

Scikit Learn

Xây dựng mô hình đơn giản

HK1, Năm học 2022 - 2023

Các loại bài toán máy học cơ bản

Dựa vào bản chất dữ liệu học và tương tác giữa bộ học và môi trường, thì learning được phân loại

- Học có giám sát (Supervised Learning): bộ học được cung cấp đầu vào và đầu ra mong muốn, mục đích học quy luật ánh xạ từ đầu vào và đầu ra. Quy trình học được tiếp tục đến khi mô hình đạt mức độ chính xác mong muốn. Cặp dữ liệu này còn được gọi là (data, label) hay (data, target) 02 loại chính
 - Phân loại (Classification): có label được chia thành một số hữu hạn nhóm
 - Hồi quy (Regression): label có giá trị thực

Các loại bài toán máy học cơ bản

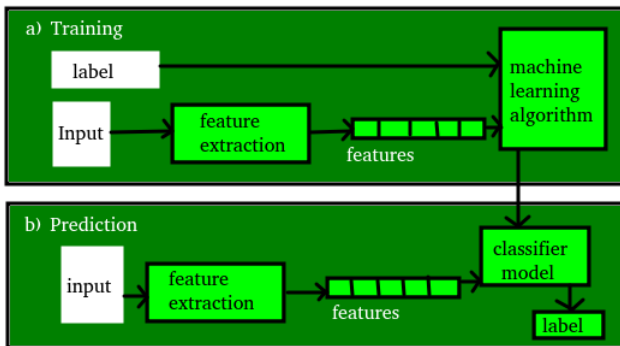
Dựa vào bản chất dữ liệu học và tương tác giữa bộ học và môi trường, thì learning được phân loại

- Học không giám sát (Unsupervised Learning): không biết đầu ra/nhãn mà chỉ có dữ liệu vào. Thuật toán sẽ dựa vào cấu trúc dữ liệu để thực hiện một công việc nào đó như phân cụm, giảm chiều dữ liệu...
- Học bán giám sát (Semi-supervised Learning): một phần dữ liệu được dán nhãn và một phần dữ liệu không có nhãn

Thuật ngữ cơ bản

- Mô hình (model) đại diện cụ thể được học từ dữ liệu bằng cách áp dụng giải thuật học máy và được gọi là giả thuyết (hypothesis)
- Thuộc tính (feature): thuộc tính có thể đo lường được của dữ liệu (data). Tập hợp các features được một tả qua một vector thuộc tính (feature vector). Các feature vector sử dụng như input của model.
- Nhãn (Target/Label): giá trị được dự đoán bởi mô hình
- Huấn luyện (Training): Ý tưởng là đưa một bộ input (features) và đầu ra dự kiến (expected labels). Sau khi huấn luyện, một mô hình sẽ được thiết lập và sử dụng nó để dự đoán đầu ra cho dữ liệu mới cùng loại với dữ liệu huấn luyện
- Dự đoán (Prediction): khi mô hình đã sẵn sàng, một tập dữ liệu đầu vào mới sẽ được đưa vào để dự đoán output.

Thuật ngữ cơ bản



Load dữ liệu

Dataset là một tập dữ liệu gồm 2 phần chính

- Features (inputs, attributes): biến dữ liệu, biểu diễn bằng một ma trận (feature matrix). Danh sách tên thuộc tính được gọi là **feature names**
- Response (target, label, output): biến đầu ra phụ thuộc vào biến thuộc tính, thường là một cột response biểu diễn bằng một vector. Các giá trị có thể cho response vector được gọi là **target names**

Tải dữ liệu dựng sẵn trong sklearn

```
1 # load dataset
2 from sklearn.datasets import load_wine, load_boston, load_digits
3 wines = load_wine();
4 # trích feature matrix (X) và response vector (y)
5 X = wines.data
6 y = wines.target
7
8 # trích ten thuộc tính và ten dịch
9 feature_names = wines.feature_names
10 target_names = wines.target_names
11 print("Feature name: ", feature_names)
12 print("Target name: ", target_names)
13 print("Kích thước của X", X.shape)
14 print("Kích thước của y", y.shape)
15 print("In 5 dòng đầu tiên", X[:5,:])
16
17 # Không phải dataset nào cũng có feature name hay target name
18 boston = load_boston()
19 print(boston.target_names) # sẽ báo lỗi
20 digits = load_digits()
21 print(digits.feature_names) # báo lỗi
```

Tải dữ liệu từ file csv

- Sử dụng thư viện **pandas**
- 02 kiểu dữ liệu quan trọng
 - Series: mảng dán nhãn một chiều có khả năng chứa bất kỳ loại dữ liệu nào.
 - DataFrame: cấu trúc dữ liệu đánh nhãn 2 chiều với các cột có thể các loại dữ liệu khác nhau.

