



coursera-  
Applied...

Qian-Han / coursera-Applied-Data-Science-with-Python

Join GitHub today

GitHub is home to over 40 million developers working together to host and review code, manage projects, and build software together.

Sign up

Branch: master **coursera-Applied-Data-Science-with-Python / Applied-Machine-Learning-in-Python / week4 / Module 4.ipynb** Find file Copy path

Qian-Han 5166es and Assignments 482c258 on 14 Jul 2017 1 contributor

778 Lines (778 tloc) 23.6 KB

You are currently looking at **version 1.0** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ](#) course resource.

## Applied Machine Learning: Module 4 (Supervised Learning, Part II)

### Preamble and Datasets

```
In [ ]: %matplotlib notebook
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.datasets import make_classification, make_blobs
from matplotlib.colors import ListedColormap
from sklearn.datasets import load_breast_cancer
from adspy_shared_utilities import load_crime_dataset

cmap_bold = ListedColormap(['#FFD700', '#0000FF', '#000000'])

# fruits dataset
fruits = pd.read_table('fruit_data_with_colors.txt')

feature_names_fruits = ['height', 'width', 'mass', 'color_score']
X_fruits = fruits[feature_names_fruits]
y_fruits = fruits['fruit_label']
target_names_fruits = ['apple', 'mandarin', 'orange', 'lemon']

X_fruits_2d = fruits[['height', 'width']]
y_fruits_2d = fruits['fruit_label']
```

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 1 页 (共 7 页)

```
# synthetic dataset for simple regression
from sklearn.datasets import make_regression
plt.figure()
plt.title('Sample regression problem with one input variable')
X_R1, y_R1 = make_regression(n_samples=100, n_features=1,
                             n_informative=1, bias=150.0,
                             noise=30, random_state=0)

plt.scatter(X_R1, y_R1, marker='o', s=50)
plt.show()

# synthetic dataset for more complex regression
from sklearn.datasets import make_friedman1
plt.figure()
plt.title('Complex regression problem with one input variable')
X_F1, y_F1 = make_friedman1(n_samples=100, n_features=7,
                             random_state=0)

plt.scatter(X_F1[:, 2], y_F1, marker='o', s=50)
plt.show()

# synthetic dataset for classification (binary)
plt.figure()
plt.title('Sample binary classification problem with two informative features')
X_C2, y_C2 = make_classification(n_samples=100, n_features=2,
                                n_redundant=0, n_informative=2,
                                n_clusters_per_class=1, flip_y=0.1,
                                class_sep=0.5, random_state=0)

plt.scatter(X_C2[:, 0], X_C2[:, 1], marker='o',
            c=y_C2, s=50, cmap=cmap_bold)
plt.show()

# more difficult synthetic dataset for classification (binary)
# with classes that are not linearly separable
X_D2, y_D2 = make_blobs(n_samples=100, n_features=2,
                        centers=8, cluster_std=1.3,
                        random_state=4)

y_D2 = y_D2 % 2
plt.figure()
plt.title('Sample binary classification problem with non-linearly separable classes')
plt.scatter(X_D2[:, 0], X_D2[:, 1], c=y_D2,
            marker='o', s=50, cmap=cmap_bold)
plt.show()

# Breast cancer dataset for classification
cancer = load_breast_cancer()
(X_cancer, y_cancer) = load_breast_cancer(return_X_y=True)

# Communities and Crime dataset
(X_crime, y_crime) = load_crime_dataset()
```

### Naive Bayes classifiers

```
In [ ]: from sklearn.naive_bayes import GaussianNB
from adspy_shared_utilities import plot_class_regions_for_classifier

X_train, X_test, y_train, y_test = train_test_split(X_C2, y_C2, random_state=0)
nbclf = GaussianNB().fit(X_train, y_train)
plot_class_regions_for_classifier(nbclf, X_train, y_train, X_test, y_test,
                                title='Gaussian Naive Bayes classifier: Dataset 1')
```

```
In [ ]: X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2,
                                                         random_state=0)

nbclf = GaussianNB().fit(X_train, y_train)
plot_class_regions_for_classifier(nbclf, X_train, y_train, X_test, y_test,
                                title='Gaussian Naive Bayes classifier: Dataset 2')
```

