

# Configuration Manual

MSc Research Project Data Analytics

## Harshitha Deenadayal Student ID: x18136231

School of Computing National College of Ireland

Supervisor: Dr. Pierpaolo Dondio

#### National College of Ireland Project Submission Sheet School of Computing



Student Name:	Harshitha Deenadayal
Student ID:	x18136231
Programme: Data Analytics	
Year:	2019
Module: MSc Research Project	
Supervisor:	Dr. Pierpaolo Dondio
Submission Due Date:	13/12/2019
Project Title:	Configuration Manual
Word Count:	1176
Page Count:	13

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	12th December 2019

#### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).□Attach a Moodle submission receipt of the online project submission, to<br/>each project (including multiple copies).□You must ensure that you retain a HARD COPY of the project, both for□

your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

## Configuration Manual

Harshitha Deenadayal x18136231

### 1 Introduction

This manual will explain system set-up including hardware and software requirements and the codes in python and pyspark programming language for the implementation of the research project: 'Distributed Intrusion Detection System for Computer Networks using Hyper-tuned Gradient Boosting algorithm on Spark Framework.'

## 2 System configuration

#### 2.1 Hardware

Processor: Intel(R) Core(TM) i7-5500U CPU@2.4GHz GPU: NVIDIA GeForce RAM:16GB Storage: 1 TB HDD; Operating system: Windows 10, 64-bit.

#### 2.2 Software for centralized environment

- Centralized environment on linux virtual machine. With Ubuntu 64-bit OS, 2 CPU processors and 11GB RAM.
- Anaconda Navigator (version 3) to launch jupyter notebook on internal explorer.
- Python(v3) programming language for data pre-processing, data mining modelling and visualization, on jupter notebook.

#### 2.3 Software for distributed environment

- Distributed environment set-up using Apache Hadoop and Apache Spark on a linux virtual machine. With Ubuntu 64-bit OS, 2 CPU processors and 11GB RAM.
- Java jdk version 8 and SSH (secure shell) keys are installed on Ubuntu OS.
- Apache Hadoop version 2.9.2 is downloaded from apache website<sup>1</sup>, installed and configured. And the path is set between hadoop and java.

 $<sup>^{1}</sup> https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-2.9.2/hadoop-2.9.2.tar.gz$ 

```
Hadoop Installation on Ubuntu
1) Java installation:
       $sudo apt install openjdk-8-jdk
       $sudo apt install rsync
2)set the path between Hadoop and java on Ubuntu and set this java file as source
$sudo nano /etc/profile
       export JAVA_HOME=/usr
       $source /etc/profile
3) configure SSH keys (secure shell)
    -Add a new user account hduser to the hadoop group and set password
       $sudo addgroup hadoopgroup
       $sudo adduser -ingroup hadoopgroup hduser
    -install ssh
       $sudo apt-get install ssh
    -enable ssh
       $sudo systemctl enable ssh
    -start ssh
       $sudo systemctl start ssh
4) Genearate keys
       $su - hduser
       $ssh-keygen -t rsa -P "
5)download hadoop from ubuntu at hadoop.apache.org
$ wget http://ftp.heanet.ie/mirrors/www.apache.org/dist/hadoop/common/hadoop-2.9.2/hadoop-2.9.2.tar.gz
       $cd /home/hduser
       $sudo mv ./hadoop-2.9.2 /usr/local
       $sudo chown -R hduser:hadoopgroup /usr/local/hadoop-2.9.2/
6) Hadoop configuration
   - login to hduser
         $nano ./.bashrc
# Hadoop config
export HADOOP_PREFIX=/usr/local/hadoop
export HADOOP_HOME=/usr/local/hadoop
export HADOOP_MAPRED_HOME=${HADOOP_HOME}
export HADOOP_COMMON_HOME=${HADOOP_HOME}
export HADOOP_HDFS_HOME=${HADOOP_HOME}
export YARN_HOME=${HADOOP_HOME}
export HADOOP_CONF_DIR=${HADOOP_HOME}/etc/hadoop
# Native path
export HADOOP COMMON LIB NATIVE DIR=${HADOOP PREFIX}/lib/native
export HADOOP_OPTS="-Djava.library.path=$HADOOP_PREFIX/lib/native"
# Java path
export JAVA_HOME="/usr"
# OS path
export PATH=$PATH:$HADOOP_HOME/bin:$JAVA_PATH/bin:$HADOOP_HOME/sbin
```

```
7) change the directory to,
$cd /usr/local/hadoop
$cd /etc/hadoop
$sudo nano core-site.xm
```

```
<configuration>
<property>
<name>fs.default.name</name>
<value>hdfs://localhost:9000</value>
</property>
</configuration>
```

