

# 8. Ensemble Learning in MLE



## 1. Introduction

- ↳ The Ensemble methods in machine learning combine the insights obtained from multiple learning models to facilitate accurate and improved decisions.
- ↳ There are two main types of ensemble learning:
  1. Bagging (Bootstrap Aggregating)
  2. Boosting.

## Ensemble learning Techniques

### Basic Ensemble Techniques :

- Max Voting
- Averaging
- Weighted average

### Advanced Ensemble Techniques :

- Stacking
- Blending
- Bagging
- Boosting

### Algorithms Based on Bagging :

- Bagging meta-estimator
- Random forest

### Algorithms Based on Boosting :

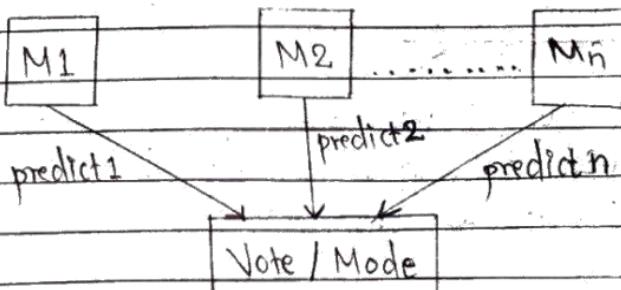
- AdaBoost

- GBM
- xGBM
- Light GBM
- Cat Boost

## 2. Basic Ensemble Techniques

### 1. Max Voting

- 6 The max voting method is generally for classification problems. In this technique, multiple models are used to make predictions for each data point.
- 6 The predictions by each model are considered as a 'vote'. The predictions which we get from majority of the models are used as the final prediction.

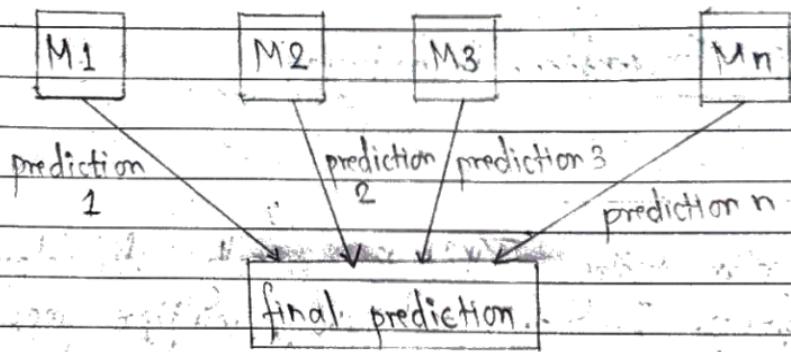


### 2. Averaging and Weighted Average Voting

- 6 Take an average of predictions from all the models

and use it to make the final prediction.

- 6 Averaging can be used for making predictions in regression problems or while calculating probabilities for classification problems.



- 6 Weighted means we have assigned the weight to models.

### Working Practically

```
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns
```

### Voting Classifier:

```
# Create a dataset using make_moons  
from sklearn.datasets import make_moons
```

```
# Get data for dataset
```

```
x, y = make_moons(n_samples=1000, noise=0.2)
```

```
# Here x is independent and y is dependent data
```

```
# Create a dataframe using x and y for tabular  
# visualization
```

```
df = {"x1": x[:, 0], "x2": x[:, 1], "y": y}
```

```
dataset = pd.DataFrame(df)
```

```
dataset.head(3)
```

```
# Visualize data
```

```
sns.scatterplot(x="x1", y="x2", data=dataset, hue=  
"y")
```

```
plt.show()
```

```
# Train test split
```

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,  
test_size=0.2, random_state=42)
```

```
# Train different models
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
from sklearn.svm import SVC
```

```
from sklearn.naive_bayes import GaussianNB
```

```
# Decision Tree Classifier
```

```
dtc = DecisionTreeClassifier()
```

```
dtc.fit(X_train, y_train)
```

```
dtc.score(X_train, y_train), dtc.score(X_test, y_test)
```

```
# SVC.
```

```
svc = SVC()
```

```
svc.fit(x-train, y-train)
```

```
svc.score(x-train, y-train), svc.score(x-test, y-test)
```

```
# Use Ensemble learning (Voting Classifier)  
from sklearn.ensemble import VotingClassifier
```

```
# Create a list of tuples with classifier names and  
# instances (estimator)
```

```
I = [("dtc", DecisionTreeClassifier()),  
     ("svc", SVC()),  
     ("gnb", GaussianNB())]
```

