

Scikit Learn

Tiền xử lý dữ liệu và mô hình hồi quy

HK1, Năm học 2022 - 2023

Chuẩn hóa dữ liệu

Chuẩn hóa dữ liệu là một bước quan trọng trong bài toán máy học.

Trung tâm hóa dữ liệu (Centering data): Đưa các điểm dữ liệu trong tập dữ liệu về xoay quanh giá trị 0 thay vì xoay quanh giá trị trung bình của tập dữ liệu

Co giãn dữ liệu (Scaling data): chuẩn hóa phạm vi của các đặc trưng dữ liệu

- Chuẩn hóa min-max (rescaling)

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

- Co giãn trung bình (mean normalization)

$$x' = \frac{x - \text{average}(x)}{\max(x) - \min(x)}$$

Chuẩn hóa dữ liệu

Co giãn dữ liệu (Scaling data): chuẩn hóa phạm vi của các đặc trưng dữ liệu

- Chính quy hóa (standardization): giúp giá trị của mỗi đặc trưng có trung bình bằng 0 và phương sai bằng 1

$$x' = \frac{x - \bar{x}}{\sigma}$$

- Vector đơn vị: biến đổi vector đặc trưng sau cho khi biến đổi có độ dài bằng 1

Chuẩn hóa dữ liệu

Mã hóa đặc trưng nhóm: thuộc tính của dữ liệu có thể lưu trữ dưới dạng chữ. Ví dụ: blue, green; small, medium, large... Có nhiều phương pháp mã hóa đặc trưng dạng nhóm, phổ biến:

- Mã hóa số (Numeric Encoding): gán mỗi giá trị của đặc trưng dạng nhóm thành số theo từng đôi một \Leftrightarrow còn gọi là *mã hóa nhãn (label encoding)* hay *mã hóa số nguyên (integer encoding)*
- Mã hóa one-hot (One-hot encoding): biến mỗi giá trị của thuộc tính thành một biến nhị phân.

Ví dụ: Thuộc tính màu sắc có 03 giá trị red, green, blue. Khi áp dụng mã hóa thì các giá trị lần lượt được biểu diễn bằng vector nhị phân là: 100, 010, 001

Chuẩn hóa min-max

```
1 from sklearn import datasets
2 from sklearn import preprocessing
3 from sklearn.model_selection import train_test_split
4
5 iris = datasets.load_iris()
6 X = iris.data
7 y = iris.target
8 X_train, X_test, y_train, y_test = train_test_split(X
    , y, test_size = 0.4, random_state = 123)
9
10 scaler = preprocessing.MinMaxScaler()
11 X_train_preprocess = scaler.fit_transform(X_train)
12 X_test_preprocess = scaler.transform(X_test)
```

Chính quy hóa

```
1 from sklearn import datasets
2 from sklearn import preprocessing
3 from sklearn.model_selection import train_test_split
4
5 iris = datasets.load_iris()
6 X = iris.data
7 y = iris.target
8 X_train, X_test, y_train, y_test = train_test_split(X
    , y, test_size = 0.4, random_state = 123)
9
10 scaler = preprocessing.StandardScaler().fit(X_train)
11 X_train_preprocess = scaler.transform(X_train)
12 X_test_preprocess = scaler.transform(X_test)
```

Mã hóa số

```
1 enc = preprocessing.OrdinalEncoder()
2 X = [['male', 'from US', 'uses Safari'], ['female', ,
3     'from Europe', 'uses Firefox']]
4 enc.fit(X)
5 enc.transform([['female', 'from US', 'uses Safari']])
```

Mã hóa one-hot

```
1 enc = preprocessing.OneHotEncoder()
2 X = [['male', 'from US', 'uses Safari'], ['female', ,
3     'from Europe', 'uses Firefox']]
4 enc.fit(X)
5 enc.transform([['female', 'from US', 'uses Safari'],
6     ['male', 'from Europe', 'uses Safari']]).toarray()
```

Bài 1 - Chuẩn hóa cho IRIS

```
1 from sklearn import preprocessing
2 # import thư viện cần thiết khác
3
4 # load dữ liệu iris
5
6 # chia dữ liệu iris theo tỉ lệ 8:2
7 X_train, X_test, y_train, y_test =
8
9 # thực hiện chuẩn hóa min-max
10 scaler = preprocessing.MixMaxScaler()
11 scaler.fit(X_train)
12 X_train1 = scaler.transform(X_train)
13 X_test1 = scaler.transform(X_test)
14
15 # thiết lập mô hình và huấn luyện mô hình trên X_train1
16 model =
17 model.fit(X_train1, y_train)
18 # thực hiện dự đoán trên X_test1
19 y_pred1 = model1.predict(X_test1)
20 # đánh giá trên y_pred1
21 score1 = accuracy_score(y_test, y_pred1)
```

Bài 2 - Chuẩn hóa cho WINE

- Sử dụng dữ liệu wine
- Chia tập dữ liệu theo tỷ lệ 8: 2
- Thực hiện chính quy hóa dữ liệu
- Cài đặt mô hình KNN và huấn luyện trên tập dữ liệu chưa chuẩn hóa
- Cài đặt mô hình KNN khác và huấn luyện trên tập dữ liệu đã chuẩn hóa
- Cài đặt mô hình SVM cho phân lớp và huấn luyện trên tập dữ liệu đã chuẩn hóa

```
1 from sklearn.svm import SVC  
2 model = SVC()
```

Bài 3 - Dataset Haberman

```
1 import pandas as pd
2 # import package can thiet khac
3 # ten cua cac cot trong dataset
4 names = ['Age', 'Year operation', 'Axillary nodes detected', ,
5           Survival status']
6 # duong dan den file du lieu
7 url = "haberman.csv"
8 # load file csv
9 dataset = pd.read_csv(url, names=names)
10 # in mot so mo ta cho dataset
11 print(dataset.describe())
12 print(dataset.head(10))
13 print(dataset.shape)
14 # chia ra feature matrix X va target y
15 array = dataset.values
16 X = array[:, :3]
17 y = array[:, 3]
18 # chia tap du lieu train, test 7:3
19 # Cai dat mo hinh SVM danh cho phan lop
20 # Danh gia viec huan luyen mo hinh
```

Bài 4 - Dataset Diabetes

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.neighbors import KNeighborsRegressor
3 from sklearn.metrics import mean_squared_error, r2_score
4 # import...
5 # load diabetes tu sklearn
6 diabetes =
7 # Hien thi so dong, so cot cua diabetes
8
9 # Tach X, y
10
11 # cho biet target thuoc mien gia tri vo han hay gioi han
12
13 # Chia du lieu theo ty le 8:2
14
15 # Thiet lap 02 mo hinh LinearRegression va KNeighborsRegressor
16
17 # Huan luyen
18
19 # Danh gia bang mean_squared_error va r2_score
20
21 # Danh gia cac mo hinh biet rang: MSE cang nho cang tot va R2
cang gan 1 cang tot
```