

Classification (2)

Import Library dan load Dataset

```
In [1]: import pandas as pd
import numpy as np

In [2]: dataset = pd.read_csv(r"C:\SEMESTER 4\LAB\Bahan Modul 10\Diabetes.csv", header=0)

Out[2]:
```

	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	20.6	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	2.278	33	1
...									
763	10	101	70	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.345	47	1
767	1	93	70	31	0	30.4	0.193	23	0

768 rows x 9 columns

```
In [3]: dataset.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    float64
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

Modelling
```

```
In [4]: #split dataset in features and target variable
feature_cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age']
X = dataset[feature_cols].values
y = dataset['Outcome'].values

In [5]: #split dataset into training set and test set
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=1)
```

Menggunakan teknik pruning

Decision Tree biasa

```
In [6]: from sklearn.tree import DecisionTreeClassifier

# Create Decision Tree Classifier object
clf = DecisionTreeClassifier()

# Train Decision Tree Classifier
clf = clf.fit(X_train, y_train)

# Predict the response for test dataset
y_pred = clf.predict(X_test)

In [7]: from sklearn import metrics

# Model Accuracy, how often is the classifier correct?
print("Accuracy", metrics.accuracy_score(y_test, y_pred))

Accuracy: 0.70128861987013

In [8]: from matplotlib import pyplot as plt
from sklearn import tree

fig = plt.figure(figsize=(25,20))
tree.plot_tree(clf, filled=True)

Out[8]:
```

Decision Tree with pruning

```
In [9]: # Create Decision Tree Classifier
clf_p = DecisionTreeClassifier(criterion='entropy', max_depth=3)

# Train Decision Tree Classifier
clf_p = clf_p.fit(X_train, y_train)

# Predict the response for test dataset
y_pred_p = clf_p.predict(X_test)

# Model Accuracy, how often is the classifier correct?
print("Accuracy", metrics.accuracy_score(y_test, y_pred_p))

Accuracy: 0.778627708277086

In [10]: from matplotlib import pyplot as plt
from sklearn import tree

fig = plt.figure(figsize=(25,20))
tree.plot_tree(clf_p, filled=True)

Out[10]:
```

Menggunakan random forest tanpa pruning

```
In [11]: from sklearn.ensemble import RandomForestClassifier

# Create Decision Tree Classifier object
rf = RandomForestClassifier()

# Train Decision Tree Classifier
rf = rf.fit(X_train, y_train)

# Predict the response for test dataset
y_pred_rf = rf.predict(X_test)

# Model Accuracy, how often is the classifier correct?
print("Accuracy", metrics.accuracy_score(y_test, y_pred_rf))

Accuracy: 0.805149895194892
```

Menampilkan gambaran dari random forest tanpa pruning

```
In [12]: fig = plt.figure(figsize=(25,20))
tree.plot_tree(rf.estimators_[0], filled=True);
```

Random Forest menggunakan pruning

```
In [14]: from sklearn.ensemble import RandomForestClassifier

# Create Decision Tree Classifier object
rf_p = RandomForestClassifier(criterion='entropy', max_depth=3, max_leaf_nodes=5, class_weight=None)

# Train Decision Tree Classifier
rf_p = rf_p.fit(X_train, y_train)

# Predict the response for test dataset
y_pred_rf_p = rf_p.predict(X_test)

print("Accuracy", metrics.accuracy_score(y_test, y_pred_rf_p))

Accuracy: 0.7963367965367965
```

Menampilkan gambaran dari Random Forest dengan pruning

```
In [15]: fig = plt.figure(figsize=(25,20))
tree.plot_tree(rf_p.estimators_[0], filled=True);
```

Perhatikan apakah hasil akurasi dengan Random Forest sebelum dan sesudah menggunakan feature importance mengalami kenaikan atau penurunan?

Hasilnya akurasi nya naik

```
In [1]: !jupyter nbconvert --to html "C:\SEMESTER 4\LAB\Bahan Modul 10\Diabetes_10.ipynb" --output-dir=""
```