### Application to a real-world dataset

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 2 页 (共 7 页)

```
In [ ]: X_train, X_test, y_train, y_test = train_test_split(X_cancer, y_cancer, random_state=0)

nbclf = GaussianNB().fit(X_train, y_train)
print('Breast cancer dataset')
print('Accuracy of GaussianNB classifier on training set: {:.2f}'
      .format(nbclf.score(X_train, y_train)))
print('Accuracy of GaussianNB classifier on test set: {:.2f}'
      .format(nbclf.score(X_test, y_test)))
```

### Ensembles of Decision Trees

#### Random forests

```
In [ ]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from adspy_shared_utilities import plot_class_regions_for_classifier_subplot

X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2,
                                                     random_state=0)

fig, subaxes = plt.subplots(1, 1, figsize=(6, 6))

clf = RandomForestClassifier().fit(X_train, y_train)
title = 'Random Forest Classifier, complex binary dataset, default settings'
plot_class_regions_for_classifier_subplot(clf, X_train, y_train, X_test,
                                         y_test, title, subaxes)

plt.show()
```

```
In [ ]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from adspy_shared_utilities import plot_class_regions_for_classifier_subplot

X_train, X_test, y_train, y_test = train_test_split(X_fruits.as_matrix(),
                                                    y_fruits.as_matrix(),
                                                    random_state=0)

fig, subaxes = plt.subplots(6, 1, figsize=(6, 32))

title = 'Random Forest, fruits dataset, default settings'
pair_list = [[0,1], [0,2], [0,3], [1,2], [1,3], [2,3]]

for pair, axis in zip(pair_list, subaxes):
    X = X_train[:, pair]
    y = y_train

    clf = RandomForestClassifier().fit(X, y)
    plot_class_regions_for_classifier_subplot(clf, X, y, None,
                                             None, title, axis,
                                             target_names_fruits)

    axis.set_xlabel(feature_names_fruits[pair[0]])
    axis.set_ylabel(feature_names_fruits[pair[1]])

plt.tight_layout()
plt.show()

clf = RandomForestClassifier(n_estimators=10,
                            random_state=0).fit(X_train, y_train)

print('Random Forest, Fruit dataset, default settings')
print('Accuracy of RF classifier on training set: {:.2f}'
      .format(clf.score(X_train, y_train)))
print('Accuracy of RF classifier on test set: {:.2f}'
      .format(clf.score(X_test, y_test)))
```

### Random Forests on a real-world dataset

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 3 页 (共 7 页)

```
In [ ]: from sklearn.ensemble import RandomForestClassifier

X_train, X_test, y_train, y_test = train_test_split(X_cancer, y_cancer, random_state=0)

clf = RandomForestClassifier(n_estimators=100,
                             random_state=0)
clf.fit(X_train, y_train)

print('Breast cancer dataset')
print('Accuracy of RF classifier on training set: {:.2f}'
      .format(clf.score(X_train, y_train)))
print('Accuracy of RF classifier on test set: {:.2f}'
      .format(clf.score(X_test, y_test)))
```

### Gradient-boosted decision trees

```
In [ ]: from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import train_test_split
from adspy_shared_utilities import plot_class_regions_for_classifier_subplot

X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2, random_state=0)

fig, subaxes = plt.subplots(1, 1, figsize=(6, 6))

clf = GradientBoostingClassifier().fit(X_train, y_train)
title = 'GBDT, complex binary dataset, default settings'
plot_class_regions_for_classifier_subplot(clf, X_train, y_train, X_test,
                                         y_test, title, subaxes)

plt.show()
```

