

ST4061 – Statistical Methods for Machine Learning II  
ST6041 – Machine Learning and Statistical Analytics II

2023-24  
Continuous Assessment 2

**Answers to Question 1**

Question	Your answer																														
1	Mean OOB RMSE for 0.001 shrinkage rate – 0.8414 Mean OOB RMSE for 0.05 shrinkage rate – 0.5053 Mean OOB RMSE for 0.01 shrinkage rate – 0.6195 Mean OOB RMSE for 0.1 shrinkage rate – 0.4933																														
2	<div>Boxplot of OOB-RMSEs across shrinkage rates</div> <table><caption>Approximate data from the boxplot</caption><tr><th>Shrinkage-values</th><th>Min</th><th>Q1</th><th>Median</th><th>Q3</th><th>Max</th></tr><tr><td>0.001</td><td>0.74</td><td>0.82</td><td>0.84</td><td>0.87</td><td>0.94</td></tr><tr><td>0.05</td><td>0.34</td><td>0.46</td><td>0.50</td><td>0.54</td><td>0.67</td></tr><tr><td>0.01</td><td>0.53</td><td>0.59</td><td>0.62</td><td>0.65</td><td>0.72</td></tr><tr><td>0.1</td><td>0.35</td><td>0.45</td><td>0.49</td><td>0.53</td><td>0.61</td></tr></table>	Shrinkage-values	Min	Q1	Median	Q3	Max	0.001	0.74	0.82	0.84	0.87	0.94	0.05	0.34	0.46	0.50	0.54	0.67	0.01	0.53	0.59	0.62	0.65	0.72	0.1	0.35	0.45	0.49	0.53	0.61
Shrinkage-values	Min	Q1	Median	Q3	Max																										
0.001	0.74	0.82	0.84	0.87	0.94																										
0.05	0.34	0.46	0.50	0.54	0.67																										
0.01	0.53	0.59	0.62	0.65	0.72																										
0.1	0.35	0.45	0.49	0.53	0.61																										

3	Shrinkage parameter recommended for training this model (with 100 trees) is 0.1. The mean OOB RMSE for models with shrinkage parameters 0.05 and 0.1 is very similar (0.5053 & 0.4933 respectively), however the spread of the OOB RMSEs is lower for the model with 0.1 as shrinkage parameter compared to that of 0.05 indicating better consistency. Since, the number of trees is not too high over here, a higher shrinkage rate (0.1) has resulted in a better training making it the final choice.
---	---