```
# Create the Voting Classifier model.
```

```
vc = VotingClassifier(I)
```

```
# You can also use weightage in terms of list:
```

```
# weights = [1, 2, 3]
```

```
vc.fit(x-train, y-train)
```

```
- vc.score(x-train, y-train), vc.score(x-test, y-test)
```

```
# Lets create a dataframe for clear visualization  
predicted = {
```

```
"DecisionTreeClassifier": dtc.predict(x-test),
```

```
"SVC": svc.predict(x-test),
```

```
"GaussianNB": gnb.predict(x-test),
```

```
"Voting Classifier": vc.predict(x-test),
```

```
}
```

```
# Create a dataframe
```

```
pd.DataFrame(predict)
```

```
# The most repeated result (mode) is final result
```

# by Voting Classifier

Voting Regressor:

```
dataset = pd.read_csv ("Maths.csv")
```

```
dataset.head(3)
```

# Check for null values

```
dataset.isnull().sum()
```

# Remove all the null values instead of filling (big  
# dataset)

```
dataset.dropna (inplace=True)
```

# Select dependent and independent data

```
x = dataset[["Number1"]]
```

```
y = dataset[["Total Sum"]]
```

# Train test split

```
from sklearn.model_selection import train_test_split
```

```
x_train, x_test, y_train, y_test = train_test_split (x,  
y, test_size=0.2, random_state=42)
```

# Train different models

```
from sklearn.linear_model import LinearRegression
```

```
from sklearn.tree import DecisionTreeRegressor
```

```
from sklearn.svm import SVR
```

# Linear Regression

```
lr = LinearRegression()
```

lr.fit(x-train, y-train)

lr.score(x-train, y-train), lr.score(x-test, y-test)

### # Decision Tree Regressor

dtr = DecisionTreeRegressor()

dtr.fit(x-train, y-train)

dtr.score(x-train, y-train), dtr.score(x-test, y-test)

### # SVR

svr = SVR()

svr.fit(x-train, y-train)

svr.score(x-train, y-train), svr.score(x-test, y-test)

### # Use Ensemble Learning (Voting Regressor)

from sklearn.ensemble import VotingRegressor

# Create a list of tuples with classifier names and

# instances (estimators)

l = [("lr", LinearRegression()), ("dtr",  
DecisionTreeRegressor()), ("svr", SVR())]

vr = VotingRegressor(l, weights=[2, 3, 4])

vr.fit(x-train, y-train)

vr.score(x-train, y-train), vr.score(x-test, y-test)

# Now lets see how Voting Regressor works by  
# comparison on dataset.

df = {

"LinearRegression": lr.predict(x-test),

"DecisionTreeRegressor": dtr.predict(x-test),

"SVR": svr.predict(x-test)

"Voting Regressor": vr.predict(x-test),

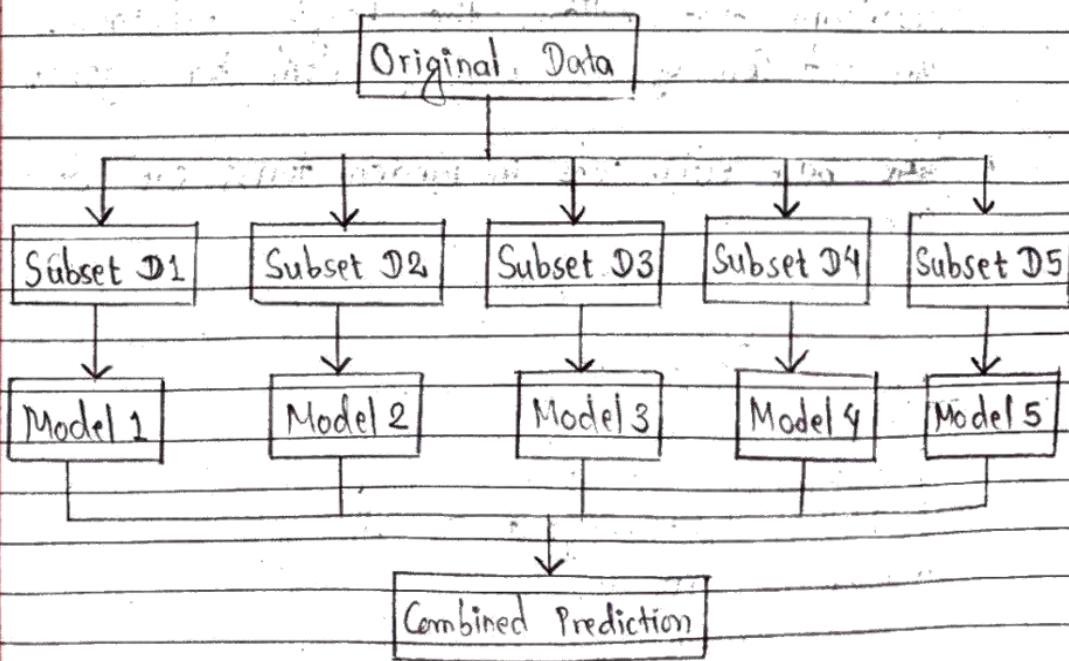
# Create a dataframe

pd.DataFrame(df).head(3)

# Voting Regressor gives average of other models

### 3. Bagging

- Bagging (or Bootstrap Aggregating) technique uses these subsets (bags) to get a fair idea of the distribution (complete set).
- The size of subsets created for bagging may be less than the original set.



↳ Bagging algorithms are:

1. Bagging meta-estimator
2. Random forest

## 1. Bagging Meta-Estimator

- ↳ Bagging meta-estimator is an ensembling algorithm that can be used for both classification (Bagging Classifier) and regression (Bagging Regressor) problems.
- ↳ It follows the typical bagging technique to make predictions.

## 2. Random Forest

- ↳ Random forest is another ensemble machine learning algorithm that follows the bagging technique. It is an extension of the bagging estimator algorithm.
- ↳ The base estimators in random forest are decision trees.

## Working Practically

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

## Classification Analysis:

