

# Performance analysis of machine learning models to predict diabetes

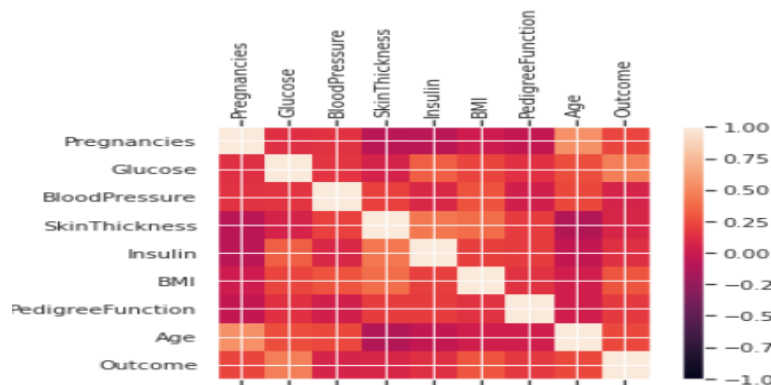
**Introduction:** Diabetes is one of the world's most common and well-known life-threatening diseases[1]. Diabetes detected early can help to decrease complications and the risk of a wide range of health issues. The purpose of this project is to predict diabetes based on various parameters using machine learning models such as logistic regression, decision tree, naïve bayes, adaboost and to evaluate their accuracy to determine which model performs best in predicting diabetes.

## Method:

### i) Dataset and preprocessing:

Pima Indian diabetes dataset has been used to classify diabetic and non-diabetic patient based on machine learning models. This dataset has a total of 768 training cases. Each training instance includes eight features and a class variable that serves as the training instance's label. Pregnancies, glucose, blood pressure, BMI, skin thickness, insulin, age, and diabetes pedigree function are the independent variables. The feature that will be predicted is the outcome, where 0 indicates that the patient does not have diabetes and 1 indicates that the patient does have diabetes. There is no null value in the dataset.

### ii) Correlation:



A significant correlation can be observed between pregnancy and age in the correlation matrix

iii) Normalization and splitting dataset: The data is normalized from 0 to 1. Following normalization, the dataset is divided into two parts: training and testing. 80% of the data is utilized to train the model, with the remaining 20% used for testing.

### iv) Applying machine learning models:

Machine learning models that are applied to predict diabetes are follows:

- a) Logistic regression: Logistic regression is a type of supervised machine learning. In this experiment, multiple learning rates were used to assess test accuracy. With a learning rate of 2, the best accuracy is 89.87%.

- b) Decision tree: Decision tree is a non-parametric classifier. The accuracy of the tree is measured at various depths. When depth=3, the decision tree has the highest accuracy of 86.08%.
- c) Naïve Bayes: The Naive Bayes model is based on the Bayes theorem. In Naïve Bayes data are categorized by computing the probability of the output depending on the attributes. Naive Bayes has an accuracy of 84.81% in predicting diabetes from a dataset.
- d) Adaboost: Adaptive Boosting, often known as "AdaBoost," focuses on classification problems and seeks to convert a group of weak classifiers into a powerful one[2].For predicting diabetes, accuracy of adaboost is 91.13% for learning rate 0.1

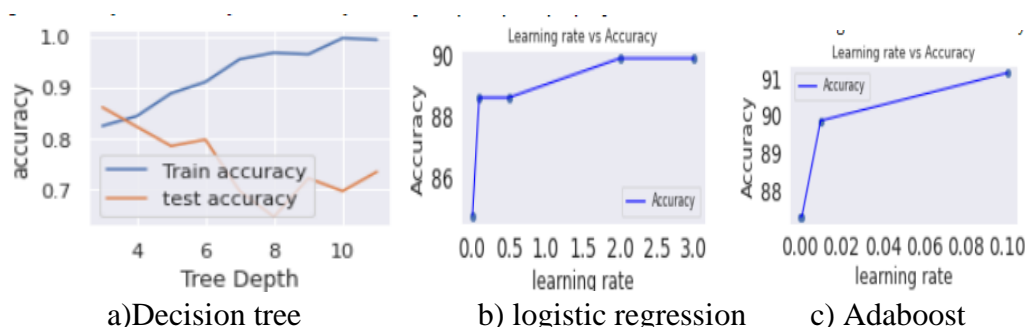


Figure 2: HyperTune decision tree, logistic regression, Adaboost for different test accuracy