\$sudo nano hdfs-site.xml

```
<configuration>
<property>
<name>dfs.replication</name>
<value>l</value> </property>
<property>
<name>dfs.name.dir</name>
<value>file:/usr/local/hadoop/hadoopdata/hdfs/namenode</value>
</property>
<property>
<name>dfs.data.dir</name>
```

```
</value>file:/usr/local/hadoop/hadoopdata/hdfs/datanode</value>
</property>
</configuration>
```

```
$sudo nano mapred-site.xml
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
        <value>yarn</value>
        </property>
</configuration>
```

```
$$sudo nano yarn-site.xml
<configuration>
  <property>
    <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle</value>
        </property>
</configuration>
```

```
8)move the main hadoop to,
$cd ../..
$cd bin
$sudo hdfs namenode -format
```

🛃 Harshitha [Running] - Oracle VM VirtualBox

File M	achine View Input Devices Help
Activit	ies 🗖 Terminal 👻 Thu 08:19 🗉
	hduser@harshitha-virtualbox: /usr/local
	File Edit View Search Terminal Help
(2)	<pre>(base) hduser@harshitha-virtualbox:-\$ start-all.sh This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh Starting namenodes on [localhost] localhost: starting namenode, logging to /usr/local/hadoop-2.9.2/logs/hadoop-hduser-namenode-harshitha-virtualbox.out localhost: starting datanode, logging to /usr/local/hadoop-2.9.2/logs/hadoop-hduser-datanode-harshitha-virtualbox.out Starting secondary namenodes [0.0.0.0]</pre>
0	<pre>bloc.or starting secondarynamemode, togging to /usr/tocat/hadoop-2.9.2/togs/hadoop-houser-secondarynamemode-harshitha-virtualbox.out starting resourcemanager, logging to /usr/local/hadoop/logs/yarn-hduser-resourcemanager-harshitha-virtualbox.out localhost: starting nodemanager, logging to /usr/local/hadoop-2.9.2/logs/yarn-hduser-nodemanager-harshitha-virtualbox.out (base) hduser@harshitha-virtualbox:-5 jps</pre>
	16176 NodeManager 15575 DataNode 15818 SecondaryNameNode 15419 NameNode
A	16494 Jps 16015 ResourceManager (base) <mark>hduser@harshitha-virtualbo</mark> x:~\$ hadoop fs -ls / Found 3 items
?	drwxr-xr-x       - hduser supergroup       0       2019-08-12       16:46       /sqoop         drwxr-wx-wx       - hduser supergroup       0       2019-08-13       13:45       /tmp         drwxr-xr-x       - hduser supergroup       0       2019-08-12       17:03       /user         (base)       hduser@harshitha-vtrtualbox:-S       cd       /usr/local       1000000000000000000000000000000000000
<u>a</u> ,	(base) hduser@harshitha-virtualbox:/usr/local\$ hadoop fs -put /home/hduser/Downloads/finaldata.csv / (base) hduser@harshitha-virtualbox:/usr/local\$ hadoop fs -ls / Found 4 items
2	Twr-rr-r       Induser supergroup       2504/282 2019-12-12 08:18 /rthaldata.csv         drwxr-xr-x       -hduser supergroup       0 2019-08-12 10:46 /sqoop         drwxr-xr-x       -hduser supergroup       0 2019-08-12 10:46 /sqoop         drwxr-xr-x       -hduser supergroup       0 2019-08-12 10:46 /sqoop         drwxr-xr-x       -hduser supergroup       0 2019-08-12 11:21 10:01         (base)       hdusershttha-vtruatbox:/usr/local5