Tải dữ liệu từ file csv

```
1 import pandas as pd
2 #Load tu file csv
3 data = pd.read_csv('weather.csv')
4 # In kích thước
5 print("Kích thước ", data.shape)
6 # in tên thuộc tính
7 print("Features: ", data.columns)
8 # lưu dữ liệu vào ma trận feature và response vector
9 X = data[data.columns[:-1]]
10 y = data[data.columns[-1]]
11 # In 05 dòng đầu của feature matrix
12 print(X.head())
13 # In 05 giá trị đầu của response vector
14 print(y.head())
```

Chia tập dữ liệu

Ở đây chúng ta sẽ bỏ qua bước tiền xử lý dữ liệu và dữ liệu lựa chọn là phù hợp.

Một khía cạnh quan trọng của các mô hình máy là xác định độ chính xác. Để làm được điều này thì mô hình cần được huấn luyện (train) sử dụng dữ liệu và sau đó dự đoán mô hình dự đoán giá trị target với cùng tập dữ liệu. Tuy nhiên phương pháp sẽ có sai sót.

Một cách khắc phục: chia tập dữ liệu thành 02 phần (01 phần dành cho huấn luyện mô hình và phần khác dành cho kiểm tra mô hình)

Chia tập dữ liệu

```
1 from sklearn.datasets import load_iris
2 iris = load_iris()
3 X = iris.data
4 y = iris.target
5
6 # package dùng chỉ dữ liệu
7 from sklearn.model_selection import train_test_split
8 # Chia tập dữ liệu thành train và test
9 # với train chiếm 60% dữ liệu, random_state là seed
   của bộ sinh số random
10 X_train, X_test, y_train, y_test = train_test_split(X
   , y, test_size=0.4, random_state = 1)
11
12 print(X_train.shape)
13 print(X_test.shape)
14 print(y_train.shape)
15 print(y_test.shape)
```

Huấn luyện mô hình

Scikit-learn cung cấp nhiều giải thuật có giao diện nhất quán để học, dự đoán độ chính xác.

Nếu sử dụng giải thuật K láng giềng (KNN - K nearest neighbors) để phân loại

```
1 # import thư viện
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn import metrics
6 # Chuẩn bị dữ liệu
7 # load dữ liệu vào biến iris
8 iris = load_iris()
9 # lưu trữ feature matrix vào biến X
10 X = iris.data
11 # lưu trữ target vector vào biến y
12 y = iris.target
13 # Chia tập dữ liệu vào tập training và test
14 X_train, X_test, y_train, y_test = train_test_split(X, y,
    test_size=0.4, random_state=123)
```

Huấn luyện mô hình

```
1 # Thiet lap mo hinh
2 knn = KNeighborsClassifier(n_neighbors=1)
3 # Qua trinh hoc
4 knn.fit(X_train, y_train)
5 # Qua trinh du doan tren tap test
6 y_pred = knn.predict(X_test)
7 # Danh gia
8 acc = metrics.accuracy_score(y_test, y_pred)
9 print('kNN model accuracy: ', acc)
10 # Du doan cho du lieu moi
11 sample = [[3.2, 5.3, 4.1, 2.8], [2.1,3.2,1.9,0.5]]
12 preds = knn.predict(sample)
13 # Chuyen tu gia tri cua target sang ten that su cua
    label
14 pred_species = [iris.target_names[p] for p in preds]
15 print('Predictions: ', pred_species)
```

Yêu cầu thực hiện

- Thay đổi cách chia dữ liệu: train - test được chia theo tỷ lệ 7:3 với `random_state = 123`
- Vẫn sử dụng KNN nhưng xây dựng các mô hình với giá trị neighbors lần lượt 1, 3, 5
- Dựa vào độ chính xác: giá trị neighbors nào thì độ chính xác cao hơn

Thay đổi model. Sử dụng model cây quyết định (Decision Tree) khác:

- Để sử dụng mô hình cần import: `from sklearn.tree import DecisionTreeClassifier`
- Thay mô hình KNN bằng mô hình: `DecisionTreeClassifier(max_depth=5)`
- Thay đổi giá trị `max_depth` thành 1, 3, 7 và xem thử giá trị nào nên lựa chọn