# ensamble-learning

December 10, 2023

```python
[104]: import numpy as np
       import pandas as pd
       import matplotlib.pyplot as plt
       from sklearn.ensemble import RandomForestClassifier, VotingClassifier,
        ↪BaggingClassifier, ExtraTreesClassifier
       from sklearn.ensemble import BaggingRegressor, RandomForestRegressor,
        ↪ExtraTreesRegressor
       from sklearn.svm import SVC
       from sklearn.linear_model import LogisticRegression, LinearRegression
       from sklearn.model_selection import GridSearchCV
       from sklearn.preprocessing import StandardScaler
```

```python
[105]: df = pd.read_csv("/kaggle/input/diabetes-dataset/diabetes.csv")
       df.head()
```

```
[105]:    Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  \
       0            6      148             72             35        0  33.6
       1            1       85             66             29        0  26.6
       2            8      183             64              0        0  23.3
       3            1       89             66             23       94  28.1
       4            0      137             40             35      168  43.1

          DiabetesPedigreeFunction  Age  Outcome
       0                     0.627   50        1
       1                     0.351   31        0
       2                     0.672   32        1
       3                     0.167   21        0
       4                     2.288   33        1
```

```python
[106]: df.shape
```

```
[106]: (768, 9)
```

```python
[107]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
```

```
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Pregnancies               768 non-null    int64
 1   Glucose                   768 non-null    int64
 2   BloodPressure             768 non-null    int64
 3   SkinThickness             768 non-null    int64
 4   Insulin                   768 non-null    int64
 5   BMI                       768 non-null    float64
 6   DiabetesPedigreeFunction  768 non-null    float64
 7   Age                       768 non-null    int64
 8   Outcome                   768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

[108]: `df.isnull().sum()`

[108]:
```
Pregnancies                 0
Glucose                     0
BloodPressure               0
SkinThickness               0
Insulin                     0
BMI                         0
DiabetesPedigreeFunction    0
Age                         0
Outcome                     0
dtype: int64
```

[109]: `df.describe()`

[109]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin |
|---|---|---|---|---|---|
| count | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 |
| mean | 3.845052 | 120.894531 | 69.105469 | 20.536458 | 79.799479 |
| std | 3.369578 | 31.972618 | 19.355807 | 15.952218 | 115.244002 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 1.000000 | 99.000000 | 62.000000 | 0.000000 | 0.000000 |
| 50% | 3.000000 | 117.000000 | 72.000000 | 23.000000 | 30.500000 |
| 75% | 6.000000 | 140.250000 | 80.000000 | 32.000000 | 127.250000 |
| max | 17.000000 | 199.000000 | 122.000000 | 99.000000 | 846.000000 |

| | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|
| count | 768.000000 | 768.000000 | 768.000000 | 768.000000 |
| mean | 31.992578 | 0.471876 | 33.240885 | 0.348958 |
| std | 7.884160 | 0.331329 | 11.760232 | 0.476951 |
| min | 0.000000 | 0.078000 | 21.000000 | 0.000000 |
| 25% | 27.300000 | 0.243750 | 24.000000 | 0.000000 |
| 50% | 32.000000 | 0.372500 | 29.000000 | 0.000000 |
| 75% | 36.600000 | 0.626250 | 41.000000 | 1.000000 |

```
max        67.100000                    2.420000      81.000000      1.000000
```

```
[110]: categorical_val = []
       continous_val = []
       for column in df.columns:
           if len(df[column].unique()) <= 10:
               categorical_val.append(column)
           else:
               continous_val.append(column)
```

```
[111]: df.columns
```

```
[111]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
              'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
             dtype='object')
```

```
[112]: feature_columns = [
           'Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness',
           'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age'
       ]

       for column in feature_columns:
           print(f"column, {column} ==> Missing zeros : {len(df.loc[df[column] ==␣
        ↪0])}")
```

```
column, Pregnancies ==> Missing zeros : 111
column, Glucose ==> Missing zeros : 5
column, BloodPressure ==> Missing zeros : 35
column, SkinThickness ==> Missing zeros : 227
column, Insulin ==> Missing zeros : 374
column, BMI ==> Missing zeros : 11
column, DiabetesPedigreeFunction ==> Missing zeros : 0
column, Age ==> Missing zeros : 0
```

```
[113]: from sklearn.impute import SimpleImputer
       fill_values = SimpleImputer(missing_values=0, strategy="mean", copy=False)
       df[feature_columns] = fill_values.fit_transform(df[feature_columns])

       for column in feature_columns:
           print(f"column,{column} ==> Missing zeros : {len(df.loc[df[column] == 0])}")
```

```
column,Pregnancies ==> Missing zeros : 0
column,Glucose ==> Missing zeros : 0
column,BloodPressure ==> Missing zeros : 0
column,SkinThickness ==> Missing zeros : 0
column,Insulin ==> Missing zeros : 0
column,BMI ==> Missing zeros : 0
```

