



DEPARTMENT OF
COMPUTER SCIENCE AND ENGINEERING

Title: Ensemble Learning



GREEN UNIVERSITY OF BANGLADESH

1 Objective(s)

- To gather knowledge of ensemble learning.

2 Ensemble Learning

Ensemble means a group of elements viewed as a whole rather than individually. An Ensemble method creates multiple models and combines them to solve it. Ensemble methods help to improve the robustness/generalization of the model. In this article, we will discuss some methods with their implementation in Python.

3 Types of Ensemble Learning

Basic Ensemble Techniques

- Max Voting
- Averaging
- Weighted Average

Advanced Ensemble Techniques

- Bagging
- Stacking
- Blending
- Boosting

Algorithms based on Bagging and Boosting

- Bagging meta-estimator
- Random Forest
- AdaBoost
- GBM

3.1 Basic Ensemble Methods

Multiple predictions are made for each data point in averaging. In this method, we take an average of predictions from all the models and use it to make the final prediction. Averaging can be used for making predictions in regression problems or while calculating probabilities for classification problems.

3.2 Averaging Method and Max Voting Method:

```
1 #import library
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import mean_squared_error
6 from sklearn import tree
7 from sklearn.linear_model import LinearRegression
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.svm import SVC
10
11 -----
12 #import dataset
```

```

13 from google.colab import files
14 f=files.upload()
15
16 -----
17 #read the dataframe
18 df=pd.read_csv("diabetes.csv")
19 df.head()

```

Output Top five rows of the data frame is shown:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Figure 1: Top five rows of the dataset

```

1 #Defining feature and corresponding class label
2 x=df.drop("Outcome",axis="columns")
3 y=df.Outcome
4
5 -----
6 #training set and testing set split
7 X_train, X_test, y_train, y_test = train_test_split(x,y, test_size=0.20)
8
9 # initializing all the model objects with default parameters
10 model_1 = LogisticRegression()
11 model_2 = SVC()
12 model_3 = tree.DecisionTreeClassifier()
13 -----
14 ## training all the model on the training dataset
15 model_1.fit(X_train, y_train)
16 model_2.fit(X_train, y_train)
17 model_3.fit(X_train, y_train)
18 -----
19 # predicting the output on the validation dataset
20 pred_1 = model_1.predict(X_test)
21 pred_2 = model_2.predict(X_test)
22 pred_3 = model_3.predict(X_test)
23 -----
24 #model score/accuracy
25 model_1.score(X_test,y_test)
26 model_2.score(X_test,y_test)
27 model_3.score(X_test,y_test)
28 -----
29 # final prediction after averaging on the prediction of all 3 models
30 pred_final = (pred_1+pred_2+pred_3)/3.0
31
32 # printing the root mean squared error between real value and predicted value
33 print(mean_squared_error(y_test, pred_final))
34 -----

```

```

35 # Making the final model using voting classifier
36 from sklearn.ensemble import VotingClassifier
37 final_model = VotingClassifier(estimators=[('lr', model_1), ('svm', model_2), ('
    dt', model_3)], voting='hard')
38 -----
39 # training all the model on the train dataset
40 final_model.fit(X_train, y_train)
41 -----
42 # predicting the output on the test dataset
43 pred_final = final_model.predict(X_test)
44 -----
45 # printing log loss between actual and predicted value
46 print(log_loss(y_test, pred_final))

```

3.3 Lab Task (Please implement yourself and show the output to the instructor)

- Implement weighted average ensemble technique on the same data set and analyse the output.

3.4 Advanced Ensemble Methods

Above are simple techniques, now let's take a look at advanced techniques for ensemble learning.

3.4.1 Bagging

It is also known as a bootstrapping method. Base models are run on bags to get a fair distribution of the whole dataset. A bag is a subset of the dataset along with a replacement to make the size of the bag the same as the whole dataset. The final output is formed after combining the output of all base models.

```

1 #import library
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import mean_squared_error
6 from sklearn import tree
7 from sklearn.linear_model import LinearRegression
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.svm import SVC
10
11 -----
12 #import dataset
13 from google.colab import files
14 f=files.upload()
15
16 -----
17 #read the dataframe
18 df=pd.read_csv("diabetes.csv")
19
20 -----
21 #Defining feature and corresponding class label
22 x=df.drop("Outcome",axis="columns")
23 y=df.Outcome
24
25 -----
26 #training set and testing set split
27 X_train, X_test, y_train, y_test = train_test_split(x,y, test_size=0.20)
28
29 -----
30 # initializing the bagging model using decision tree as base model
31 from sklearn.ensemble import BaggingClassifier

```

```

32 |
33 | bag_model = BaggingClassifier(
34 |     base_estimator=tree.DecisionTreeClassifier(),
35 |     n_estimators=100,
36 |     max_samples=0.8,
37 |     oob_score=True,
38 |     random_state=0
39 | )
40 | bag_model.fit(X_train, y_train)
41 | bag_model.oob_score_
42 | -----
43 | bag_model.score(X_test, y_test)
44 | -----
45 | #k-fold cross validation
46 | from sklearn.model_selection import cross_val_score
47 | scores = cross_val_score(tree.DecisionTreeClassifier(), x, y, cv=5)
48 | scores

```

4 Conclusion

Based on the focused objective(s), we will learn about the different ensemble methods and comparison between ensemble methods.

5 Lab Exercise (Submit as a report)

- Implement stacking ensemble and AdaBoost ensemble on the same dataset

6 Policy

Copying from internet, classmate, seniors, or from any other source is strongly prohibited. 100% marks will be *deducted* if any such copying is detected.