

Consider the following dataset: Iris Dataset The iris dataset is a simple and beginner-friendly dataset that contains information about the flower petal and sepal sizes. The dataset has 3 classes with 50 instances in each class, therefore, it contains 150 rows with only 4 columns. Requirement: Implement any appropriate machine-learning algorithm on the dataset to provide insights on the dataset.

1) we first import the necessary libraries that will help us to do analysis on our dataset

```
In [12]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report
```

2) we then load the iris data and convert it to data frame as it provides an easy-to-use interface for data manipulation, visualization, and analysis.

```
In [19]: from sklearn.datasets import load_iris
iris = load_iris()
df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
df['target'] = iris.target
print(iris)
```

```
{'data': array([[5.1, 3.5, 1.4, 0.2],
 [4.9, 3. , 1.4, 0.2],
 [4.7, 3.2, 1.3, 0.2],
 [4.6, 3.1, 1.5, 0.2],
 [5. , 3.6, 1.4, 0.2],
 [5.4, 3.9, 1.7, 0.4],
 [4.6, 3.4, 1.4, 0.3],
 [5. , 3.4, 1.5, 0.2],
 [4.4, 2.9, 1.4, 0.2],
 [4.9, 3.1, 1.5, 0.1],
 [5.4, 3.7, 1.5, 0.2],
 [4.8, 3.4, 1.6, 0.2],
 [4.8, 3. , 1.4, 0.1],
 [4.3, 3. , 1.1, 0.1],
 [5.8, 4. , 1.2, 0.2],
 [5.7, 4.4, 1.5, 0.4],
 [5.4, 3.9, 1.3, 0.4],
 [5.1, 3.5, 1.4, 0.3],
 [5.7, 3.8, 1.7, 0.3],
 [5.1, 3.8, 1.5, 0.3],
 [5.4, 3.4, 1.7, 0.2],
 [5.1, 3.7, 1.5, 0.4],
 [4.6, 3.6, 1. , 0.2],
 [5.1, 3.3, 1.7, 0.5],
 [4.8, 3.4, 1.9, 0.2],
 [5. , 3. , 1.6, 0.2],
 [5. , 3.4, 1.6, 0.4],
 [5.2, 3.5, 1.5, 0.2],
 [5.2, 3.4, 1.4, 0.2],
 [4.7, 3.2, 1.6, 0.2],
 [4.8, 3.1, 1.6, 0.2],
 [5.4, 3.4, 1.5, 0.4],
 [5.2, 4.1, 1.5, 0.1],
 [5.5, 4.2, 1.4, 0.2],
 [4.9, 3.1, 1.5, 0.2],
 [5. , 3.2, 1.2, 0.2],
 [5.5, 3.5, 1.3, 0.2],
 [4.9, 3.6, 1.4, 0.1],
 [4.4, 3. , 1.3, 0.2],
 [5.1, 3.4, 1.5, 0.2],
 [5. , 3.5, 1.3, 0.3],
 [4.5, 2.3, 1.3, 0.3],
 [4.4, 3.2, 1.3, 0.2],
 [5. , 3.5, 1.6, 0.6],
 [5.1, 3.8, 1.9, 0.4],
 [4.8, 3. , 1.4, 0.3],
 [5.1, 3.8, 1.6, 0.2],
 [4.6, 3.2, 1.4, 0.2],
 [5.3, 3.7, 1.5, 0.2],
 [5. , 3.3, 1.4, 0.2],
 [7. , 3.2, 4.7, 1.4],
 [6.4, 3.2, 4.5, 1.5],
 [6.9, 3.1, 4.9, 1.5],
 [5.5, 2.3, 4. , 1.3],
 [6.5, 2.8, 4.6, 1.5],
 [5.7, 2.8, 4.5, 1.3],
 [6.3, 3.3, 4.7, 1.6],
 [4.9, 2.4, 3.3, 1. ],
 [6.6, 2.9, 4.6, 1.3],
 [5.2, 2.7, 3.9, 1.4],
 [5. , 2. , 3.5, 1. ],
 [5.9, 3. , 4.2, 1.5],
 [6. , 2.2, 4. , 1. ],
 [6.1, 2.9, 4.7, 1.4],
```