```
column,DiabetesPedigreeFunction ==> Missing zeros : 0
column,Age ==> Missing zeros : 0
```

[114]:
```python
X = df[feature_columns]
y = df.Outcome
```

[115]:
```python
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
 ↪random_state=42)
```

[116]:
```python
from sklearn.metrics import confusion_matrix, accuracy_score,
 ↪classification_report


def evaluate(model, X_train, X_test, y_train, y_test):
    y_test_pred = model.predict(X_test)
    y_train_pred = model.predict(X_train)

    print("TRAINIG RESULTS: \n===============================")
    clf_report = pd.DataFrame(classification_report(y_train, y_train_pred,
 ↪output_dict=True))
    print(f"CONFUSION MATRIX:\n{confusion_matrix(y_train, y_train_pred)}")
    print(f"ACCURACY SCORE:\n{accuracy_score(y_train, y_train_pred):.4f}")
    print(f"CLASSIFICATION REPORT:\n{clf_report}")

    print("TESTING RESULTS: \n===============================")
    clf_report = pd.DataFrame(classification_report(y_test, y_test_pred,
 ↪output_dict=True))
    print(f"CONFUSION MATRIX:\n{confusion_matrix(y_test, y_test_pred)}")
    print(f"ACCURACY SCORE:\n{accuracy_score(y_test, y_test_pred):.4f}")
    print(f"CLASSIFICATION REPORT:\n{clf_report}")
```

Bagging Algorithms

[117]:
```python
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
rf_clf = RandomForestClassifier(random_state=42, n_estimators=1000)
rf_clf.fit(X_train, y_train)
evaluate(rf_clf, X_train, X_test, y_train, y_test)
```

```
TRAINIG RESULTS:
===============================
CONFUSION MATRIX:
[[349   0]
 [  0 188]]
ACCURACY SCORE:
1.0000
```

```
CLASSIFICATION REPORT:
                 0       1   accuracy   macro avg   weighted avg
precision      1.0     1.0        1.0         1.0            1.0
recall         1.0     1.0        1.0         1.0            1.0
f1-score       1.0     1.0        1.0         1.0            1.0
support      349.0   188.0        1.0       537.0          537.0
TESTING RESULTS:
==============================
CONFUSION MATRIX:
[[123  28]
 [ 29  51]]
ACCURACY SCORE:
0.7532
CLASSIFICATION REPORT:
                    0          1   accuracy   macro avg   weighted avg
precision    0.809211   0.645570   0.753247    0.727390       0.752538
recall       0.814570   0.637500   0.753247    0.726035       0.753247
f1-score     0.811881   0.641509   0.753247    0.726695       0.752878
support    151.000000  80.000000   0.753247  231.000000     231.000000
```

```python
[118]: from sklearn.ensemble import BaggingClassifier
       from sklearn.ensemble import RandomForestClassifier
       from sklearn.metrics import accuracy_score
       from sklearn.model_selection import train_test_split
       from sklearn.tree import DecisionTreeClassifier

       # Create a BaggingClassifier
       base_classifier = DecisionTreeClassifier()
       bagging_clf = BaggingClassifier(base_classifier, n_estimators=10,␣
        ↪random_state=42)

       # Fit the BaggingClassifier
       bagging_clf.fit(X_train, y_train)

       # Calculate and store accuracy scores for Bagging Classifier
       bagging_scores = {
           'Train': accuracy_score(y_train, bagging_clf.predict(X_train)),
           'Test': accuracy_score(y_test, bagging_clf.predict(X_test)),
       }

       # Calculate and store accuracy scores for Bagging Classifier
       scores = {
           'Bagging Classifier': {
               'Train': accuracy_score(y_train, bagging_clf.predict(X_train)),
               'Test': accuracy_score(y_test, bagging_clf.predict(X_test)),
           },
       }
```

```python
# Calculate and store accuracy scores for Random Forest
scores['Random Forest'] = {
    'Train': accuracy_score(y_train, rf_clf.predict(X_train)),
    'Test': accuracy_score(y_test, rf_clf.predict(X_test)),
}
```

Boosting Algorithms

```
[119]: from sklearn.ensemble import AdaBoostClassifier

       ada_boost_clf = AdaBoostClassifier(n_estimators=30)
       ada_boost_clf.fit(X_train, y_train)
       evaluate(ada_boost_clf, X_train, X_test, y_train, y_test)
```