**Result:** Adaboost has the highest accuracy of 91.13% among four machine learning models. To predict diabetes from the dataset, logistic regression, naive bayes, and decision tree all perform well. Precision, recall and f1 score of adaboost are 0.91,0.88, 0.89 respectively which are higher than remaining three models. The feature importance of the attributes was analyzed in the AdaBoost model, and it was discovered that glucose has the most significant influence in the model for prediction. Comparison among various models is shown below

	Model Name	Test Accuracy	Precision	Recall	f1-score
0	Logistic Regression	89.873418	0.890678	0.841481	0.860177
1	Naive Bayes	84.810127	0.822293	0.835185	0.828012
2	Decision tree	86.076000	0.836054	0.865926	0.846602
3	Adaboost	91.139241	0.910686	0.881481	0.894000

Table 1: Comparison among four models

**Conclusion:** The results show that there is no significant difference in accuracy between the four models. On the dataset, all models perform well, with adaboost performing the best. Because of the nature of the data in the dataset, it is not possible to distinguish between type 1 and type 2 databases. The goal for the future is to explore various features and determine the type of diabetes. Furthermore, these four models can be integrated with other machine learning models to improve prediction accuracy.

## Bibliography:

- 1) Butt UM, Letchmunan S, Ali M, Hassan FH, Baqir A, Sherazi HHR. Machine Learning Based Diabetes Classification and Prediction for Healthcare Applications. *J Healthc Eng.* 2021 Sep 29;2021:9930985. doi: 10.1155/2021/9930985. PMID: 34631003; PMCID: PMC8500744.
- 2)N. H. Taz, A. Islam and I. Mahmud, "A Comparative Analysis of Ensemble Based Machine Learning Techniques for Diabetes Identification," *2021 2nd International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST)*, 2021, pp. 1-6, doi: 10.1109/ICREST51555.2021.9331036.



# Performance analysis of machine learning models to predict diabetes

## import necessary libraries

```
In [ ]: # This Python 3 environment comes with many helpful analytics libraries installed
# It is derived by the kaggle/python docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
# for ML model
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import cross_val_score, cross_validate
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.ensemble import AdaBoostClassifier

sns.set(style='white')

# Input data files are available in the read-only "../input/" directory
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# You can write up to 200MB to the current directory (/kaggle/working/) that gets preserved as output when you create a version using "Save & Run All"
# You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session
```

## Load dataset

```
In [ ]: data = pd.read_csv("../input/pima-indians-diabetes-database/diabetes.csv")
data.head()
```

```
Out [ ]: 
```

	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	0	88	67	0	0	22.3	0.672	32	1
3	1	80	44	23	94	28.1	0.167	21	0
4	0	137	40	35	160	43.1	2.288	33	1

## Find Correlation

```
In [ ]: correlations = data.corr(method = 'pearson') # Correlations between all pairs of attributes
# Print the correlations
print(correlations)
# Print the correlation matrix
print(correlations)
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073635	0.127683	-0.033323	0.543441	0.231898
Glucose	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337	0.263914	0.466881
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.089933	0.281905	0.041265	0.239526	0.095068
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.426769	0.366979	0.189929	-0.119970	0.074762
Insulin	-0.073635	0.331357	0.089933	0.426769	1.000000	0.119769	0.185071	0.042163	0.136548
BMI	0.127683	0.221071	0.281905	0.366979	0.119769	1.000000	0.140547	0.036242	0.292696
DiabetesPedigreeFunction	-0.033323	0.137337	0.041265	0.189929	0.185071	0.140547	1.000000	0.035661	0.178844
Age	0.543441	0.263914	0.239526	-0.119970	-0.042163	0.036242	0.035661	1.000000	0.238256
Outcome	0.231898	0.466881	0.095068	0.074762	0.136548	0.292696	0.178844	0.238256	1.000000