[5.6, 2.9, 3.6, 1.3],
[6.7, 3.1, 4.4, 1.4],
[5.6, 3. , 4.5, 1.5],
[5.8, 2.7, 4.1, 1.],
[6.2, 2.2, 4.5, 1.5],
[5.6, 2.5, 3.9, 1.1],
[5.9, 3.2, 4.8, 1.8],
[6.1, 2.8, 4. , 1.3],
[6.3, 2.5, 4.9, 1.5],
[6.1, 2.8, 4.7, 1.2],
[6.4, 2.9, 4.3, 1.3],
[6.6, 3. , 4.4, 1.4],
[6.8, 2.8, 4.8, 1.4],
[6.7, 3. , 5. , 1.7],
[6. , 2.9, 4.5, 1.5],
[5.7, 2.6, 3.5, 1.],
[5.5, 2.4, 3.8, 1.1],
[5.5, 2.4, 3.7, 1.],
[5.8, 2.7, 3.9, 1.2],
[6. , 2.7, 5.1, 1.6],
[5.4, 3. , 4.5, 1.5],
[6. , 3.4, 4.5, 1.6],
[6.7, 3.1, 4.7, 1.5],
[6.3, 2.3, 4.4, 1.3],
[5.6, 3. , 4.1, 1.3],
[5.5, 2.5, 4. , 1.3],
[5.5, 2.6, 4.4, 1.2],
[6.1, 3. , 4.6, 1.4],
[5.8, 2.6, 4. , 1.2],
[5. , 2.3, 3.3, 1.],
[5.6, 2.7, 4.2, 1.3],
[5.7, 3. , 4.2, 1.2],
[5.7, 2.9, 4.2, 1.3],
[6.2, 2.9, 4.3, 1.3],
[5.1, 2.5, 3. , 1.1],
[5.7, 2.8, 4.1, 1.3],
[6.3, 3.3, 6. , 2.5],
[5.8, 2.7, 5.1, 1.9],
[7.1, 3. , 5.9, 2.1],
[6.3, 2.9, 5.6, 1.8],
[6.5, 3. , 5.8, 2.2],
[7.6, 3. , 6.6, 2.1],
[4.9, 2.5, 4.5, 1.7],
[7.3, 2.9, 6.3, 1.8],
[6.7, 2.5, 5.8, 1.8],
[7.2, 3.6, 6.1, 2.5],
[6.5, 3.2, 5.1, 2.],
[6.4, 2.7, 5.3, 1.9],
[6.8, 3. , 5.5, 2.1],
[5.7, 2.5, 5. , 2.],
[5.8, 2.8, 5.1, 2.4],
[6.4, 3.2, 5.3, 2.3],
[6.5, 3. , 5.5, 1.8],
[7.7, 3.8, 6.7, 2.2],
[7.7, 2.6, 6.9, 2.3],
[6. , 2.2, 5. , 1.5],
[6.9, 3.2, 5.7, 2.3],
[5.6, 2.8, 4.9, 2.],
[7.7, 2.8, 6.7, 2.],
[6.3, 2.7, 4.9, 1.8],
[6.7, 3.3, 5.7, 2.1],
[7.2, 3.2, 6. , 1.8],
[6.2, 2.8, 4.8, 1.8],
[6.1, 3. , 4.9, 1.8],

```
[6.4, 2.8, 5.6, 2.1],
[7.2, 3. , 5.8, 1.6],
[7.4, 2.8, 6.1, 1.9],
[7.9, 3.8, 6.4, 2. ],
[6.4, 2.8, 5.6, 2.2],
[6.3, 2.8, 5.1, 1.5],
[6.1, 2.6, 5.6, 1.4],
[7.7, 3. , 6.1, 2.3],
[6.3, 3.4, 5.6, 2.4],
[6.4, 3.1, 5.5, 1.8],
[6. , 3. , 4.8, 1.8],
[6.9, 3.1, 5.4, 2.1],
[6.7, 3.1, 5.6, 2.4],
[6.9, 3.1, 5.1, 2.3],
[5.8, 2.7, 5.1, 1.9],
[6.8, 3.2, 5.9, 2.3],
[6.7, 3.3, 5.7, 2.5],
[6.7, 3. , 5.2, 2.3],
[6.3, 2.5, 5. , 1.9],
[6.5, 3. , 5.2, 2. ],
[6.2, 3.4, 5.4, 2.3],
[5.9, 3. , 5.1, 1.8]]), 'target': array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]), 'frame': None, 'target_na
mes': array(['setosa', 'versicolor', 'virginica'], dtype='<U10'), 'DESCR': '.. _iris_dat
aset:\n\nIris plants dataset\n-----\n\n**Data Set Characteristics:**\n\n:
Number of Instances: 150 (50 in each of three classes)\n      :Number of Attributes: 4 nu
meric, predictive attributes and the class\n      :Attribute Information:\n          - sepal
length in cm\n          - sepal width in cm\n          - petal length in cm\n          - petal