### Gradient boosted decision trees on the fruit dataset

```
In [ ]: X_train, X_test, y_train, y_test = train_test_split(X_fruits.as_matrix(),
                                                          y_fruits.as_matrix(),
                                                          random_state=0)

fig, subaxes = plt.subplots(6, 1, figsize=(6, 32))

pair_list = [[0,1], [0,2], [0,3], [1,2], [1,3], [2,3]]

for pair, axis in zip(pair_list, subaxes):
    X = X_train[:, pair]
    y = y_train

    clf = GradientBoostingClassifier().fit(X, y)
    plot_class_regions_for_classifier_subplot(clf, X, y, None,
                                             None, title, axis,
                                             target_names_fruits)

    axis.set_xlabel(feature_names_fruits[pair[0]])
    axis.set_ylabel(feature_names_fruits[pair[1]])

plt.tight_layout()
plt.show()

clf = GradientBoostingClassifier().fit(X_train, y_train)

print('GBDT, Fruit dataset, default settings')
print('Accuracy of GBDT classifier on training set: {:.2f}'
      .format(clf.score(X_train, y_train)))
print('Accuracy of GBDT classifier on test set: {:.2f}'
      .format(clf.score(X_test, y_test)))
```

### Gradient-boosted decision trees on a real-world dataset

```
In [ ]: from sklearn.ensemble import GradientBoostingClassifier

X_train, X_test, y_train, y_test = train_test_split(X_cancer, y_cancer, random_state=0)

clf = GradientBoostingClassifier(random_state=0)
clf.fit(X_train, y_train)
```

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 4 页 (共 7 页)

```
print('Breast cancer dataset (learning_rate=0.1, max_depth=3)')
print('Accuracy of GBDT classifier on training set: {:.2f}'
      .format(clf.score(X_train, y_train)))
print('Accuracy of GBDT classifier on test set: {:.2f}'\n'
      .format(clf.score(X_test, y_test)))

clf = GradientBoostingClassifier(learning_rate=0.01, max_depth=2, random_state=0)
clf.fit(X_train, y_train)

print('Breast cancer dataset (learning_rate=0.01, max_depth=2)')
print('Accuracy of GBDT classifier on training set: {:.2f}'
      .format(clf.score(X_train, y_train)))
print('Accuracy of GBDT classifier on test set: {:.2f}'
      .format(clf.score(X_test, y_test)))
```

### Neural networks

#### Activation functions

```
In [ ]: xrange = np.linspace(-2, 2, 200)

plt.figure(figsize=(7, 6))

plt.plot(xrange, np.maximum(xrange, 0), label='relu')
plt.plot(xrange, np.tanh(xrange), label='tanh')
plt.plot(xrange, 1 / (1 + np.exp(-xrange)), label='logistic')
plt.legend()
plt.title('Neural network activation functions')
plt.xlabel('Input value (x)')
plt.ylabel('Activation function output')

plt.show()
```

#### Neural networks: Classification

##### Synthetic dataset 1: single hidden layer

```
In [ ]: from sklearn.neural_network import MLPClassifier
from adspy_shared_utilities import plot_class_regions_for_classifier_subplot

X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2, random_state=0)

fig, subaxes = plt.subplots(3, 1, figsize=(6, 18))

for units, axis in zip([4, 8, 16], subaxes):
    nnclf = MLPClassifier(hidden_layer_sizes=(units, ),
                         solver='lbfgs',
                         random_state=0).fit(X_train, y_train)
    plot_class_regions_for_classifier_subplot(nnclf, X_train, y_train,
                                             X_test, y_test, title, axis)

plt.tight_layout()
```

##### Synthetic dataset 1: two hidden layers (60 hidden units)

```
In [ ]: from adspy_shared_utilities import plot_class_regions_for_classifier

X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2, random_state=0)

nnclf = MLPClassifier(hidden_layer_sizes=(10, 10, 10),
                      solver='lbfgs',
                      random_state=0).fit(X_train, y_train)

plot_class_regions_for_classifier(nnclf, X_train, y_train, X_test, y_test,
                                title='Dataset 1: Neural net classifier, 2 layers, 10/10/10 units')
```

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 5 页 (共 7 页)

### Regularization parameter: alpha

```
In [ ]: X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2, random_state=0)

fig, subaxes = plt.subplots(4, 1, figsize=(6, 23))

for this_alpha, axis in zip([0.01, 0.1, 1.0, 5.0], subaxes):
    nnclf = MLPClassifier(solver='lbfgs',
                         hidden_layer_sizes=(10, 10, 10),
                         activation='tanh',
                         random_state=0).fit(X_train, y_train)

    title = 'Dataset 2: NN classifier, alpha = {:.3f}'.format(this_alpha)

    plot_class_regions_for_classifier_subplot(nnclf, X_train, y_train,
                                             X_test, y_test, title, axis)

plt.tight_layout()
```