• Apache spark 2.4.3 is downloaded and installed<sup>2</sup>.

```
Apache Spark
1) Install Spark under /usr/local
        $ sudo cp ~/Downloads/spark-2.4.3-bin-hadoop2.7.tgz .
$ sudo tar xzf spark-2.4.3-bin-hadoop2.7.tgz
2) Create a symbolic link named spark
        $ sudo ln -s spark-2.4.3-bin-hadoop2.7 spark
3) Change the ownership of the files
        $ sudo chown -R hduser:hadoopgroup spark
4) Download from website and install Anaconda with Python 2.7 for Linux
        $ sudo bash Anaconda2-2019.03-Linux-x86_64.sh
5) Spark environment variables to the .bashrc file from /home/hduser
export SPARK_HOME=/usr/local/spark
export PATH=$PATH:$SPARK_HOME/bin
export PYSPARK_PYTHON=/usr/local/anaconda2/bin/python2
export PYSPARK_DRIVER_PYTHON=jupyter
export PYSPARK_DRIVER_PYTHON_OPTS="notebook"
6) Source bashrc
        $ source ./.bashrc
```

<sup>&</sup>lt;sup>2</sup>https://spark.apache.org/releases/spark-release-2-4-3.html

🏏 Harsh	itha [Running] - Oracle VM VirtualBox
File Ma	achine View Input Devices Help
Activiti	es 🖻 Terminal 👻 👘 Thu 08:25 🛛
	hduser@harshitha-virtualbox: /usr/local/spark
• 🔁	File Edit View Search Terminal Help
	(base) hduser@harshitha-virtualbox:~\$ cd /usr/local/spark
	(base) hduser@harshitha-virtualbox:/usr/local/spark\$ sbin/start-master.sh
	starting org.apache.spark.deploy.master.master, logging to /usr/local/spark/logs/spark-nduser-org.apache.spark.deploy.master (base) huusendharshitha-virtualbox:/usr/local/sparks ins
	16176 NodeManager
•	16849 Jps
	155/5 DataNode 15818 SecondarvNameNode
	15419 NameNode
	16812 Master
-	16015 ResourceManager
	(Dase) nousergnarshitha-virtuaioox:/usr/usr/usa/sparks pyspark [T 88:25:66 748 NotehookAnol The port 8888 is already in use trying another port.
	[1 08:25:07.073 NotebookApp] JupyterLab extension Loaded from /home/hduser/anaconda3/lib/python3.7/site-packages/jupyterLab
_	[I 08:25:07.073 NotebookApp] JupyterLab application directory is /home/hduser/anaconda3/share/jupyter/lab
-8-	[I 08:25:07.094 NotebookApp] Serving notebooks from local directory: /usr/local/spark-2.4.3-bin-hadoop2.7
A	[1 06:25:07.094 NOTEDOOKAPJ] THE JUDYLET NOTEDOOK IS TUMINING AT: [T 08:25:07.095 NOTEDOOKAD] http://localbost.18889/170ken=23d9a8bbda128efb5b51b389e85ca188d1d5f15b946a11fa
	[1 08:25:07.096 NotebookApp] or http://127.0.0.1:8889/?token=23d9a8bbda128efb5b51b389e85ca188d1d5f15b946a11fa
	[I 08:25:07.096 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
	[C 08:25:07.407 NotebookApp]
	To access the notebook, open this file in a browser:
a	file:///home/hduser/.local/share/jupyter/runtime/nbserver-16874-open.html
	Or copy and paste one of these URLs:
	or http://10.dtinost.ass//tuben=2309a8b00a128eFb351b389e85ca188d1051150940a11fa

• Anaconda with python 2.7 on spark environment, to launch jupyter notebook on internal explorer.

### 3 Project Development

The development of this project happens in several stages namely: data pre-processing, feature engineering, modelling, evaluation and cross validation. The experiments are conducted on both centralized and distributed environments to compare the time-taken to train the models. In this research project, data pre-processing stage includes exploratory data analysis and data transformations. Two steps of feature engineering is carried out and they are: Feature reduction and feature selection. Feature reduction is based on association rule mining technique where feature selection is based on variable importance graph plotted using random forest. The dataset is split into train and test samples. Data mining models are build using various algorithms like SVM, decision tree, random forest, gradient boosting and an ensemble of all these four models. Parameters of all algorithms are hyper tuned to estimate the best fit to achieve best accuracy. All the models are trained with best-fit estimators using train dataset and validated on unknown test instances. A confusion matrix is plotted for measuring the performance of each model. Using the pre-defined function of python accuracy, precision, recall and F1 score are calculated. Finally, all the models are cross validated to ensure there are no overfitting issues.