width in cm\n          - class:\n          - Iris-Setosa\n          - Iris-Ver
sicolour\n          - Iris-Virginica\n          \n      :Summary Statistics:\n
\n      =====\n      in  Max    Mean    SD    Class Correlation\n      =====\n
===== \n      sepal length:   4.3   7.9   5.84   0.83   0.7826\n      sepal width
h:   2.0  4.4   3.05   0.43  -0.4194\n      petal length:   1.0  6.9   3.76   1.76   0.
9490 (high!)\n      petal width:   0.1  2.5   1.20   0.76   0.9565 (high!)\n      =====
===== \n\n      :Missing Attribute Values:
None\n      :Class Distribution: 33.3% for each of 3 classes.\n      :Creator: R.A. Fisher\n
:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)\n      :Date: July, 1988\n\nThe fa
mous Iris database, first used by Sir R.A. Fisher. The dataset is taken\nfrom Fisher\'s
paper. Note that it\'s the same as in R, but not as in the UCI\nMachine Learning Reposit
ory, which has two wrong data points.\n\nThis is perhaps the best known database to be f
ound in the\npattern recognition literature. Fisher\'s paper is a classic in the field
and\nis referenced frequently to this day. (See Duda & Hart, for example.) The\ndata s
et contains 3 classes of 50 instances each, where each class refers to a\ntype of iris p
lant. One class is linearly separable from the other 2; the\nlatter are NOT linearly se
parable from each other.\n\n.. topic:: References\n\n      - Fisher, R.A. "The use of multi
ple measurements in taxonomic problems"\n      Annual Eugenics, 7, Part II, 179-188 (193
6); also in "Contributions to\n      Mathematical Statistics" (John Wiley, NY, 1950).\n
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis.\n      (Q32
7.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.\n      - Dasarthy, B.V. (198
0) "Nosing Around the Neighborhood: A New System\n      Structure and Classification Rule
for Recognition in Partially Exposed\n      Environments". IEEE Transactions on Pattern
Analysis and Machine\n      Intelligence, Vol. PAMI-2, No. 1, 67-71.\n      - Gates, G.W. (1
972) "The Reduced Nearest Neighbor Rule". IEEE Transactions\n      on Information Theor
y, May 1972, 431-433.\n      - See also: 1988 MLC Proceedings, 54-64. Cheeseman et al"s AU
TOCLASS II\n      conceptual clustering system finds 3 classes in the data.\n      - Many, m
any more ...', 'feature_names': ['sepal length (cm)', 'sepal width (cm)', 'petal length
```

```
(cm)', 'petal width (cm)']], 'filename': 'iris.csv', 'data_module': 'sklearn.datasets.data'}  
a'}
```