```
TRAINIG RESULTS:
===============================
CONFUSION MATRIX:
[[310  39]
 [ 51 137]]
ACCURACY SCORE:
0.8324
CLASSIFICATION REPORT:
                         0            1   accuracy    macro avg   weighted avg
precision         0.858726     0.778409   0.832402     0.818567       0.830607
recall            0.888252     0.728723   0.832402     0.808488       0.832402
f1-score          0.873239     0.752747   0.832402     0.812993       0.831056
support         349.000000   188.000000   0.832402   537.000000     537.000000
TESTING RESULTS:
===============================
CONFUSION MATRIX:
[[123  28]
 [ 27  53]]
ACCURACY SCORE:
0.7619
CLASSIFICATION REPORT:
                         0            1   accuracy    macro avg   weighted avg
precision         0.820000     0.654321   0.761905     0.737160       0.762622
recall            0.814570     0.662500   0.761905     0.738535       0.761905
f1-score          0.817276     0.658385   0.761905     0.737830       0.762249
support         151.000000    80.000000   0.761905   231.000000     231.000000
```

```
[120]: scores['AdaBoost'] = {
           'Train': accuracy_score(y_train, ada_boost_clf.predict(X_train)),
           'Test': accuracy_score(y_test, ada_boost_clf.predict(X_test)),
       }
```

```
[121]: from sklearn.ensemble import GradientBoostingClassifier

       grad_boost_clf = GradientBoostingClassifier(n_estimators=100, random_state=42)
       grad_boost_clf.fit(X_train, y_train)
       evaluate(grad_boost_clf, X_train, X_test, y_train, y_test)
```

```
TRAINIG RESULTS:
===============================
CONFUSION MATRIX:
[[342    7]
 [ 19 169]]
ACCURACY SCORE:
0.9516
CLASSIFICATION REPORT:
                     0           1  accuracy    macro avg  weighted avg
precision     0.947368    0.960227  0.951583     0.953798      0.951870
recall        0.979943    0.898936  0.951583     0.939439      0.951583
f1-score      0.963380    0.928571  0.951583     0.945976      0.951194
support     349.000000  188.000000  0.951583   537.000000    537.000000
TESTING RESULTS:
===============================
CONFUSION MATRIX:
[[116  35]
 [ 26  54]]
ACCURACY SCORE:
0.7359
CLASSIFICATION REPORT:
                     0           1  accuracy    macro avg  weighted avg
precision     0.816901    0.606742  0.735931     0.711821      0.744119
recall        0.768212    0.675000  0.735931     0.721606      0.735931
f1-score      0.791809    0.639053  0.735931     0.715431      0.738906
support     151.000000   80.000000  0.735931   231.000000    231.000000
```

```
[122]: scores['Gradient Boosting'] = {
            'Train': accuracy_score(y_train, grad_boost_clf.predict(X_train)),
            'Test': accuracy_score(y_test, grad_boost_clf.predict(X_test)),
       }
```

```
[123]: from sklearn.ensemble import VotingClassifier
       from sklearn.linear_model import LogisticRegression
       from sklearn.tree import DecisionTreeClassifier
       from sklearn.svm import SVC
       from sklearn.model_selection import train_test_split
       from sklearn.metrics import accuracy_score

       # Define classifiers
       log_reg = LogisticRegression(solver='liblinear')
```

```
tree = DecisionTreeClassifier()
svm_clf = SVC(gamma='scale')

estimators = [('Logistic', log_reg), ('Tree', tree), ('SVM', svm_clf)]

voting = VotingClassifier(estimators=estimators)
voting.fit(X_train, y_train)
evaluate(voting, X_train, X_test, y_train, y_test)
```

```
TRAINIG RESULTS:
===============================
CONFUSION MATRIX:
[[327  22]
 [ 82 106]]
ACCURACY SCORE:
0.8063
CLASSIFICATION REPORT:
                    0          1  accuracy   macro avg  weighted avg
precision    0.799511   0.828125  0.806331    0.813818      0.809529
recall       0.936963   0.563830  0.806331    0.750396      0.806331
f1-score     0.862797   0.670886  0.806331    0.766841      0.795610
support    349.000000 188.000000  0.806331  537.000000    537.000000
TESTING RESULTS:
===============================
CONFUSION MATRIX:
[[131  20]
 [ 36  44]]
ACCURACY SCORE:
0.7576
CLASSIFICATION REPORT:
                    0          1  accuracy   macro avg  weighted avg
precision    0.784431   0.687500  0.757576    0.735966      0.750862
recall       0.867550   0.550000  0.757576    0.708775      0.757576
f1-score     0.823899   0.611111  0.757576    0.717505      0.750206
support    151.000000  80.000000  0.757576  231.000000    231.000000
```

[124]:
```
scores['Voting'] = {
        'Train': accuracy_score(y_train, voting.predict(X_train)),
        'Test': accuracy_score(y_test, voting.predict(X_test)),
    }
```

[125]:
```
scores_df = pd.DataFrame(scores)
scores_df.plot(kind='barh', figsize=(15, 8))
```

[125]: <Axes: >