

机器学习DAY05

数据集划分

对于分类问题训练集和测试集的划分不应该用整个样本空间的特定百分比作为训练数据，而应该在其每一个类别的样本中抽取特定百分比作为训练数据。sklearn模块提供了数据集划分相关方法，可以方便的划分训练集与测试集数据，使用不同数据集训练或测试模型，达到提高分类可信度。

数据集划分相关API：

```
1 import sklearn.model_selection as ms
2
3 训练输入，测试输入，训练输出，测试输出 = \
4     ms.train_test_split(
5         输入集，输出集，test_size=测试集占比，random_state=随机种子)
6
```

案例：

```
1 import numpy as np
2 import sklearn.model_selection as ms
3 import sklearn.naive_bayes as nb
4 import matplotlib.pyplot as mp
5 data = np.loadtxt('./data/multiple1.txt', unpack=False, dtype='u20',
6 delimiter=',')
7 print(data.shape)
8 x = np.array(data[:, :-1], dtype=float)
9 y = np.array(data[:, -1], dtype=float)
10 # 划分训练集和测试集
11 train_x, test_x, train_y, test_y = \
12     ms.train_test_split(x, y, test_size=0.25, random_state=7)
13 # 朴素贝叶斯分类器
14 model = nb.GaussianNB()
15 # 用训练集训练模型
16 model.fit(train_x, train_y)
17
18 l, r = x[:, 0].min() - 1, x[:, 0].max() + 1
19 b, t = x[:, 1].min() - 1, x[:, 1].max() + 1
20 n = 500
21 grid_x, grid_y = np.meshgrid(np.linspace(l, r, n), np.linspace(b, t, n))
22 samples = np.column_stack((grid_x.ravel(), grid_y.ravel()))
23 grid_z = model.predict(samples)
24 grid_z = grid_z.reshape(grid_x.shape)
25
26 pred_test_y = model.predict(test_x)
27 # 计算并打印预测输出的精确度
28 print((test_y == pred_test_y).sum() / pred_test_y.size)
29
30 mp.figure('Naive Bayes Classification', facecolor='lightgray')
31 mp.title('Naive Bayes Classification', fontsize=20)
32 mp.xlabel('x', fontsize=14)
```

```

32 mp.ylabel('y', fontsize=14)
33 mp.tick_params(labelsize=10)
34 mp.pcolormesh(grid_x, grid_y, grid_z, cmap='gray')
35 mp.scatter(test_x[:,0], test_x[:,1], c=test_y, cmap='brg', s=80)
36 mp.show()

```

交叉验证

由于数据集的划分有不确定性，若随机划分的样本正好处于某类特殊样本，则得到的训练模型所预测的结果的可信度将受到质疑。所以需要进行多次交叉验证，把样本空间中的所有样本均分成n份，使用不同的训练集训练模型，对不同的测试集进行测试时输出指标得分。sklearn提供了交叉验证相关API：

```

1 import sklearn.model_selection as ms
2 指标值数组 = \
3     ms.cross_val_score(模型, 输入集, 输出集, cv=折叠数, scoring=指标名)

```

案例：使用交叉验证，输出分类器的精确度：

```

1 # 划分训练集和测试集
2 train_x, test_x, train_y, test_y = \
3     ms.train_test_split(
4         x, y, test_size=0.25, random_state=7)
5 # 朴素贝叶斯分类器
6 model = nb.GaussianNB()
7 # 交叉验证
8 # 精确度
9 ac = ms.cross_val_score(model, train_x, train_y, cv=5, scoring='accuracy')
10 print(ac.mean())
11 #用训练集训练模型
12 model.fit(train_x, train_y)

```

交叉验证指标

1. 精确度(accuracy)：分类正确的样本数/总样本数
2. 查准率(precision_weighted)：针对每一个类别，预测正确的样本数比上预测出来的样本数
3. 召回率(recall_weighted)：针对每一个类别，预测正确的样本数比上实际存在的样本数
4. f1得分(f1_weighted)：
 - 2x查准率x召回率/(查准率+召回率)

