

```
In [12]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report

In [13]: data = pd.read_csv('r'D:\Download\archive (8)\IRIS.csv')

In [14]: data.head()

Out[14]:   sepal_length  sepal_width  petal_length  petal_width  species
0         5.1         3.5         1.4         0.2  Iris-setosa
1         4.9         3.0         1.4         0.2  Iris-setosa
2         4.7         3.2         1.3         0.2  Iris-setosa
3         4.6         3.1         1.5         0.2  Iris-setosa
4         5.0         3.6         1.4         0.2  Iris-setosa

In [15]: data.info

Out[15]: <bound method DataFrame.info of      sepal_length  sepal_width  petal_length  petal_width  species
0         5.1         3.5         1.4         0.2  Iris-setosa
1         4.9         3.0         1.4         0.2  Iris-setosa
2         4.7         3.2         1.3         0.2  Iris-setosa
3         4.6         3.1         1.5         0.2  Iris-setosa
4         5.0         3.6         1.4         0.2  Iris-setosa
...         ...         ...         ...         ...         ...
145        6.7         3.0         5.2         2.3  Iris-virginica
146        6.3         2.5         5.8         1.9  Iris-virginica
147        6.5         3.0         5.2         2.0  Iris-virginica
148        6.2         3.4         5.4         2.3  Iris-virginica
149        6.9         3.0         5.1         1.8  Iris-virginica
[150 rows x 5 columns]>

In [11]: from sklearn.datasets import load_iris
data = load_iris()
print(data.DESCR)
..._iris_dataset:

Iris plants dataset
-----
**Data Set Characteristics:**

 :Number of Instances: 150 (50 in each of three classes)
 :Number of Attributes: 4 numeric, predictive attributes and the class
 :Attribute Information:
   - sepal length in cm
   - sepal width in cm
   - petal length in cm
   - petal width in cm
   - class:
     - Iris-Setosa
     - Iris-Versicolour
     - Iris-Virginica

 :Summary Statistics:

=====
      Min    Max   Mean     SD   Class Correlation
=====
sepal length:  4.3  7.9  5.84   0.83   0.7826
sepal width:   2.0  4.4  3.85   0.43  -0.4184
petal length:   1.0  6.9  5.78   1.76   0.9490 (high)
petal width:    0.1  2.5  1.20   0.78   0.9565 (high)
=====

 :Missing Attribute Values: None
 :Class Distribution: 33.3% for each of 3 classes.
 :Creator: R.A. Fisher
 :Donor: Richard Marshall (MARSHALL@PLUMBio.arc.nasa.gov)
 :Date: July, 1988

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken
from Fisher's paper. Note that it's the same as in R, but not as in the UCI
Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the
pattern recognition literature. Fisher's paper is a classic in the field and
is referenced frequently to this day. (See Duda & Hart, for example.) The
data set contains 3 classes of 50 instances each, where each class refers to a
type of Iris plant. One class is linearly separable from the other 2; the
latter are NOT linearly separable from each other.

... topic:: References

- Fisher, R.A. "The use of multiple measurements in taxonomic problems"
  Annual Eugenics, 7, Part II: 297-388 (1936); also in "Contributions to
  Mathematical Statistics" (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis.
  (Q227.D83) John Wiley & Sons. ISBN 0-471-22661-1. See page 219.
- Dasarthy, B.V. (1989) "Nosing Around the Neighborhood: A New System
  Structure and Classification Rule for Recognition in Partially Exposed
  Environments". IEEE Transactions on Pattern Analysis and Machine
  Intelligence, Vol. 9(MC-2), No. 5, 67-71.
- Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule". IEEE Transactions
  on Information Theory, May 1972, 431-433.
- See also: 1988 IJCC Proceedings, 34-64. Chesselan et al.'s AUTOCLASS II
  conceptual clustering system finds 3 classes in the data.
- Many, many more ...