```
In [ ]: import numpy as np

# plot correlation matrix
cm = plt.figure(figsize=(10,10))
cm = data.corr()
ax = fig.add_subplot(111)
ax = cm
ax.set_xticklabels(cm.columns)
ax.set_yticklabels(cm.columns)
ax.set_xlabel('')
ax.set_ylabel('')
ax.set_xticks(cm.columns)
ax.set_yticks(cm.columns)
ax.set_xticklabels(cm.columns, rotation=90) # rotate x-tick labels by 90 degrees
ax.set_yticklabels(cm.columns)
plt.show()
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Split training and test set

```
In [ ]: # 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 = 89)
```

## Naive Bayes classifier

```
In [ ]: # Naive Bayes
from sklearn.naive_bayes import GaussianNB
nb = GaussianNB()
nb.fit(x_train, y_train)
print("Accuracy of Naive Bayes Algorithm: ", nb.score(x_test, y_test))
ac_nb = nb.score(x_test, y_test)
print("Accuracy of Naive Bayes Algorithm: ", ac_nb)
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Decision tree classifier

```
In [ ]: from sklearn.tree import DecisionTreeClassifier
def evaluation(y):
    m = DecisionTreeClassifier(max_depth=3)
    m.fit(x_train, y_train)
    y_pred = m.predict(x_test)
    accuracy = round((np.sum(y_pred == y_test)) / len(y_test), 2)
    precision = round((np.sum(y_pred == y_test) / len(y_pred)), 2)
    recall = round((np.sum(y_pred == y_test) / len(y_test)), 2)
    f_score = round((2 * precision * recall) / (precision + recall), 2)
    return accuracy, precision, recall, f_score
depth_list = range(1, 12)
log = np.zeros((2, len(depth_list)))

def validation_graph(log, list_x_label):
    plt.figure(figsize=(10,10))
    plt.xlabel('accuracy')
    plt.ylabel('precision')
    sns.lineplot(log, list_x_label, x=list_x_label, y=list_x_label)
    test_ac = []
    for i, depth in enumerate(depth_list):
        model = DecisionTreeClassifier(max_depth=depth)
        model.fit(x_train, y_train)
        predict_train = model.predict(x_train)
        predict_test = model.predict(x_test)
        result_train = evaluation(predict_train, y_train)
        result_test = evaluation(predict_test, y_test)
        test_ac.append(result_test)
        log[0][i], log[1][i] = result_train[0], result_test[0]

    change_in_accuracy = []
    validation_graph(log, depth_list, "Tree Depth")
    print(test_ac)
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Classification report and plotting decision tree

```
In [ ]: # we pick the parameter that perform the best from the above output to predict the test value
def print_report(result):
    print("accuracy:", result[0], "\nprecision:", result[1], "\nrecall:", result[2], "\nf-score:", result[3], "\nresult:", result[4])

model = DecisionTreeClassifier(max_depth=3)
model.fit(x_train, y_train)
predict_train = model.predict(x_train)
predict_test = model.predict(x_test)
result_train = evaluation(predict_train, y_train)
result_test = evaluation(predict_test, y_test)
result_output(result)
ac_tree = result_test[0]
precision = result_test[1]
recall = result_test[2]
f_score = result_test[3]
print(classification_report(y_test, predict_test, target_names=["Train accuracy", "Test accuracy"]))
sns.lineplot(log, list_x_label, x=list_x_label, y=list_x_label)
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Feature importance of decision tree

```
In [ ]: diabetes_features = (x for i, x in enumerate(x1.columns) if i != 8)
print("Feature Importance: ", diabetes_features)

Decision Tree Feature Importances:
0.0 0.58157958 0.0 0.10291882 0.02815081
0.0 0.28735081
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Logistic regression

