

Scikit Learn

Giảm chiều ma trận đặc trưng, tìm kiếm lưới và lưu
mô hình

HK2, Năm học 2019 - 2020

Giảm chiều dữ liệu - Dimensionality Reduction

Dimensionality Reduction (giảm chiều dữ liệu) là một trong những kỹ thuật quan trọng trong Máy học.

Các vector đặc trưng trong các bài toán thực tế có thể có số chiều (cột) lớn. Nếu thực hiện lưu trữ và tính toán trực tiếp trên dữ liệu có số chiều cao này thì sẽ gặp khó khăn về việc lưu trữ và tốc độ tính toán.

Vì vậy giảm số chiều dữ liệu là một bước quan trọng trong nhiều bài toán.

Đây cũng được coi là phương pháp nén dữ liệu.

Một trong những cách có thể sử dụng là Principal Component Analysis (PCA)

Giảm chiều dữ liệu - Dimensionality Reduction

```
1 from sklearn import datasets
2 from sklearn.decomposition import PCA
3 from sklearn.model_selection import train_test_split
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.metrics import accuracy_score
6
7 wine = datasets.load_wine()
8 X, y = wine.data, wine.target
9 X_train_org, X_test_org, y_train, y_test = train_test_split(X, y,
10     test_size = 0.3, random_state = 10)
# khai tao PCA
11 pca = PCA(n_components=10, random_state=10)
# huan luyen PCA
12 pca.fit(X_train_org)
# chuyen doi cho tap train va test
13 X_train = pca.transform(X_train_org)
14 X_test = pca.transform(X_test_org)
```

Giảm chiều dữ liệu - Dimensionality Reduction

```
1 model = KNeighborsClassifier()
2 model.fit(X_train, y_train)
3 preds = model.predict(X_test)
4 score = accuracy_score(y_test, preds)
5 print(score)
```

Tìm tham số tốt cho mô hình

Các mô hình đang sử dụng thường bao gồm nhiều hyperparameters và tùy thuộc vào bài toán cụ thể sẽ có các hyperparameters thích hợp.

Để tìm hyperparameters phù hợp cho từng bài toán cụ thể, cách đơn giản là thử nhiều biến khác nhau.

- Random Search: từ những biến cần tìm chọn ngẫu nhiên giá trị cho các hyperparameter để đánh giá độ chính xác của mô hình
- Grid Search: giả dụ có 02 hyperparameter có các giá trị có thể nhận từ 0-9. Grid Search sẽ ghép lần lượt từng giá trị của hyperparam 1 với hyperparam 2 để tính toán độ chính xác của model.

Tìm tham số tốt cho mô hình

```
1 from sklearn import datasets
2 from sklearn.model_selection import train_test_split
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.metrics import accuracy_score
6
7 irisX, irisY = datasets.load_iris(return_X_y=True)
8 X_train, X_test, y_train, y_test = train_test_split(
9     irisX,irisY,test_size=0.4,
10    random_state = 1)
11
12 parameters = {'n_neighbors':[1,2,3,4,5,6,7,8,9,10]}
13 knn = KNeighborsClassifier()
14 # scoring co the: 'neg_mean_squared_error', 'r2'
15 grid_result = clf.fit(X_train, y_train)
```

Tìm tham số tốt cho mô hình

```
1 print(grid_result)
2 print(grid_result.best_params_)
3 means = clf.cv_results_['mean_test_score']
4 print(means)
5
6 model = grid_result.best_estimator_
7
8 preds = model.predict(X_test)
9 print(accuracy_score(y_test, preds))
```

Tìm mô hình

```
1 from sklearn import datasets
2 from sklearn.model_selection import train_test_split
3 from sklearn.svm import SVC
4 from sklearn.linear_model import SGDClassifier
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.metrics import accuracy_score
8 from sklearn.externals import joblib
9
10 X, y = datasets.load_iris(return_X_y = True)
11 X_train, X_test, y_train, y_test = train_test_split(X
12 , y, test_size = 0.4,
random_state = 12)
```

Tìm mô hình

```
1 models = []
2 models.append((‘SupportVectorClassifier’, SVC()))
3 models.append((‘SGDClassifier’, SGDClassifier()))
4 models.append((‘RandomForest’, RandomForestClassifier()
    )))
5 models.append((‘KNN’, KNeighborsClassifier()))
6
7 results = []
8 scoring = ‘accuracy’
9 for name, model in models:
10     model.fit(X_train, y_train)
11     pred = model.predict(X_test)
12     result = accuracy_score(y_test, pred)
13     results.append((name, result))
14     joblib.dump(model, name + ‘.pkl’)
```

Tìm mô hình

```
1 sample = [[3.2 , 5.3 , 4.1 , 2.8] , [2.1 ,3.2 ,1.9  
 ,0.5]]  
2 for name, result in results:  
3     model = joblib.load(name+'.pkl')  
4     pred = model.predict(sample)  
5     print('Model %s co do chinh xac %f' %(name, result  
 ))  
6     print(pred)
```

Bài 1: Digits

Sử dụng tập dữ liệu `digits` của `sklearn` (`load_digits`), thực hiện các yêu cầu sau

- Kiểm tra có bao nhiêu đặc trưng/thuộc tính đối với tập dữ liệu trên (số cột của ma trận đặc trưng)
- Train và test được chia theo tỷ lệ: 8:2
- Sử dụng `PCA` để giảm chiều dữ liệu xuống 16
- Sử dụng mô hình KNN để đánh giá việc học trên tập dữ liệu và tập dữ liệu đã giảm chiều: việc giảm chiều có tác dụng cải thiện độ chính xác không?

Bài 2: Titanic

Sử dụng titanic từ file *titanic_pre.csv* (file có header nên không cần truyền names cho nó).

Giả sử mô hình SVC có 02 hyperparameter: C có giá trị 1 hoặc 10, kernel có giá trị 'linear' hoặc 'rbf'

- Train và test set được chia theo tỷ lệ: 8:2
- Thực hiện chuẩn hóa dữ liệu (chọn 1 trong các cách đã được giới thiệu)
- Sử dụng GridSearchCV để tìm mô hình tốt nhất
- Đánh giá mô hình trên test
- Lưu lại mô hình đó với tên file *svc_titanic.pkl*

Bài 3: Fertility Diagnosis

Sử dụng dữ liệu từ file *fertility_Diagnosis.csv* (file có header nên không cần truyền names cho nó).

- Xác định cần mô hình phân lớp hay hồi quy?
- Chia tập dữ liệu 8:2
- Chuẩn hóa dữ liệu (tùy ý)
- Giả sử có 03 mô hình cần xem xét: KNN, DecisionTree và RandomForest (sử dụng tham số mặc định). In đánh giá của các mô hình trên và cho biết mô hình/những mô hình nào được xem xét chọn