

Klasyfikacja i regresja: Walidacja krzyżowa

28 listopada 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 Prawidłowa realizacja walidacji krzyżowej 10p.◇

Treść

Walidacja krzyżowa jest metodą estymacji jakości algorytmów uczących na małych zbiorach danych. Pozwala ona także na odpowiednie dobranie hiperparametrów danego algorytmu uczącego. Jednak jej błędne wykorzystanie może doprowadzić do błędnych wniosków. Zadaniem do wykonania podczas zajęć jest analiza dwóch scenariuszów walidacji krzyżowej:

1. Selekcja najlepszych parametrów na zbiorze uczącym i uruchomienie walidacji krzyżowej,
2. Selekcja najlepszych parametrów podczas walidacji krzyżowej.

Należy zastanowić się, czy oba scenariusze są poprawne. Jeżeli nie, to który jest błędny.

W celu realizacji zadania wykonaj następujące kroki:

1. Zapoznaj się z kodem zamieszczonym w pliku `compare_cv_scenarios.py` (plik jest dostępny na stronie przedmiotu oraz wydrukowany poniżej).
2. Zapoznaj się z opisem metod `SelectFromModel` oraz `SelectKBest`.
3. Uruchom eksperyment dla syntetycznych danych. Wypróbuj różne ustawienia liczby przykładów, liczby cech oraz liczby wybieranych cech. Liczbę iteracji (`fold`) pozostaw bez zmian.
4. Czy udało się zaobserwować nieprawidłowości?
5. Odpowiednio rozszerz kod, aby wyniki otrzymać wielokrotnie i odpowiednio je uśrednić.
6. Wyciągnij ostateczne wnioski.

```
1 import numpy as np
2 from sklearn.metrics import confusion_matrix
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import cross_val_predict
5 from sklearn.pipeline import Pipeline
6 from sklearn.feature_selection import SelectFromModel,
  SelectKBest
7
8
9 def create_random_data(n_instances, n_features):
10     X = 2*np.random.rand(n_instances, n_features) - 1
11     scores = 2*np.random.rand(n_instances) - 1
12     where_positive = scores > 0
13     where_negative = scores <= 0
14     scores[where_positive] = 1
15     scores[where_negative] = 0
```

```

16     y = np.array(scores.ravel().tolist(), dtype=np.int32)
17     return X, y
18
19 def print_evaluation(y, predicted):
20     print("True: {0}/{1} = {2}".format(sum(y), len(y), float(sum
21         (y))/len(y)))
22     print("Confusion matrix")
23     cm = confusion_matrix(y, predicted)
24     print(cm)
25     correct = cm[0, 0] + cm[1, 1]
26     incorrect = cm[1, 0] + cm[0, 1]
27     print(float(correct)/(correct + incorrect))
28
29 def selection_outside_cv(X, y, folds, sfm = SelectFromModel(
30     LogisticRegression())):
31     print("Scenario: Selection + CV(LR)")
32     sfm.fit(X, y)
33     X_selected = sfm.transform(X)
34     predicted = cross_val_predict(LogisticRegression(),
35         X_selected, y, cv=folds)
36     print_evaluation(y, predicted)
37
38 def selection_inside_cv(X, y, folds, sfm = SelectFromModel(
39     LogisticRegression())):
40     print("Scenario: CV(Selection, LR)")
41     clf = Pipeline([
42         ('feature_selection', sfm),
43         ('classification', LogisticRegression())
44     ])
45     predicted = cross_val_predict(clf, X, y, cv=folds)
46     print_evaluation(y, predicted)
47
48 def main():
49     n_instances = 1000
50     n_features = 1000
51     folds = 10
52     n_top_features = 2
53
54     X, y = create_random_data(n_instances, n_features)
55
56     #use one of the methods for feature selection
57     #sfm = SelectFromModel(LogisticRegression())
58     sfm = SelectKBest(k=n_top_features)
59
60     selection_outside_cv(X, y, folds, sfm)
61     selection_inside_cv(X, y, folds, sfm)
62
63 if __name__ == "__main__":
64     main()

```

compare_cv_scenarios.py