In [13]: data.keys()

Out[13]: dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_module', 'columns'])

In [14]: print(data.data[5])
print(data.feature_names)

[[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]]
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']

In [15]: print(data.target)
print(data.target_names)

[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
 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]
['setosa' 'versicolor' 'virginica']

In [17]: df = pd.DataFrame(data.data)
df.head()

Out[17]:   0  1  2  3
0  5.1 3.5 1.4 0.2
1  4.9 3.0 1.4 0.2
2  4.7 3.2 1.3 0.2
3  4.6 3.1 1.5 0.2
4  5.0 3.6 1.4 0.2

In [18]: df.columns = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width']
df.head()

Out[18]:   sepal_length  sepal_width  petal_length  petal_width
0         5.1         3.5         1.4         0.2
1         4.9         3.0         1.4         0.2
2         4.7         3.2         1.3         0.2
3         4.6         3.1         1.5         0.2
4         5.0         3.6         1.4         0.2

In [19]: target = pd.DataFrame(data.target)
target = target.rename(columns = {'target':''})
print(target.head())
print('all target column have 3 value: ',target.target.unique())

target
0      0
1      0
2      0
3      0
4      0
all target column have 3 value:  [0 1 2]

In [20]: df = pd.concat([df, target], axis = 1)
df.head()

Out[20]:   sepal_length  sepal_width  petal_length  petal_width  target
0         5.1         3.5         1.4         0.2      0
1         4.9         3.0         1.4         0.2      0
2         4.7         3.2         1.3         0.2      0
3         4.6         3.1         1.5         0.2      0
4         5.0         3.6         1.4         0.2      0

In [21]: df.describe()

Out[21]:   sepal_length  sepal_width  petal_length  petal_width  target
count    150.000000    150.000000    150.000000    150.000000    150.000000
mean       5.843333    3.057333    3.758000    1.198333    1.000000
std       0.828066    0.438666    1.765298    0.762238    0.819232
min       4.300000    2.000000    1.000000    0.100000    0.000000
25%      5.100000    2.800000    1.600000    0.300000    0.000000
50%      5.800000    3.000000    3.450000    1.300000    1.000000
75%      6.400000    3.300000    5.100000    1.800000    2.000000
max       7.900000    4.400000    6.900000    2.500000    2.000000

In [22]: print('see what type of data you have:\n',df.dtypes)
print('-----')
print('check for missing values:\n',df.isnull().sum())

see what type of data you have:
sepal_length    float64
sepal_width     float64
petal_length     float64
petal_width     float64
target          int32
dtype: object
-----
check for missing values:
sepal_length    0
sepal_width     0
petal_length    0
petal_width     0
target          0
dtype: int64

In [23]: plt.rcParams['figure.figsize'] = [18,5]
sns.heatmap(df.corr(), annot = True);

sepal_length  1 -0.12 0.87 0.82 0.78
sepal_width -0.12 1 -0.43 -0.37 -0.43
petal_length 0.87 -0.43 1 0.96 0.95
petal_width 0.82 -0.37 0.96 1 0.96
target 0.78 -0.43 0.95 0.96 1

In [24]: x_index = 0
y_index = 1

formatter = plt.FuncFormatter(lambda i, *args: data.target_names[int(i)])
plt.figure(figsize=(8,6))
plt.scatter(data.data[:, x_index], data.data[:, y_index], c=data.target)
plt.colorbar(ticks=[0, 1, 2], format=formatter)
plt.xlabel(data.feature_names[x_index])
plt.ylabel(data.feature_names[y_index])
plt.tight_layout()
plt.show()

In [25]: x_index = 2
y_index = 3