The above explained project development process is carried out similarly on distributed (apache hadoop and apache spark) environment as well and time-taken(in-minutes) to train the models are compared with centralized environment without hadoop and spark set-up.

## 4 Code in Python

#### 4.1 Libraries

In [ ]:	<pre>import numpy as np import pandas as pd import seaborn as sns from pprint import pprint import statistics</pre>	#to performs operations on #for data manipulations #visualiation library for #provides data structures #for calculation	n array statistical graph in a input interpreter form
	<pre>from sklearn.preprocessing import LabelEncoder from sklearn.multioutput import MultiOutputClas from sklearn.model_selection import train_test from sklearn.model_selection import cross_val_s from sklearn.model_selection import GridSearchC</pre>	sifier split core, KFold,StratifiedKFolo V	#for converting cat variables into numeric #for hybrid model d #for cross validation #for hyper-parameter tuning
	<pre>from sklearn.tree import DecisionTreeClassifier from sklearn.svm import SVC from sklearn.ensemble import RandomForestClassi</pre>	fier,GradientBoostingClass	ifier
	<pre>from sklearn.metrics import confusion_matrix from sklearn.metrics import precision_score,red from sklearn.metrics import accuracy_score</pre>	call_score,f1_score	#for model evaluation
T- 1 1.			
TU [ ]:	<pre>%stime #Visualization import matplotlib.pyplot as plt</pre>		
	<pre>#Config for plotting graph fig_size = plt.rcParams["figure.figsize"] fig_size[0] = 50 fig_size[1] = 20 plt.rcParams["figure.figsize"] = fig_size plt.rcParams.update({'font.size': 50})</pre>		
	<pre>plt.rcParams.update({'font.size': 50})</pre>		

In the above lines of codes, all the required libraries are imported and parameters for visualization graph are set<sup>3</sup>.

#### 4.2 Data pre-processing



In the above lines of code, dataset is imported. Some Data exploration is performed, followed by transformation of all categorical variables to numeric.

<sup>&</sup>lt;sup>3</sup>https://docs.python.org/3/library/

#### 4.3 Feature Reduction

```
In [ ]: #Feature reduction: remove high correlated independent attributes
# Create correlation matrix
corr_matrix = dataset.corr().abs()
# Select upper triangle of correlation matrix
upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(np.bool))
# Find features with correlation greater than 0.7
to_drop = [column for column in upper.columns if any(upper[column] > 0.7)]
# Drop features
dataset.drop(to_drop, axis=1, inplace=True)
In [ ]: to_drop # following attributes are dropped
In [ ]: dataset.shape #verify the dimension of the dataset
In [ ]: dataset.corr() #checking the correlation between independent variables again
In [ ]: sns.heatmap(dataset.corr()) #verify correlation again
```

In the above lines of codes, dataset is checked for high correlated independent variables and a heat map is plotted for visualization. Highly correlated variables are then dropped.

#### 4.4 Feature Selection



In the above lines of codes, a variable importance graph is plotted using random forest algorithm. Based on the significance of variables in predicting class of unknown instances, only top 10 significant variables are selected for training data mining models.

#### 4.5 Support Vector Machine

In [ ]:	#SVC with hyper parameter tuning
	<pre>svc_param = {'C':[1,10,100,1000],'gamma':[1,0.1,0.01,0.001], 'kernel':['rbf']} #select a range of hyper-parameters svc_grid = GridSearchCV(SVC(),svc_param,n_jobs=-1,refit = True, verbose=2, cv = 10) #run all iterations svc_grid, fit(predictors_train,result_train) svc_best = svc_grid.best estimator svc_best.fit(predictors_train,result_train) svc_pred = svc_best.predict(predictors_test) #validated using test data</pre>
In [ ]:	svc_confusion_matrix =confusion_matrix(result_test,svc_pred)       #draw a confusion to find TP,TN,FP and FN         print (svc_confusion_matrix)       sns.heatmap(svc_confusion_matrix)
	<pre>plt.rcParams["figure.figsize"] = fig_size plt.rcParams.update{{'font.size': 50}} plt.title("SVC Confusion Matrix") print ("Precision score "+str(precision_score(result_test,svc_pred))) #find precision, recall, fl score and Accuracy print ("Recall score "+str(fl score(result_test,svc_pred))) print ("Fl score "+str(fl score(result_test,svc_pred))) print ("Accuracy "+str(faccuracy_score(result_test,svc_pred)))</pre>
	<pre>svck10 = cross_val_score(svc_best, predictors_test, result_test, cv=10) # 10-fold cross-validation svc_avg10 = statistics.mean(svck10) print(svck10) print(svc_avg10)</pre>