### **R code for Question 1**

```
## Question 1
rm(list=ls())
require(gbm)
require(ISLR)
df = na.omit(Hitters)
df$Salary = log(df$Salary)
rates = c(0.001, 0.05, 0.01, 0.1)
L = length(rates)
set.seed(4061)

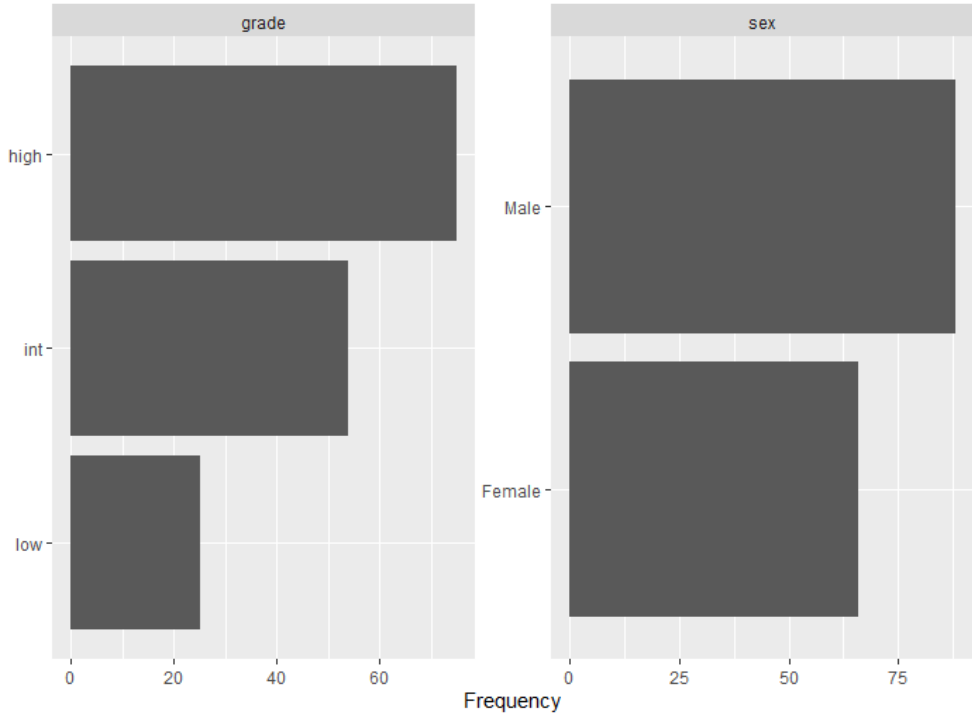
n = nrow(df)
B = 100

OOB_RMSEs = matrix(NA, nrow = B, ncol = L)
colnames(OOB_RMSEs) = rates
for (i in 1:L) {
  for (j in 1:B){
    idxs = sample(1:n, n, replace=TRUE)
    X_train = df[idxs,]
    X_test = df[-idxs,]
    Y_test = df[-idxs,]$Salary
    gbm_tree = gbm(Salary~., data=X_train, distribution = 'gaussian', shrinkage = rates[i])
    test_preds = predict(gbm_tree, X_test)
    OOB_RMSEs[j, i] = sqrt(mean((test_preds - Y_test)^2))
  }
}
head(OOB_RMSEs)

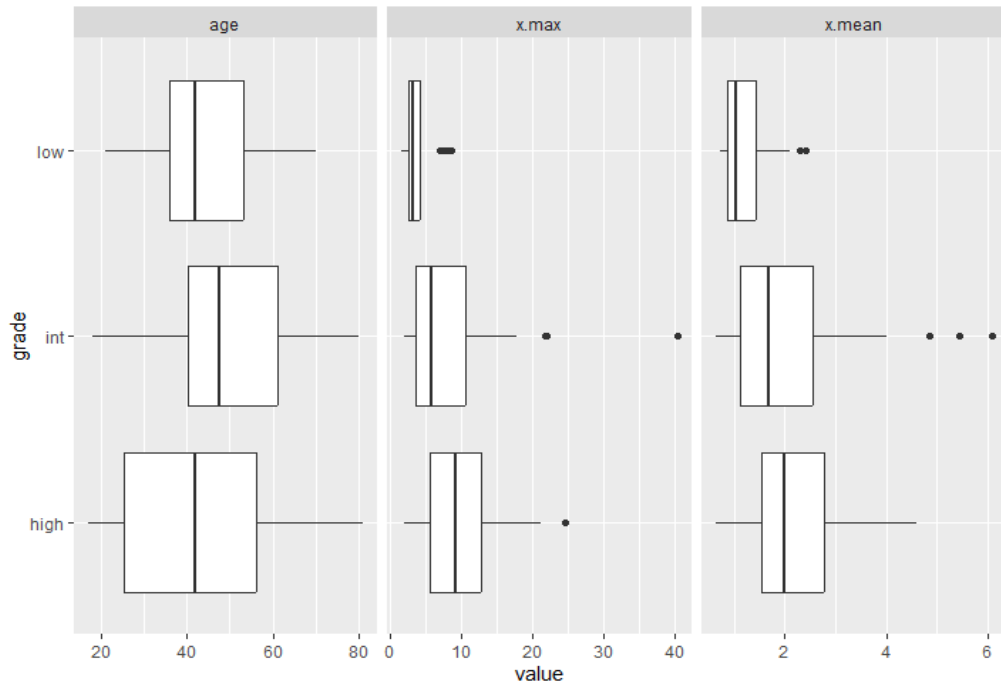
## Question 1(1)
OOB_RMSEs_mean = apply(OOB_RMSEs, 2, mean)
round(OOB_RMSEs_mean, 4)

## Question 1(2)
par(mfrow=c(1,1))
boxplot(OOB_RMSEs, main="Boxplot of OOB-RMSEs across shrinkage
rates", xlab="Shrinkage-values", ylab="Bootstrap RMSE-estimates", col = 'cyan')
```

