

Task 1

Read the dataset

Code:

```
setwd('/Users/duwonha/Desktop/University/Second year/T2/Big data/Assignment2')  
drugdata <- read.csv('drug200.csv', stringsAsFactors = TRUE)
```

Task 2

Dataset split into training and test sets

Code:

```
training_size <- floor(0.8 * nrow(drugdata))  
set.seed(101)  
train_int <- sample(seq_len(nrow(drugdata)), size = training_size)  
trainingSet <- drugdata[train_int, ]  
testSet <- drugdata[-train_int, ]
```

Task 3

Importing library

Code:

```
library(ISLR)  
data(package = "ISLR")  
require(tree)
```

Tree Construction

Code:

```
tree_accuracy <- tree(formula = Drug ~ Age+Sex+BP+Cholesterol+Na_to_K, data =  
trainingSet)
```

Tree Plot

Code:

```
plot(tree_accuracy)  
text(tree_accuracy, pretty = 0)
```

Accuracy Calculation

Code:

```
tree_pred = predict(tree_accuracy, drugdata[-train_int,], type = 'class')  
with(drugdata[-train_int,], table(tree_pred, Drug))
```

Task 4

Cross-validation and Tree Construction

Code:

```
drug.cv = cv.tree(tree_accuracy, FUN = prune.misclass)  
drug.cv > plot(drug.cv)  
drug.cv
```

#Once you run the drug.cv, you can see deviation of size 6 is 1 which means size 6 is the most accurate one. Size 6 will be used for making prune. (best = 6)

```
$size
```

```
[1] 6 4 3 2 1
```

```
$dev
```

```
[1] 1 19 26 43 84
```

```
prune.drug = prune.misclass(tree_accuracy, best = 6)
```

Tree plot

Code:

```
plot(prune.drug)
```

```
text(prune.drug, pretty=0)
```

Accuracy Calculation

Code:

```
tree_pred = predict(prune.drug, drugdata[-train_int,], type = "class")
```

```
with(drugdata[-train_int,], table(tree_pred, Drug))
```

Screenshot