formatter = plt.FuncFormatter(lambda i, *args: data.target_names[int(i)])
plt.figure(figsize=(8,6))
plt.scatter(data.data[:, x_index], data.data[:, y_index], c=data.target)
plt.colorbar(ticks=[0, 1, 2], format=formatter)
plt.xlabel(data.feature_names[x_index])
plt.ylabel(data.feature_names[y_index])
plt.tight_layout()
plt.show()

In [26]: x_index = 2
y_index = 3

formatter = plt.FuncFormatter(lambda i, *args: data.target_names[int(i)])
plt.figure(figsize=(8,6))
plt.scatter(data.data[:, x_index], data.data[:, y_index], c=data.target)
plt.colorbar(ticks=[0, 1, 2], format=formatter)
plt.xlabel(data.feature_names[x_index])
plt.ylabel(data.feature_names[y_index])
plt.tight_layout()
plt.show()

In [27]: x_index = 2
y_index = 3

formatter = plt.FuncFormatter(lambda i, *args: data.target_names[int(i)])
plt.figure(figsize=(8,6))
plt.scatter(data.data[:, x_index], data.data[:, y_index], c=data.target)
plt.colorbar(ticks=[0, 1, 2], format=formatter)
plt.xlabel(data.feature_names[x_index])
plt.ylabel(data.feature_names[y_index])
plt.tight_layout()
plt.show()

In [28]: for feat in ['sepal_length', 'sepal_width', 'petal_length', 'petal_width']:
df[feat].hist(ecp='red',color='green')
plt.suptitle(feat)
plt.show()
print('-----')

sepal_length
-----
sepal_width
-----
petal_length
-----
petal_width
-----

In [29]: X = df.copy()
y = X.pop('target')
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1, stratify = y)

In [31]: print('How split:\n')
print(X_train, X_train.shape)
print(y_train, y_train.shape)
print(X_test, X_test.shape)
print(y_test, y_test.shape)

How split:
X_train (128, 4)
y_train (128,)
X_test (30, 4)
y_test (30,)

In [32]: scaler = StandardScaler()
X_train = pd.DataFrame(scaler.fit_transform(X_train), columns=X_train.columns)
X_test = pd.DataFrame(scaler.transform(X_test), columns=X_test.columns)

In [33]: print(df.target.value_counts(normalize=True))
print('\n The baseline prediction for this model is 1/3')

0    0.333333
1    0.333333
2    0.333333
Name: target, dtype: float64

The baseline prediction for this model is 1/3

Logistic Regression model

In [34]: lg = LogisticRegression()
lg.fit(X_train, y_train)
print('accuracy of model: ',lg.score(X_test, y_test))
cv = cross_val_score(lg, X_train, y_train, cv=10)
print('accuracy of model after CVS: ',np.mean(cv))

accuracy of model:  0.9656666666666667
accuracy of model after CVS:  0.9699999999999998

In [35]: df_coef = pd.DataFrame(lg.coef_, columns=X_train.columns)
df_coef

Out[35]:   sepal_length  sepal_width  petal_length  petal_width
0    -1.02746    1.001818   -1.036991   -1.667970
1    0.402982   -0.323432   -0.277761   -0.650011
2    0.699764   -0.678306   2.114653   2.317889

In [36]: predict = lg.predict(X_test)
compare = pd.DataFrame({'actual': y_test, 'predicted': predict})
compare['compar.reset_index(drop = True)']

Out[36]:   actual  predicted
0      2         2
1      0         0
2      1         1
3      0         0
4      0         0
5      0         0
6      2         2
7      2         2
8      2         2
9      1         1

In [37]: pd.DataFrame(confusion_matrix(y_test, predict, labels=[2, 1, 0]),
                      index=[2, 1, 0], columns=[2, 1, 0])

Out[37]:   2  1  0
2  9  1  0
1  0 10  0
0  0  0 10

In [38]: print(classification_report(y_test, predict))

precision    recall  f1-score   support

0       1.00      1.00      1.00       10
1       1.00      1.00      1.00       10
2       1.00      1.00      1.00       10

accuracy   0.97      0.97      0.97       30
macro avg   0.97      0.97      0.97       30
```