## Answers to Question 2

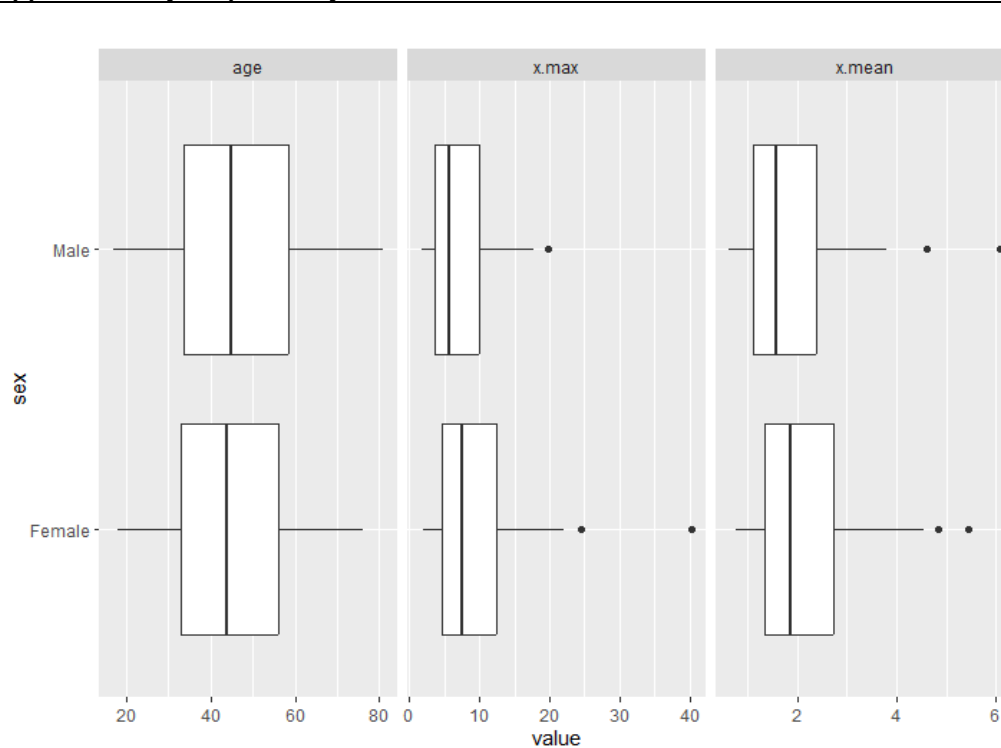
Question	Your answer
1	<p>“Sex” is the only categorical feature here with 2 levels (“Female”, “Male”). Grade is also a categorical column, however it’s not a feature here but the outcome (labels).</p>
2	 <p>The figure consists of two horizontal bar charts. The left chart, titled 'grade', shows the frequency of three categories: 'high' (approx. 70), 'int' (approx. 55), and 'low' (approx. 25). The right chart, titled 'sex', shows the frequency of two categories: 'Male' (approx. 88) and 'Female' (approx. 65). Both charts have a y-axis with category labels and an x-axis labeled 'Frequency' with numerical scales (0-60 for grade, 0-75 for sex).</p> <p>There are more number of Males (approx. 88) compared to the number of females (approx. 65) in the data indicating a gender imbalance in the data present. Also, there is a clear imbalance in the data recorded for the different grade of tumours with category “high” dominating and “low” having the least amount of samples present.</p>

3



Overall, the values for x.max increase with change from grades from low to int to high respectively. A similar trend is observed in x.mean across grades as well. There is a wide age distribution for “high” grade tumors, whereas for “int” and “low” grade most number of samples are in the age group of 40-60 and 30-50 approximately respectively.

