

Klasyfikacja i regresja: Wstęp do biblioteki `scikit-learn`

24 października 2018

Opis pliku z zadaniami

Wszystkie zadania na zajęciach będą przekazywane w postaci plików `.pdf`, sformatowanych podobnie do tego dokumentu. Zadania będą różnego rodzaju. Za każdym razem będą one odpowiednio oznaczone:

- Zadania do wykonania na zajęciach oznaczone są symbolem \triangle – nie są one punktowane, ale należy je wykonać w czasie zajęć.
- Punktowane zadania do wykonania na zajęciach oznaczone są symbolem \diamond – należy je wykonać na zajęciach i zaprezentować prowadzącemu.
- Zadania do wykonania w domu oznaczone są symbolem \star – są one punktowane, należy je dostarczyć w sposób podany przez prowadzącego i w wyznaczonym terminie (zwykle przed kolejnymi zajęciami).

1 Zapoznanie się z biblioteką sklearn



Treść

W tym ćwiczeniu przedstawiony i omówiony jest krótki skrypt w języku python. Na podstawie danych uczących zostanie wytrenowany klasyfikator potrafiący rozróżnić cyfry na prostych obrazach o rozmiarze 8×8 . Następnie klasyfikator zostanie wykorzystany do klasyfikacji pojedynczego przykładu.

Należy wykonać następujące kroki:

1. Zapoznaj się ze stroną główną projektu sklearn:

<http://scikit-learn.org/>

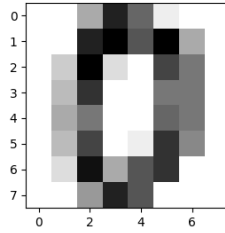
- Przewodnik użytkownika dostępny jest tutaj:
http://scikit-learn.org/stable/user_guide.html
- Bieżące ćwiczenie zostało opracowane na podstawie tutorialu:
<http://scikit-learn.org/stable/tutorial/basic/tutorial.html>
- Dodatkowo warto zapoznać się z typami danych obsługiwanymi przez sklearn oraz proponowanymi metodami wczytywania danych: <http://scikit-learn.org/stable/datasets/index.html#external-datasets>

2. Zaimportuj bibliotekę zawierającą przykładowe dane. Wczytaj plik zawierający obrazy cyfr oraz wypisz dane.

```
1 from sklearn import datasets
2
3 #Load the digits dataset
4 digits = datasets.load_digits()
5
6 #Print examples with their features
7 print(digits.data)
8
9 #Array with output (decision) values
10 digits.target
11
12 #The first example
13 digits.images[0]
```

3. Powyższe dane są zapisane wektorowo. Można jednak łatwo je przedstawić jako bitmapę.

```
1 import matplotlib.pyplot as plt
2
3 #Display the first digit
4 plt.figure(1, figsize=(3, 3))
5 plt.imshow(digits.images[0], cmap=plt.cm.gray_r,
6             interpolation='nearest')
7 plt.show()
```



4. Wytrenuj Twój pierwszy klasyfikator używając wszystkich przykładów oprócz ostatniego.

```

1 from sklearn import svm
2
3 #Initialize the SVM classifier
4 clf = svm.SVC(gamma=0.001, C=100.)
5
6 #Train the classifier on training data
7 clf.fit(digits.data[:-1], digits.target[:-1])

```

5. Zastosuj klasyfikator do ostatniego przykładu i sprawdź czy predykcja jest poprawna (wizualnie oraz sprawdzając prawdziwą etykietę w danych).

```

1 #Predict the last digit
2 clf.predict(digits.data[-1:])
3
4 #The digit looks like
5 #Display the last digit
6 plt.figure(1, figsize=(3, 3))
7 plt.imshow(digits.images[-1], cmap=plt.cm.gray_r,
8             interpolation='nearest')
9 plt.show()
10 #The true label of the last image
11 digits.target[-1]

```

6. Nauczony klasyfikator można zapisać wykorzystując bibliotekę `pickle` lub `joblib`.

```

1 import pickle
2
3 s = pickle.dumps(clf)
4 clf2 = pickle.loads(s)
5 clf2.predict(digits.data[-1:])
6
7 from sklearn.externals import joblib
8 joblib.dump(clf, '/home/username/clf.joblib')
9 clf3 = joblib.load('/home/username/clf.joblib')
10 clf3.predict(digits.data[-1:])

```