在交叉验证过程中，针对每一次交叉验证，计算所有类别的查准率、召回率或者f1得分，然后取各类别相应指标值的平均数，作为这一次交叉验证的评估指标，然后再将所有交叉验证的评估指标以数组的形式返回调用者。

```

1  # 交叉验证
2  # 精确度
3  ac = ms.cross_val_score( model, train_x, train_y, cv=5, scoring='accuracy')
4  print(ac.mean())
5  # 查准率
6  pw = ms.cross_val_score( model, train_x, train_y, cv=5,
7  scoring='precision_weighted')
8  print(pw.mean())
9  # 召回率
10 rw = ms.cross_val_score( model, train_x, train_y, cv=5,
11 scoring='recall_weighted')
12 print(rw.mean())
13 # f1得分
14 fw = ms.cross_val_score( model, train_x, train_y, cv=5,
15 scoring='f1_weighted')
16 print(fw.mean())

```

混淆矩阵

每一行和每一列分别对应样本输出中的每一个类别，行表示实际类别，列表示预测类别。

	A类别	B类别	C类别
A类别	5	0	0
B类别	0	6	0
C类别	0	0	7

上述矩阵即为理想的混淆矩阵。不理想的混淆矩阵如下：

	A类别	B类别	C类别
A类别	3	1	1
B类别	0	4	2
C类别	0	0	7

查准率 = 主对角线上的值 / 该值所在列的和

召回率 = 主对角线上的值 / 该值所在行的和

获取模型分类结果的混淆矩阵的相关API：

```

1  import sklearn.metrics as sm
2  混淆矩阵 = sm.confusion_matrix(实际输出, 预测输出)

```

案例：输出分类结果的混淆矩阵。

```

1 #输出混淆矩阵并绘制混淆矩阵图谱
2 cm = sm.confusion_matrix(test_y, pred_test_y)
3 print(cm)
4 mp.figure('Confusion Matrix', facecolor='lightgray')
5 mp.title('Confusion Matrix', fontsize=20)
6 mp.xlabel('Predicted Class', fontsize=14)
7 mp.ylabel('True Class', fontsize=14)
8 mp.xticks(np.unique(pred_test_y))
9 mp.yticks(np.unique(test_y))
10 mp.tick_params(labelsize=10)
11 mp.imshow(cm, interpolation='nearest', cmap='jet')
12 mp.show()

```

分类报告

sklearn.metrics提供了分类报告相关API，不仅可以得到混淆矩阵，还可以得到交叉验证查准率、召回率、f1得分的结果，可以方便的分析出哪些样本是异常样本。

```

1 # 获取分类报告
2 cr = sm.classification_report(实际输出, 预测输出)

```

案例：输出分类报告：

```

1 # 获取分类报告
2 cr = sm.classification_report(test_y, pred_test_y)
3 print(cr)

```

决策树分类

决策树分类模型会找到与样本特征匹配的叶子节点然后以投票的方式进行分类。在样本文件中统计了小汽车的常见特征信息及小汽车的分类，使用这些数据基于决策树分类算法训练模型预测小汽车等级。

汽车价格	维修费用	车门数量	载客数	后备箱	安全性	汽车级别

案例：基于决策树分类算法训练模型预测小汽车等级。

1. 读取文本数据，对每列进行标签编码，基于随机森林分类器进行模型训练，进行交叉验证。

```

1 import numpy as np
2 import sklearn.preprocessing as sp
3 import sklearn.ensemble as se
4 import sklearn.model_selection as ms
5
6 data = np.loadtxt('../data/car.txt', delimiter=',', dtype='U10')
7 data = data.T
8 encoders = []
9 train_x, train_y = [], []
10 for row in range(len(data)):
11     encoder = sp.LabelEncoder()
12     if row < len(data) - 1:
13         train_x.append(encoder.fit_transform(data[row]))
14     else:

```

```

15     train_y = encoder.fit_transform(data[row])
16     encoders.append(encoder)
17     train_x = np.array(train_x).T
18     # 随机森林分类器
19     model = se.RandomForestClassifier(max_depth=6, n_estimators=200,
20                                     random_state=7)
21     print(ms.cross_val_score(model, train_x, train_y, cv=4,
22                             scoring='f1_weighted').mean())
21     model.fit(train_x, train_y)