3) Split the dataset into training and testing sets as it allows you to evaluate the model's performance, prevent overfitting, tune hyperparameters, and perform cross-validation.

```
In [3]: X_train, X_test, y_train, y_test = train_test_split(df[iris.feature_names], df['target'])  
  
X_train
```

```
Out[3]:
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
61	5.9	3.0	4.2	1.5
92	5.8	2.6	4.0	1.2
112	6.8	3.0	5.5	2.1
2	4.7	3.2	1.3	0.2
141	6.9	3.1	5.1	2.3
...
9	4.9	3.1	1.5	0.1
103	6.3	2.9	5.6	1.8
67	5.8	2.7	4.1	1.0
117	7.7	3.8	6.7	2.2
47	4.6	3.2	1.4	0.2

112 rows × 4 columns

```
In [4]: X_test
```

Out [4] :

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
114	5.8	2.8	5.1	2.4
62	6.0	2.2	4.0	1.0
33	5.5	4.2	1.4	0.2
107	7.3	2.9	6.3	1.8
7	5.0	3.4	1.5	0.2
100	6.3	3.3	6.0	2.5
40	5.0	3.5	1.3	0.3
86	6.7	3.1	4.7	1.5
76	6.8	2.8	4.8	1.4
71	6.1	2.8	4.0	1.3
134	6.1	2.6	5.6	1.4
51	6.4	3.2	4.5	1.5
73	6.1	2.8	4.7	1.2
54	6.5	2.8	4.6	1.5
63	6.1	2.9	4.7	1.4
37	4.9	3.6	1.4	0.1
78	6.0	2.9	4.5	1.5
90	5.5	2.6	4.4	1.2
45	4.8	3.0	1.4	0.3
16	5.4	3.9	1.3	0.4
121	5.6	2.8	4.9	2.0
66	5.6	3.0	4.5	1.5
24	4.8	3.4	1.9	0.2
8	4.4	2.9	1.4	0.2
126	6.2	2.8	4.8	1.8
22	4.6	3.6	1.0	0.2
44	5.1	3.8	1.9	0.4
97	6.2	2.9	4.3	1.3
93	5.0	2.3	3.3	1.0
26	5.0	3.4	1.6	0.4
137	6.4	3.1	5.5	1.8
84	5.4	3.0	4.5	1.5
27	5.2	3.5	1.5	0.2
127	6.1	3.0	4.9	1.8
132	6.4	2.8	5.6	2.2
59	5.2	2.7	3.9	1.4
18	5.7	3.8	1.7	0.3
83	6.0	2.7	5.1	1.6

```
In [6]: y_train
```

```
Out[6]: 61      1
          92      1
          112     2
           2      0
          141     2
          ..
           9      0
          103     2
           67     1
          117     2
           47     0
Name: target, Length: 112, dtype: int32
```

```
In [7]: y_test
```

```
Out[7]: 114     2
          62     1
          33     0
          107    2
           7     0
          100    2
          40     0
          86     1
          76     1
          71     1
          134    2
          51     1
          73     1
          54     1
          63     1
          37     0
          78     1
          90     1
          45     0
          16     0
          121    2
          66     1
          24     0
           8     0
          126    2
          22     0
          44     0
          97     1
          93     1
          26     0
          137    2
          84     1
          27     0
          127    2
          132    2
          59     1
          18     0
          83     1
Name: target, dtype: int32
```

4) Train the model

```
In [8]: knn = KNeighborsClassifier(n_neighbors=1)
         knn.fit(X_train, y_train)
```

Out [8]:

```
▼ KNeighborsClassifier  
KNeighborsClassifier(n_neighbors=1)
```

5) Test the model and generate a report

In [9]:

```
y_pred = knn.predict(X_test)  
print(classification_report(y_test, y_pred, target_names=iris.target_names))
```

	precision	recall	f1-score	support
setosa	1.00	1.00	1.00	13
versicolor	1.00	0.94	0.97	16
virginica	0.90	1.00	0.95	9
accuracy			0.97	38
macro avg	0.97	0.98	0.97	38
weighted avg	0.98	0.97	0.97	38

6) classification report interpretation

The classification report provides information about the performance of the K-Nearest Neighbors (KNN) classifier on the iris dataset.

1. Precision: Precision measures the proportion of true positive instances among the positive instances that were predicted by the model. In the case of the iris dataset, the precision values for each class are:
 - Setosa: 1.00
 - Versicolor: 1.00
 - Virginica: 0.90 A precision value of 1.00 indicates that all positive instances predicted by the model were true positives, while a value of 0.90 means that 90% of the positive instances predicted by the model were true positives, and 10% were false positives.
1. Recall: Recall measures the proportion of true positive instances that were correctly identified by the model. In the case of the iris dataset, the recall values for each class are:
 - Setosa: 1.00
 - Versicolor: 0.94
 - Virginica: 0.95 A recall value of 1.00 indicates that all true positive instances were correctly identified by the model, while a value of 0.94 and 0.95 means that 94% and 95% of the true positive instances were correctly identified, respectively.
1. F1-score: The F1-score is the harmonic mean of precision and recall, and it provides a balanced measure of both. In the case of the iris dataset, the F1-score values for each class are:
 - Setosa: 1.00
 - Versicolor: 0.97
 - Virginica: 0.95 A F1-score value of 1.00 indicates that the model perfectly balanced precision and recall, while a value of 0.97 and 0.95 means that the model achieved a good balance between precision and recall for the versicolor and virginica classes, respectively.
1. Support: Support represents the number of samples in each class. In the case of the iris dataset, the support values for each class are:

- Setosa: 13
- Versicolor: 16
- Virginica: 9

The overall accuracy of the KNN classifier is 0.97, which means that 97% of the predictions made by the model were correct. The macro average values for precision, recall, and F1-score are 0.97, 0.98, and 0.97, respectively, indicating that the model performed well across all classes.