4



	<p>The distribution of age for both males and females are very similar to each other with a median age value slightly higher for males. Overall, the values for x.max column is slightly on the higher side for females compared to males with both having wider distribution and a higher median value. For the x.mean column, values are right skewed for both females and males with a wider distribution for females and a higher median value compared to males. Highest x.mean value is for a male (approx.. 6), which is an outlier according to the graph.</p> <p>However, one important factor to remember here is the number of data points for males and females is NOT the same as the data consists of more males recorded than females indicating a clear gender imbalance.</p>						
5	<p>p-value for Wilcoxon test between Age and Sex – <b>0.7259</b> p-value for Wilcoxon test between Age and X.Max – <b>0.07627</b> p-value for Wilcoxon test between Age and X.Mean – <b>0.07691</b></p> <p>The Wilcoxon test between Age and Sex yielded a p-value of &gt;0.05 meaning there is no statistically significant difference between distribution of age of males and that of females. (Fail to reject the Null hypothesis)</p> <p>The Wilcoxon test between Age and X.Max yielded a p-value of &gt;0.05 meaning there is no statistically significant difference between distribution of X.Max of males and that of females. (Fail to reject the Null hypothesis)</p> <p>The Wilcoxon test between Age and X.Mean yielded a p-value of &gt;0.05 meaning there is no statistically significant difference between distribution of X.Mean of males and that of females. (Fail to reject the Null hypothesis)</p>						
6	<p>Mean of Age after Min-Max Scaling – <b>0.4465</b> Mean of X.Mean after Min-Max Scaling – <b>0.2477</b> Mean of X.Max after Min-Max Scaling – <b>0.1655</b></p>						
7	Overall Error quoted by neuralnet function used - <b>7.543743</b> (seed set as 4061)						
8	<p>(i) Overall Accuracy – <b>0.92208 (92.208%)</b></p> <p>(ii) Class-wise Sensitivity for this fit –</p> <table><tr><td>“high”</td><td>“int”</td><td>“low”</td></tr><tr><td><b>0.9733</b> <b>(97.33%)</b></td><td><b>0.963</b> <b>(96.3%)</b></td><td><b>0.68</b> <b>(68%)</b></td></tr></table>	“high”	“int”	“low”	<b>0.9733</b> <b>(97.33%)</b>	<b>0.963</b> <b>(96.3%)</b>	<b>0.68</b> <b>(68%)</b>
“high”	“int”	“low”					
<b>0.9733</b> <b>(97.33%)</b>	<b>0.963</b> <b>(96.3%)</b>	<b>0.68</b> <b>(68%)</b>					
9	<p><b>26</b> columns identified by the correlation filter to remove with 0.95 as cutoff –</p> <p>[‘compactness1’, ‘compactness2’, ‘sphericity’, ‘l.major’, ‘major.axis.length’, ‘l.minor’, ‘minor.axis.length’, ‘l.least’, ‘least.axis.length’, ‘RMS’, ‘mean_HIST’, ‘sum.avg_GLCM’, ‘auto.corr_GLCM’, ‘var_HIST’, ‘joint.var_GLCM’, ‘sum.var_GLCM’, ‘joint.max_GLCM’, ‘energy_GLCM’, ‘entropy_GLCM’, ‘homogeneity_GLCM’, ‘inv.diff.mom_GLCM’, ‘diff.entropy_GLCM’, ‘dissimilarity_GLCM’, ‘homogeneity.norm_GLCM’, ‘contrast_GLCM’, ‘inv.diff.mom.norm_GLCM’]</p>						

## **R code for Question 2**

```
## Question 2
require(caret)
require(neuralnet)
require(DataExplorer)
df = read.csv(file="uws.csv", stringsAsFactors=TRUE)
subdf = df[,c("grade","sex","age","x.mean","x.max")]
y = df$grade
x = df
x$grade = NULL
```

```
## Question 2(1)
## sex is a categorical feature
str(df)
```

```
## Question 2(2)
plot_bar(df)
```

```
## Question 2(3)
plot_boxplot(subdf, by = 'grade')
```

```
## Question 2(4)
plot_boxplot(subdf, by = 'sex')
```

```
## Question 2(5)
age = subdf$age
sex = subdf$sex
x_max = subdf$x.max
x_mean = subdf$x.mean
```

```
wilcox.test(age ~ sex, alternative = "two.sided")
wilcox.test(x_max ~ sex, alternative = "two.sided")
wilcox.test(x_mean ~ sex, alternative = "two.sided")
```

```
## Question 2(6)
conversion <- function(x){
  # function recoding levels into numerical values
  if(is.factor(x)){
    levels(x)
    return(as.numeric(x))
  } else {
    return(x)
  }
}
```

```

    }
  }
  scaling <- function(x){
    # function applying normalization to [0,1] scale
    mins = min(x,na.rm=TRUE)
    maxs = max(x,na.rm=TRUE)
    return((x-mins)/(maxs-mins))
  }
  df_inter = data.frame(lapply(df,conversion))
  df_scaled = data.frame(lapply(df_inter,scaling))
  means = apply(df_scaled[, c("age", "x.mean", "x.max")], 2, mean)
  round(means, 4)
  df_scaled$grade = NULL

## Question 2(7)
set.seed(4061)
mod = neuralnet(y~., data = df_scaled, hidden=c(5), linear.output = FALSE)
(error = mod$result.matrix["error",])

## Question 2(8)
col_names = colnames(mod$response)
final_preds = as.factor(col_names[max.col(predict(mod, df_scaled))])
cf_mat = caret::confusionMatrix(final_preds, y)
(overall_accuracy = cf_mat$overall[1])
specificity_class = cbind(cf_mat$byClass[1],cf_mat$byClass[2],cf_mat$byClass[3])
colnames(specificity_class) = c("high","int","low")
round(specificity_class, 4)

## Question 2(9)
x$sex = as.numeric(x$sex)
correlation_matrix <- cor(x)
cols = colnames(x)
pairs = c()
for (i in cols) {
  for (j in cols) {
    if (i!=j & abs(correlation_matrix[i,j]) > 0.95) {
      pairs = c(pairs, c(i,j))
    }
  }
}
cat("Columns to remove : ", unique(pairs))
unique(pairs)

```