```
# Create a dataset using make_moons  
from sklearn.datasets import make_moons
```

```
# Get data for dataset
```

```
x, y = make_moons(n_samples=1000, noise=0.2)
```

```
# Here x is independent and y is dependent data
```

```
# Create a dataframe using x and y for tabular  
# visualization
```

```
df = {"x1": [:, 0], "x2": [:, 1], "y": y}
```

```
dataset = pd.DataFrame(df)
```

```
dataset.head(8)
```

```
# Visualize data
```

```
sns.scatterplot(x="x1", y="x2", data=dataset,  
hue="y")  
plt.show()
```

```
# Train test split
```

```
from sklearn.model_selection import train_test_split
```

```
x_train, x_test, y_train, y_test = train_test_split(x,  
y, test_size=0.2, random_state=42)
```

```
# Classification analysis using Bagging Meta-Estimator  
# technique
```

```
from sklearn.ensemble import BaggingClassifier
```

```
# We need an estimator for BaggingClassifier
```

```
from sklearn.svm import SVC
```

# Check for the score of our estimator

```
SVC = SVC()
```

```
SVC.fit(x-train, y-train)
```

```
SVC.score(x-train, y-train), SVC.score(x-test, y-test)
```

# Create the Bagging Classifier model

```
bc = BaggingClassifier( estimator=SVC(), n_estimators=50 )
```

```
bc.fit(x-train, y-train)
```

```
bc.score(x-train, y-train), bc.score(x-test, y-test)
```

# Classification analysis using Random Forest technique  
from sklearn.ensemble import RandomForestClassifier

```
rfc = RandomForestClassifier(n_estimators=50)
```

```
rfc.fit(x-train, y-train)
```

```
rfc.score(x-train, y-train), rfc.score(x-test, y-test)
```

## Regression Analysis:

```
dataset = pd.read_csv("Maths.csv")
```

```
dataset.head(3)
```

# Check for null values

```
dataset.isnull().sum()
```

# Remove all the null values

```
dataset.dropna(inplace=True)
```

```
# Select dependent and independent data
```

```
x = dataset[["Number 3"]]
```

```
y = dataset[["Total Sum"]]
```

```
# Train test split
```

```
from sklearn.model_selection import train_test_split
```

```
x_train, x_test, y_train, y_test = train_test_split(x,  
y, test_size=0.2, random_state=42)
```

```
# Regression analysis using Bagging Meta-Estimator
```

```
# technique
```

```
from sklearn.ensemble import BaggingRegressor
```

```
# We need an estimator for Bagging Regressor
```

```
from sklearn.linear_model import LinearRegression
```

```
# Check for the score of our estimator.
```

```
lr = LinearRegression()
```

```
lr.fit(x_train, y_train)
```

```
lr.score(x_train, y_train), lr.score(x-test, y-test)
```

```
# Create the BaggingRegressor model
```

```
br = BaggingRegressor(estimator=LinearRegression(),  
n_estimators=50)
```

```
br.fit(x_train, y_train)
```

```
br.score(x_train, y_train), br.score(x-test, y-test)
```

```
# Regression analysis using Random Forest technique
```

```
from sklearn.ensemble import RandomForestRegressor
```

rfr = Random Forest Regressor (n\_estimators = 50)  
rfr.fit (x-train, y-train)

rfr.score (x-train, y-train), rfr.score (x-test, y-test)