### Hypertune logistic regression for different learning rate

```
In [ ]: from sklearn.linear_model import LogisticRegression
learning_rate = [0.01, 0.1, 0.5, 1, 3]
ac = []
for i in learning_rate:
    lr = LogisticRegression(C=1)
    lr.fit(x_train, y_train)
    acc_lr = lr.score(x_test, y_test)
    ac.append(acc_lr)

for learning_rate = 0.01 test accuracy: 84.8391265827847 %
for learning_rate = 0.1 test accuracy: 88.4875849370885 %
for learning_rate = 0.5 test accuracy: 88.4875849370885 %
for learning_rate = 1 test accuracy: 89.8734777215899 %
for learning_rate = 3 test accuracy: 89.8734777215899 %
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Feature importance for logistic regression

```
In [ ]: lr = LogisticRegression(C=2)
lr.fit(x_train, y_train)
acc_lr = lr.score(x_test, y_test)
predictions = lr.predict(x_test)

feature_names = ["Pregnancies", "Glucose", "BloodPressure", "SkinThickness", "Insulin", "BMI", "DiabetesPedigreeFunction", "Age"]
feature_importance = pd.DataFrame(feature_names, columns = ["feature"])

importance = lr.coef_
importance.sort_values(by="importance", ascending=False)
feature_importance["importance"] = importance
feature_importance.sort_values(by="importance", ascending=False)
print(feature_importance)
ax = feature_importance.plot.barh(x="feature", y="importance")
plt.show()
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## AdaBoost Classifier

### hypertune AdaBoost for different learning rate

```
In [ ]: learning_rate = [0.001, 0.001, 0.01, 0.1]
ac2 = []
for i in learning_rate:
    ada = AdaBoostClassifier(learning_rate=i, n_estimators=30)
    ada.fit(x_train, y_train)
    pred = ada.predict(x_test)
    ac_ada = ada.score(x_test, y_test)
    ac2.append(ac_ada)
    print("for learning_rate = ", i, "Accuracy: ", ac_ada)
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## feature importance of AdaBoost classifier

```
In [ ]: feature_names = ["Pregnancies", "Glucose", "BloodPressure", "SkinThickness", "Insulin", "BMI", "DiabetesPedigreeFunction", "Age"]
feature_importance = pd.DataFrame(feature_names, columns = ["feature"])

importance = ada.feature_importances_
importance.sort_values(by="importance", ascending=False)
feature_importance["importance"] = importance
feature_importance.sort_values(by="importance", ascending=False)
print(feature_importance)
ax = feature_importance.plot.barh(x="feature", y="importance")
plt.show()
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns

## Comparing Accuracy of all models

```
In [ ]: all_accuracy = [{"Logistic Regression", ac_log, precision_log, recall_log, f_score_log}, {"Naive Bayes", ac_nb, precision_nb, recall_nb, f_score_nb}, {"Decision tree", ac_tree, precision_t, accuracy_summary}, {"AdaBoost", ac_ada, precision_ada, recall_ada, f_score_ada}]

accuracy_summary = pd.DataFrame(data=all_accuracy, columns=["Model Name", "Test Accuracy", "Precision", "Recall", "f1-score"])

accuracy_summary
```

```
Out [ ]: 
```

	Model Name	Test Accuracy	Precision	Recall	f1-score
0	Logistic Regression	89.873477	0.890767	0.841461	0.866177
1	Naive Bayes	84.839127	0.822229	0.893186	0.859032
2	Decision tree	88.487585	0.880000	0.880000	0.880000
3	AdaBoost	91.139245	0.908697	0.881481	0.895000

## Plotting the accuracy

```
In [ ]: accomp.array([ac_log, ac_nb, ac_tree, ac_ada])
fig = plt.figure(figsize=(10,10))
plt.scatter(x=learning_rate, y=accuracy_summary, color=["purple", "green", "red", "blue"])
plt.xticks(rotation=45)
plt.xlabel("Accuracy")
```

```
Out [ ]: 
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	160.0	43.1	2.288	33	1
6	3	78.0	50.0	32.0	88.0	31.0	0.248	26	0
8	2	197.0	70.0	45.0	543.0	30.5	0.158	53	1
13	1	189.0	60.0	23.0	846.0	30.1	0.398	59	1

392 rows x 9 columns