In the above lines of code, hyper parameters of SVC algorithm are tuned. 'rbf' kernel was selected since dataset was non-linear. SVC model is trained with train dataset with number of iterations equal to all combinations of hyper parameters and best fit parameters is selected using 'estimator' function. It is then, validated using test dataset. Performance of this model is calculated using confusion matrix. Metric such as precision, recall, f1 score and accuracy are used for performance evaluation.

#### 4.6 Decision Tree

In [ ]:	#decision tree classifier	
	<pre>dt = DecisionTreeClassifier()</pre>	<pre>#base classifier</pre>
	<pre>sample_split_range = list(range(2, 50)) dt param = dict(min_samples_split_sample_split_range)</pre>	#set the parameter
	dt grid = GridSearchCV(dt, dt param, cv=10,n jobs=-1,verbose=2,	<pre>scoring='accuracy') #run all the iterations</pre>
	dt_grid.fit(predictors_train, result_train)	#fit the model
	<pre>dt_best = dt_grid.best_estimator</pre>	#select the best set of hyper-parameters
	dt_pred = dt_best_predict(predictors_test)	#test the model
In [ ]:	<pre>confusion_matrix_dt =confusion_matrix(result_test,dt_pred) </pre>	
	print (confusion_matrix_dt)	
	plt.title("Decision Tree Confusion Matrix")	
	<pre>print ("Precision "+str(precision_score(result_test,dt_pred)))</pre>	
	<pre>print ("Recall "+str(recall_score(result_test,dt_pred))) print ("El score "+str(fl score(result_test_dt_pred)))</pre>	
	print ("Accuracy "+str(accuracy score(result test, dt pred)))	
	<pre>dtk10 = cross_val_score(dt_best, predictors_test, result_test, dt_avg10 = statistics_mean(dtk10)</pre>	cv=10) # 10-fold cross-validation
	print(dtk10)	
	print(dt_avg10)	

In the above lines of code, hyper parameters of decision tree algorithm are tuned. Decision tree model is trained with train dataset with number of iterations equal to all combinations of hyper parameters and best fit parameters is selected using 'estimator' function. It is then, validated using test dataset. Performance of this model is calculated using confusion matrix. Metric such as precision, recall, f1 score and accuracy are used for performance evaluation. A 10-fold cross-validation of the accuracy of the model is performed to ensure there is no over-fitting.

#### 4.7 Random Forest



In the above lines of code, hyper parameters of random forest algorithm are tuned. Random forest model is trained with train dataset with number of iterations equal to all combinations of hyper parameters and best fit parameters is selected using 'estimator' function. It is then, validated using test dataset. Performance of this model is calculated using confusion matrix. Metric such as precision, recall, f1 score and accuracy are used for performance evaluation. A 10-fold cross-validation of the accuracy of the model is performed to ensure there is no over-fitting.

#### 4.8 Gradient Boosting

In [ ]:	#firstly tune the parameter for gradient boosting	
	<pre>min_samples_leaf= [20, 50,100,150] learning_rate= [0.5, 0.4, 0.3, 0.2, 0.1, 0.05, 0.02, 0.01] max_depth= [4, 6, 8]</pre>	#number of splits at leaf nodes #contribution of each tree towards learning #number of internal node
	<pre>tune1=[] for item in learning_rate:     plot1 = GradientBoostingClassifier(learning_rate=item)     plot1.fit(predictors_train,result_train)     tune1.append(plot1.score(predictors_test,result_test))</pre>	<pre>#plotting graph of accuracy vs learning rate</pre>
	<pre>plt.plot(learning_rate,tune1) plt.title("Learning_rate")</pre>	
In [ ]:	<pre>tune2=[] for item in max_depth:     graph2 = GradientBoostingClassifier(max_depth=item)     graph2.fit(predictors_train, result_train)     tune2.append(graph2.score(predictors_test, result_test))     print(item) plt.plct(max_depth,tune2) plt.title("Max_depth")</pre>	<pre># plotting graph of accuracy vs max-depth</pre>
In [ ]:	<pre>tune3=[] for item in min_samples_leaf:     graph3 = GradientBoostingClassifier(min_samples_leaf=item)     graph3.fit(predictors_train, result_train)     tune3.append(graph3.score(predictors_test, result_test))     print(item) plt.plot(min_samples_leaf,tune3) plt.title("Min_samples_leaf")</pre>	<pre># plotting graph of accuracy vs min_samples_leaf</pre>