### The effect of different choices of activation function

```
In [ ]: X_train, X_test, y_train, y_test = train_test_split(X_D2, y_D2, random_state=0)

fig, subaxes = plt.subplots(3, 1, figsize=(6, 18))

for this_activation, axis in zip(['logistic', 'tanh', 'relu'], subaxes):
    nnclf = MLPClassifier(solver='lbfgs',
                         hidden_layer_sizes=(10, 10, 10),
                         activation=this_activation,
                         random_state=0).fit(X_train, y_train)

    title = 'Dataset 2: NN classifier, 2 layers 10/10, {} \
activation function'.format(this_activation)

    plot_class_regions_for_classifier_subplot(nnclf, X_train, y_train,
                                             X_test, y_test, title, axis)

plt.tight_layout()
```

### Neural networks: Regression

```
In [ ]: from sklearn.neural_network import MLPRegressor
from sklearn.preprocessing import MinMaxScaler

X_train, X_test, y_train, y_test = train_test_split(X_cancer, y_cancer, random_state=0)
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

clf = MLPRegressor(hidden_layer_sizes=(100, 100),
                   solver='lbfgs',
                   random_state=0).fit(X_train_scaled, y_train)

print('Breast cancer dataset')
print('Accuracy of NN classifier on training set: {:.2f}'
      .format(clf.score(X_train_scaled, y_train)))
print('Accuracy of NN classifier on test set: {:.2f}'
      .format(clf.score(X_test_scaled, y_test)))
```

```
In [1]: find . -maxdepth 1 -not -type d
./addresses.csv
./train.csv
./Module 2.ipynb
./Assignment 3.ipynb
./Assignment 4.ipynb
./Assignment 1.ipynb
./test.csv
./CommunityPreprocessedData.txt
./adspy_shared_utilities.py
./Module 3.ipynb
./fruit_data_with_colors.txt
./Assignment 4.ipynb
./Assignment 2.ipynb
./mushrooms.csv
./Classifier Visualization.ipynb
./latimes.csv
./Module 1.ipynb
```

```
In [7]: !python readyonly
```

```
In [8]: !cp ./Module 2.ipynb readyonly/Module 2.ipynb
```

```
In [ ]: !cp target '2.ipynb' is not a directory
```

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 6 页 (共 7 页)

```
scaler = MinMaxScaler()

X_train, X_test, y_train, y_test = train_test_split(X_cancer, y_cancer, random_state=0)
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

clf = MLPClassifier(hidden_layer_sizes=(100, 100),
                   solver='lbfgs',
                   random_state=0).fit(X_train_scaled, y_train)

print('Breast cancer dataset')
print('Accuracy of NN classifier on training set: {:.2f}'
      .format(clf.score(X_train_scaled, y_train)))
print('Accuracy of NN classifier on test set: {:.2f}'
      .format(clf.score(X_test_scaled, y_test)))
```

```
In [1]: find . -maxdepth 1 -not -type d
./addresses.csv
./train.csv
./Module 2.ipynb
./Assignment 3.ipynb
./Assignment 4.ipynb
./Assignment 1.ipynb
./test.csv
./CommunityPreprocessedData.txt
./adspy_shared_utilities.py
./Module 3.ipynb
./fruit_data_with_colors.txt
./Assignment 4.ipynb
./Assignment 2.ipynb
./mushrooms.csv
./Classifier Visualization.ipynb
./latimes.csv
./Module 1.ipynb
```

```
In [7]: !python readyonly
```

```
In [8]: !cp ./Module 2.ipynb readyonly/Module 2.ipynb
```

```
In [ ]: !cp target '2.ipynb' is not a directory
```

https://github.com/Qian-Han/coursera-Applied-Data-Science-...Applied-Machine-Learning-in-Python/week4/Module4/204.ipynb

2020/4/18 下午8:48  
第 7 页 (共 7 页)