIF.
X.alt <- model.matrix(~.^2, data.frame(OD.alt))
# Variance-covariance matrix of 46-run design. Assuming sigma^2 = 1
var.eff.one <- diag(solve(t(X.alt)%*%X.alt))
# Set the left margin of plot be 1
par(oma=c(0,1,0,0))
#create horizontal bar chart to display each VIF value
barplot(nrow(X.alt)*var.eff.one, main = "VIF Values", horiz = TRUE, col = "cyan", las=1, cex.names=0.4,
#add vertical line at 5
abline(v = 5, lwd = 3, lty = 2)
```

From the VIF plot above, we have there's no factor has VIF greater than 5. Thus, we don't have a severe multi-collinearity or aliasing issue.

3. Compare

- We have the fractional factorial design only suffers from multi-collinearity problem of these three pairs of aliased interactions. It will be a good design if `mtry:classwt`, `ntree:replace`, `replace:classwt`, `ntree:mtry`, `replace:mtry`, and `ntree:classwt` are not potentially important effects.
- However, there's no serious multi-collinearity problem or aliasing issue for the D-optimal design. But the variances of coefficient estimates of both main effects and two-parameter interactions are inflated trivially due to collinearity.

```
# We compare the run size of each runs for both designs  
summary(as.data.frame(my.design))  
summary(alternative.design$design)
```

- On practical term, since both design require 32 runs and there's no huge difference between the proportion of two levels, we have the time it takes to run the experiments are similar for both designs.

Question 3. Recommend one experimental design between the two options in Question 1. Motivate your decision.

- We select the D-optimal Design since there's no series multi-collinearity problem or aliasing issue for all main effects and two-factor interactions.
- Also, D-optimal designs minimize the variance- covariance matrix of the coefficient's estimates and thus they are attractive to develop linear models that allow to make precise inferences. Our goal is to find the most important tuning parameters.

Question 4. Using a commercial software, the TAs and I came up with the experimental design shown in Table 2. How does your recommended design in the previous question compare with this one?

```
table2 <- read.csv("Data/Alternative_Experimental_Design.csv")  
head(table2)
```

4.1 Color Map for Table 2 Design

```
# Convert factors to -1 and 1
table2 <- sapply(table2,function(x)(as.integer(as.factor(x)) -2) )
table2[,3][table2[,3]==0] <- 1
X.alt <- model.matrix(~.^2-1, data.frame(table2[,,-1]))

# Create color map on pairwise correlations.
contrast.vectors.correlations.alt <- cor(X.alt)
corrplot(contrast.vectors.correlations.alt, type = "full", addgrid.col = "gray",
         tl.col = "black", tl.srt = 90, method = "color", tl.cex=0.5)
```

- The experimental design shown in Table2 is 7 factor with 3 levels, however, the design we proposed in Question 1 has only 2 levels for each factor.
- From the color map above, we have the design in Table 2 has lighter cells for aliasing among main effects with other main effects and there's no alias between main effects and two-parameter interactions. However, since the cells at bottom right corner are darker, we have the degree of aliasing is generally greater than these for our D-optimal design for two-parameter interactions. It is reasonable since proposed design only has 22 runs, which is less than our D-optimal design with 32 runs.
- Since the proposed design is a 3-level design, it allows quadratic model, which includes x^2 as predictors. However, our 2-level D-optimal design doesn't support that. But our design has better performance on estimating two-parameter interactions.

Question 5. Collect data using your recommended design in Question 3.

```
# import function cv.rf(design, y, X)
source("CrossValidation_RandomForest/CrossValidation_RF.R")
# import dataset Cardiovascular
load("Data/cardiovascular.Rdata")
```

- We need to randomized the run order of the experiments to ensure the independence of errors.

```
# convert the design from factor to numerical values
design.numerical <- as.data.frame(sapply(alternative.design$design, function(x) as.numeric(as.character(x))))
# randomize the run order
design.randomized <- design.numerical[sample(1:nrow(design.numerical)),]
head(design.randomized)
```

```
result <- cv.rf(design.randomized, y, X)
# save a copy of result thus no need to run cv.rf again
result.copy <- result
print(result)
```

Question 6. Conduct a detailed data analysis. What are the influential tuning parameters? What is the final model that links the tuning parameters to the cross-validation accuracy? Does the final model provide a good fit to the data?

```
# Converting numerical variables back to factors
str(result)
result[,c("ntree", "mtry", "replace",
          "nodesize", "classwt", "cutoff", "maxnodes")] <- lapply(result[,c("ntree", "mtry", "replace",
                                                                              "nodesize", "classwt", "cutoff",
                                                                              "maxnodes")], function(x) as.factor(x))

# Re-encoding Factors
levels(result$ntree) = c(-1,1)
levels(result$mtry) = c(-1,1)
levels(result$replace) = c(-1,1)
levels(result$nodesize) = c(-1,1)
levels(result$classwt) = c(-1,1)
levels(result$cutoff) = c(-1,1)
levels(result$maxnodes) = c(-1,1)
result <- as.data.frame(sapply(result, function(x) as.numeric(as.character(x)))))
head(result)
```

A model with interaction terms with optimal design

```
# Running Linear Model
cv.model <- lm(CV~.^2, data=result.copy)
summary(cv.model)
```

As seen from the output above, there's no significant factors except the intercept under the level of $\alpha = 0.05$. The difference between the Multiple and Adjusted R^2 values is substantial, which indicates that the high R^2 value could be due to the inclusion of several independent factors. We would thus need a reduced model.

Bayes Information Criterion (BIC)

We use BIC to avoid overfitting.

$$BIC = n \ln\left(\frac{RSS}{n}\right) + \ln(n)k$$

where n is the number of observations, k is the number of predictors, and \$RSS\$ stands for residual square sum.

- The best model has the smallest BIC value.
- For n larger than 8, $\log(n) > 2$, and so this is a greater penalty for complexity than AIC.
- BIC favors simpler models than AIC.