Above lines of codes, plots Accuracy vs. learning-rate, Accuracy vs. max-depth and Accuracy vs. min-samples-leaf graphs to determine the best values of hyper-parameters for high accuracy.



In the above lines of code, best fit values which are obtained based on graphs plotted are set for hyper-parameters. Random forest model is trained with train dataset with number of iterations equal to all combinations of hyper parameters and best fit parameters is selected using 'estimator' function. It is then, validated using test dataset. Performance of this model is calculated using confusion matrix. Metric such as precision, recall, f1 score and accuracy are used for performance evaluation. A 10-fold cross-validation of the accuracy of the model is performed to ensure there is no over-fitting.

#### 4.9 Hybrid



Using the above lines of codes, an ensemble model of all four models is trained to enhance the overall performance of data mining model. This model is then cross-validated with 10 folds.

## 5 Code in Pyspark

Code written in pyspark for distributed environment (Apache Hadoop and Apache Spark) is similar to code written python, except for minor changes.



In the above lines of codes, a connection to spark is established first. Then, all the required libraries are loaded.



In the above lines of code, dataset is imported from HDFS. Columns names transformed by converted into lowercase, removing special characters and white spaces. All the categorical variables in dataset are converted into numeric.

And the rest of the codes are similar to code written in python and are based on feature engineering performed earlier. The main purpose of experiment is to measure the time taken for training the models.



In [ ]: #SVC with hyper parameter tuning
svc\_param = {'C':[1,10,100,1000],'gamma':[1,0.1,0.01], 'kernel':['rbf']} #select a range of hyper-parameters
svc\_grid = GridSearchCV(SVC(),svc\_param,n\_jobs=-1,refit = True, verbose=2, cv=10) #run all iterations
svc\_grid.fit(predictors\_train,result\_train) #fit the model
svc\_best = svc\_grid.best estimator #stelect the best fit
svc\_best.fit(predictors\_train,result\_train) #train model with best fit parameters
svc\_pred = svc\_best.predict(predictors\_test) #train model with best fit parameters
svc\_pred = svc\_onfusion\_matrix(result\_test,svc\_pred) #draw a confusion to find TP,TN,FP and FN
print (svc\_confusion\_matrix)
print ("Precision score "+str(frecult\_score(result\_test,svc\_pred))) #find precision, recall, fl score, Accuracy
print ("Recall score "+str(frecult\_test,svc\_pred)))
print (svc\_confusion\_matrix)
print (svc\_confusion\_matrix(result\_test,svc\_pred)) #draw a confusion to find TP,TN,FP and FN
print (svc\_confusion\_matrix)
print (svc\_confusion\_matrix(result\_test,svc\_pred)))
print (svc\_confusion\_matrix(result\_test,svc\_pred)))
print (svc\_confusion\_matrix(result\_test,svc\_pred))
print (svc\_confusion\_matrix)
print (svc\_confusion\_matrix(result\_test,svc\_pred)) #draw a confusion to find TP,TN,FP and FN
print (svc\_confusion\_matrix(result\_test,svc\_pred)))
print (svc\_confusion\_matrix(result\_test,svc\_pred)))
print (svc\_confusion\_matrix(result\_test,svc\_pred)))
print (svc\_confusion\_matrix(result\_test,svc\_pred))) #find precision, recall, fl score, Accuracy
print ("Precision score" +str(frecision\_score(result\_test,svc\_pred))) #find precision, recall, fl score, Accuracy
print ("Recall score" +str(recall\_score(result\_test,svc\_pred))) #find precision, recall, fl score, Accuracy
print ("Recall score" +str(recall\_score(result\_test,svc\_pred))) #find precision, recall, fl score, Accuracy
print ("Recall score" +str(recall\_score(result\_test,svc\_pred)))
print ("Accuracy" +str(accuracy\_score(result\_test,svc\_pred)))

