

Київський національний університет імені Тараса Шевченка
Факультет інформаційних технологій

Кафедра програмних систем і технологій

ЗВІТ

з практичної роботи № 3

Тема: «Алгоритми класифікації із застосуванням бібліотек»

Дисципліна «Спеціалізоване програмування автоматизованих систем»

Підготував:
студент гр. ПЗ-33(1)
Мішак Максим

Перевірила:
Ніколаєнко Анастасія Юріївна

Завдання :

Класифікатор Методу k -найближчих сусідів (відстань Хеммінга). Дані з файлу «seeds.csv».

Хід роботи :

Для виконання роботи будемо використовувати наступний перелік бібліотек :

- **pandas**
- **matplotlib**
- **seaborn**
- **sklearn**

Для класифікації будемо використовувати готове рішення з бібліотеки **sklearn** , а саме **KNeighborsClassifier** . Використання цього методу потребує знання деяких нюансів , до прикладу розділення наданого датасету відбувається наступним чином :

```
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.2,  
random_state=42)
```

цю конструкцію можна записати будь яким зручним способом - але такий запис рахується канонічним . Для інфографіки будемо використовувати елемент бібліотеки **sklearn.metrics classification_report** , **confusion_matrix**. Це нам потрібно для виведення результатів обробки програмою датасету .

Згідно методичних вказівок , згідно з моїм варіантом я маю датасет з назвою seeds.csv . Структура даних в нього виглядає так :

```
"V1", "V2", "V3", "V4", "V5", "V6", "V7", "class"  
15.26, 14.84, 0.871, 5.763, 3.312, 2.221, 5.22, 1
```

Маємо данні , які репрезентують вектори , деякі з яких суттєво відрізняються значеннями . Для того , щоб ці значення не заважали класифікатору працювати будемо виконувати нормалізацію даних :

```
scaler = StandardScaler()  
scaler.fit(X_train)  
X_train = scaler.transform(X_train)  
X_test = scaler.transform(X_test)
```

```
knn.fit(X_test, y_test)
```

Це дозволяє “форматувати” данні до приблизного одного вигляду та значення .

Код застосунку

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sb
```

```
from pandas.plotting import scatter_matrix
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.preprocessing import StandardScaler
```

```
from sklearn.metrics import classification_report,  
confusion_matrix ,ConfusionMatrixDisplay
```