```
library(leaps)
Best_Subset <- regsubsets(CV~.^2, data = result,
                          nbest = 1, nvmax=NULL,
                          force.in = NULL, force.out=NULL,
                          method = "exhaustive",
                          really.big=F)
summary_best_subset <- summary(Best_Subset)
plot(summary_best_subset$bic, xlab = "Number of Variables", ylab = "BIC", type = "l")
min = which.min(summary_best_subset$bic)
points(min, summary_best_subset$bic[min], col = "red", cex = 2, pch = 20)
```

We now select the model with the lowest BIC value.

```
### Select Model with least BIC score
# We first select the best model according to BIC.
modelwith.minimum.BIC <- which.min(summary_best_subset$bic)
best.model <- summary_best_subset$which[modelwith.minimum.BIC,][-1]
# Only keep the predictors are indicated by 'TRUE' and our response 'goldDiff'
keep <- names(best.model[best.model==T])
# Print the selected factors
a <- paste(keep, sep=" ", collapse=" + ")
a
```

We now use these predictors to construct a reduced model.


```
# construct a reduced model based on BIC criteria
cv.model.bic <- lm(CV~classwt + cutoff + maxnodes + classwt:cutoff + classwt:maxnodes, data=result)
# T-test
summary(cv.model.bic)
```

The difference between the Multiple and Adjusted R^2 values is now lower. Under $\alpha = 0.05$, all predictors, classwt, cutoff, maxnodes, the interaction between classwt and cutoff, and the interaction between classwt and maxnodes are significant factors. We now analyse the residuals.

Analysis of Residuals

```
# Analysis of Residuals
cv.res <- cv.model.bic$residuals
cv.fit <- cv.model.bic$fitted.values
par(mfrow = c(1,3))
qqnorm(cv.res)
qqline(cv.res)
plot(cv.fit,cv.res, main = "Residuals vs Fitted Values",
     xlab = "Predicted", ylab = "Residuals")

N <- nrow(result)
plot(x = 1:N, y = cv.res, xlab = "Run Order",
     ylab = "Residuals")
```

We see that the points on the normal Q-Q Plot adhere to the qqline except at extremities. The residuals in the residuals vs fitted are evenly distributed across each fitted value. And the distribution of residuals in the residuals vs run order plots are fairly homoskedastic. Thus, we say for our final modal the assumption of constant variance, independence, and normality are satisfied.

Running Confirmation Tests

Since we can run 35 tests in total and we have already run 32 times, we have 3 runs left for confirmation experiment.

The final fitted equation is

$$\hat{y} = 0.622765 - 0.048193E + 0.028372F + 0.015964G + 0.023461EF + 0.037590EG$$

where \hat{y} is the predicted response and A=ntree, B=mtry, C=replace, D=nodesize, E=classwt, F=cutoff, and G=maxnodes.

Using this equation, we can find the settings that maximize the porosity percentage. To this end, we use the 'optim' function in R. We first define an objective function which we want to maximize.

```
obj_func <- function(x){
  pred.y <- 0.622765 - 0.048193 * x[5] + 0.028372 * x[6] + 0.015964*x[7] + 0.023461 *x[5]*x[6] + 0.037590*x[5]*x[7]
  return(-1*pred.y) # Because the 'optim' function minimizes.
}
```

Now, we use the *optim* function in R.

```
# First, use optim() to find the optimal setting according to our model
optim(par = c(-1, -1, -1, -1,-1,-1,-1), fn = obj_func, lower = -1, upper = 1, method = "L-BFGS-B")
```

```
# Then, construct a dataframe of 7 columns, 1 row, with the optimal setting above
design.confirm <- data.frame(ntree = -1,
                             mtry = -1,
                             replace = -1,
                             nodesize = -1,
                             classwt = -1,
                             cutoff = 1,
                             maxnodes = -1)
head(design.confirm)
```

```
# Then, Get the prediction interval of the optimal setting based on our model
predict(cv.model.bic, newdata = design.confirm, interval = "predict")
```

```
# Convert the optimal setting to numerical values
design.confirm <- data.frame(ntree = rep(100,3),
                             mtry = rep(2,3),
                             replace = rep(0,3),
                             nodesize = rep(1,3),
                             classwt = rep(0.5,3),
                             cutoff = rep(0.8,3),
                             maxnodes = rep(10,3))
head(design.confirm)
```

```
# Lastly, run cv.rf() three times with optimal setting (with numerical data) to confirm
confirm.result <- cv.rf(design.confirm, y, X)
confirm.result$CV
mean(confirm.result$CV)
```

We have the results of confirmation experiments are in the range of our prediction interval (0.608844, 0.7895945).

Moreover, the mean of the results (0.6992193) is close to the fitted value 0.6970846.

5. Statement of Contribution

Sebastian: Wrote and compiled the report.

Ovie: Helped with analysis and report-writing

Qianli: Led the design and analysis

6. References

1. Vazquez, Alan. Final Project Spring 2022. STATS 101B.
2. By: IBM Cloud Education. "What Is Random Forest?" *IBM*,
<https://www.ibm.com/cloud/learn/random-forest>.
3. Vazquez, Alan. Data Description Spring 2022. STATS 101B.
4. Vazquez, Alan. Lecture 10_1.pdf. STATS 101A. "Bayes Information Criterion (BIC)"