In [ ]:	<pre>from sklearn.tree import DecisionTreeClassifier #decision tree classifier</pre>	
	<pre>dt = DecisionTreeClassifier()</pre>	#base classifier
	<pre>sample_split_range = list(range(2, 50))</pre>	#set the parameter
	<pre>dt_param = dict(min_samples_split=sample_split_range) dt_param = dict(min_samples_split=sample_split_range)</pre>	2
	<pre>dt_grid = GridSearch(V(dt, dt_param, cv=10,n_jobs=-1,verbose dt_grid fit(predictore_train_result_train)</pre>	=2, scoring= accuracy) #run all the iterations
	dt_bost = dt_grid_bost_ostimator	#ril the model #soloct the best set of byper parameters for learning
	dt best fit(predictors train result train)	#train the model
	dt pred = dt best.predict(predictors test)	#test the model
	<pre>confusion_matrix_dt =confusion_matrix(result_test,dt_pred) print (confusion_matrix_dt) print ("Precision_"+str[precision_score(result_test,dt_pred]) print ("Recall "+str(result_score(result_test,dt_pred])) print ("Fi score "+str(fl_score(result_test,dt_pred))) print ("Accuracy "+str(accuracy_score(result_test,dt_pred)))</pre>	))

In [ ]:	<pre>from sklearn.ensemble import RandomForestC from pprint import pprint #RandormForest classifier with and without</pre>	assifier #provides data structures in a input interpreter form hyper parameter tuning
	<pre>rf = RandomForestClassifier() pprint(rf.get_params()) rf_param = {     'bootstrap': [True],     'max_depth': [80, 90, 100, 110],     'max_features': [2, 3],     'min_samples_leaf': [3, 4, 5],     'min_samples_split': [8, 10, 12],     'n_estimators': [100, 200, 300, 1000] }</pre>	<pre>#create base rf classifier #check the default parameters #tune the parameters #sampling are drawn with replacement #number of internal node #number of attributes considered for splitting #number of splits at leaf nodes #number of splits at each internal node</pre>
	<pre>rf_grid = GridSearchCV(estimator = rf, para cv = 10, n_jobs = rf_grid.fit(predictors_train, result_train) rf_best = rf_grid.best_estimator_ rf_final = rf_best.fit(predictors_train,res rf_pred = rf_final.predict(predictors_test)</pre>	<pre>am_grid = rf_param, -1, verbose = 2) #iterations all combination of hyper-parameter #fit the model # find the best estimator sult_train) ##train the model with best fit estimators #validate the model using test data</pre>
	<pre>print ("Precision score "+str(precision_sco print ("Recall score "+str(recall_score[res print ("Fl score "+str(fl_score(result_test print ("Accuracy "+str(accuracy score[resu])</pre>	pre(result_test,rf_pred))) sult_test,rf_pred))) :,rf_pred))) :. t test,rf pred)))

In []:	<pre>from sklearn.multioutput import MultiOutputClassifier from sklearn.ensemble import VotingClassifier hybrid_classifier = MultiOutputClassifier(VotingClassifier(estimators=</pre>
In [ ]:	<pre>hybrid_classifier.fit(predictors_train,result_train) #fit the hybrid model pred_hybrid=hybrid_classifier.predict(predictors_test) #validate the hybrid</pre>
In [ ]:	<pre>model1=[] for item in pred_hybrid:model1.append(item[0]) #add hybrid model in the table</pre>
In [ ]:	<pre>confusion_matrix_hybrid =confusion_matrix(np.array(result_test['class']),np.array(modell)) print (confusion_matrix_hybrid) print ('Precision "+str(precision score(np.array(result_test['class']),np.array(modell)))) print ("Recall "+str(recall_score(np.array(result_test['class']),np.array(modell)))) print ('Fl score "+str(fl_score(np.array(result_test['class']),np.array(modell)))) print ("Accuracy "+str(accuracy_score(np.array(result_test['class']),np.array(modell))))</pre>