```
#read the data from file
```

```
DataSet =
```

```
pd.read_csv('/Users/macbook/Desktop/SPAS3/dataset/se  
ds.csv')
```

```
knn = KNeighborsClassifier(n_neighbors=5,  
metric='hamming')
```

```
#set data as X and y
```

```
X = DataSet.iloc[:, :-1].values
```

```
y = DataSet.iloc[:, 7].values
```

```
#spliting the data fot 2 frames
```

```
c
```

```
#scale and transform datda
```

```
scaler = StandardScaler()
```

```
scaler.fit(X_train)
```

```
X_train = scaler.transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

```
knn.fit(X_test, y_test)

#make a "prediction" for the data
y_pred = knn.predict(X_test)

#print the default data
print(DataSet)

#printing of data shapes
print(X_train.shape, X_test.shape, y_train.shape,
y_test.shape)

#printing reports confusion matrix and class report
print(classification_report(y_test, y_pred))

cm = confusion_matrix(y_test, y_pred)

print(cm)

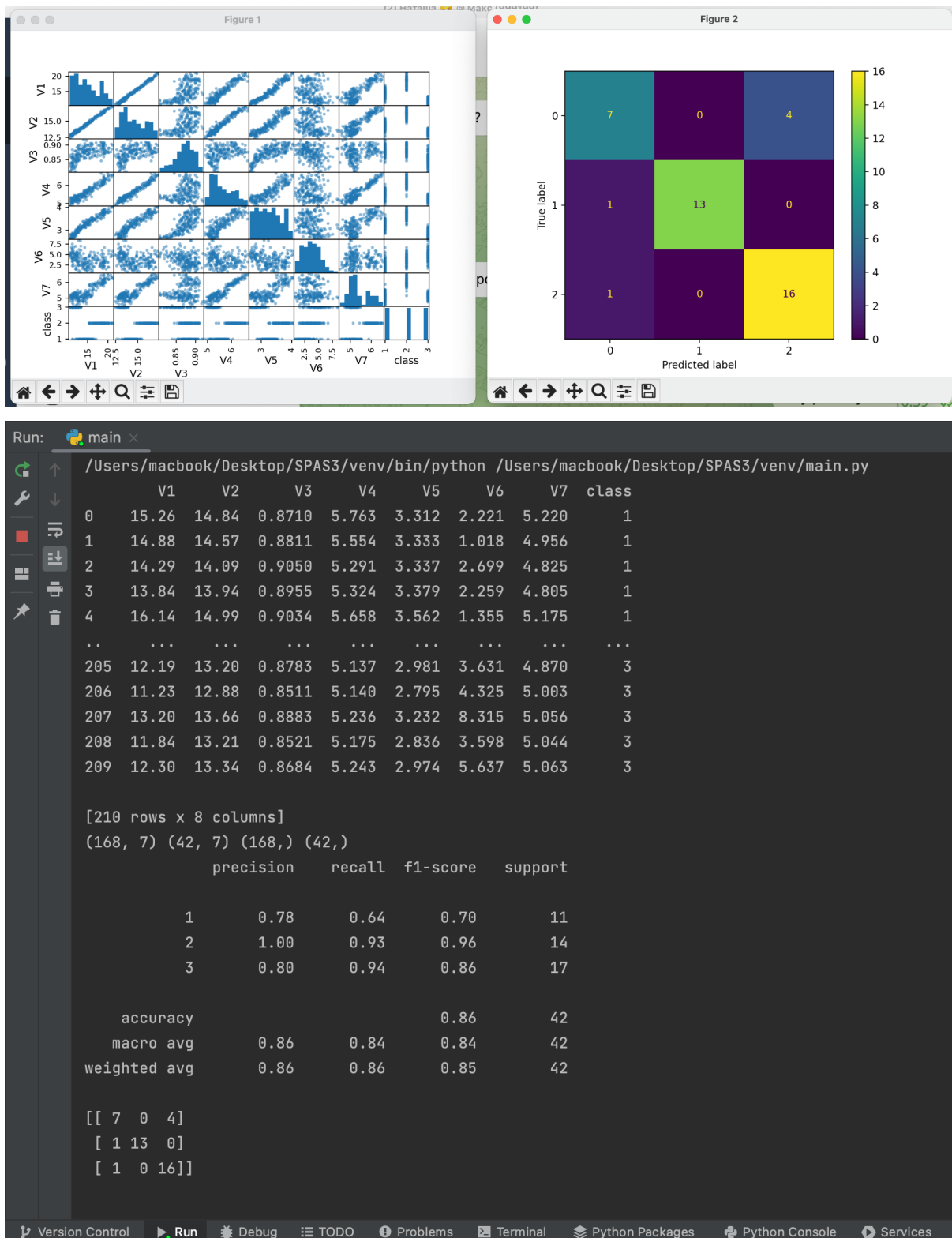
#plotting a scatter matrix
scatter_matrix(DataSet)

#plotting the confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm)

disp.plot()

plt.show()
```

Скріншоти виконання програми



Висновок :

Було опрацьовано методи класифікації бібліотеки sklearn , а саме метод knn. В результаті можемо побачити , що при теперішніх налаштуваннях

функції пошуку найблищого сусіда середня ефективність є 86% . Така точність була досягнута шляхом емпіричного підбору параметра кількості найблищих сусідів . Для кожного окремого класу це показник варіюється , так як різні класи можуть приблизно співпадати по показникам .