```

1. 自定义测试集，使用已训练的模型对测试集进行测试，输出结果。

```

1  data = [
2      ['high', 'med', '5more', '4', 'big', 'low', 'unacc'],
3      ['high', 'high', '4', '4', 'med', 'med', 'acc'],
4      ['low', 'low', '2', '4', 'small', 'high', 'good'],
5      ['low', 'med', '3', '4', 'med', 'high', 'vgood']]
6
7  data = np.array(data).T
8  test_x, test_y = [], []
9  for row in range(len(data)):
10     encoder = encoders[row]
11     if row < len(data) - 1:
12         test_x.append(encoder.transform(data[row]))
13     else:
14         test_y = encoder.transform(data[row])
15  test_x = np.array(test_x).T
16  pred_test_y = model.predict(test_x)
17  print((pred_test_y == test_y).sum() / pred_test_y.size)
18  print(encoders[-1].inverse_transform(test_y))
19  print(encoders[-1].inverse_transform(pred_test_y))

```

验证曲线

验证曲线：模型性能 = f(超参数)

验证曲线所需API：

```

1  train_scores, test_scores = ms.validation_curve(
2      model,                    # 模型
3      输入集, 输出集,
4      'n_estimators',          # 超参数名
5      np.arange(50, 550, 50),  # 超参数序列
6      cv=5                     # 折叠数
7  )

```

train_scores的结构:

参数取值	第一次cv	第二次cv	第三次cv	第四次cv	第五次cv
50	0.91823444	0.91968162	0.92619392	0.91244573	0.91040462
100	0.91968162	0.91823444	0.91244573	0.92619392	0.91244573
...

test_scores的结构与train_scores的结构相同。

案例：在小汽车评级案例中使用验证曲线选择较优参数。

```
1 # 获得关于n_estimators的验证曲线
2 model = se.RandomForestClassifier(max_depth=6, random_state=7)
3 n_estimators = np.arange(50, 550, 50)
4 train_scores, test_scores = ms.validation_curve(model, train_x, train_y,
5 'n_estimators', n_estimators, cv=5)
6 print(train_scores, test_scores)
7 train_means1 = train_scores.mean(axis=1)
8 for param, score in zip(n_estimators, train_means1):
9     print(param, '->', score)
10
11 mp.figure('n_estimators', facecolor='lightgray')
12 mp.title('n_estimators', fontsize=20)
13 mp.xlabel('n_estimators', fontsize=14)
14 mp.ylabel('F1 Score', fontsize=14)
15 mp.tick_params(labelsize=10)
16 mp.grid(linestyle=':')
17 mp.plot(n_estimators, train_means1, 'o-', c='dodgerblue', label='Training')
18 mp.legend()
19 mp.show()
```

```
1 # 获得关于max_depth的验证曲线
2 model = se.RandomForestClassifier(n_estimators=200, random_state=7)
3 max_depth = np.arange(1, 11)
4 train_scores, test_scores = ms.validation_curve(
5     model, train_x, train_y, 'max_depth', max_depth, cv=5)
6 train_means2 = train_scores.mean(axis=1)
7 for param, score in zip(max_depth, train_means2):
8     print(param, '->', score)
9
10 mp.figure('max_depth', facecolor='lightgray')
11 mp.title('max_depth', fontsize=20)
12 mp.xlabel('max_depth', fontsize=14)
13 mp.ylabel('F1 Score', fontsize=14)
14 mp.tick_params(labelsize=10)
15 mp.grid(linestyle=':')
16 mp.plot(max_depth, train_means2, 'o-', c='dodgerblue', label='Training')
17 mp.legend()
18 mp.show()
```

学习曲线

学习曲线：模型性能 = f(训练集大小)

学习曲线所需API：

```
1 __, train_scores, test_scores = ms.learning_curve(
2     model,          # 模型
3     输入集, 输出集,
4     train_sizes=[0.9, 0.8, 0.7], # 训练集大小序列
5     cv=5            # 折叠数
6 )
```

train_scores的结构:

案例：在小汽车评级案例中使用学习曲线选择训练集大小最优参数。

```
1 # 获得学习曲线
2 model = se.RandomForestClassifier( max_depth=9, n_estimators=200,
   random_state=7)
3 train_sizes = np.linspace(0.1, 1, 10)
4 _, train_scores, test_scores = ms.learning_curve(
5     model, x, y, train_sizes=train_sizes, cv=5)
6 test_means = test_scores.mean(axis=1)
7 for size, score in zip(train_sizes, train_means):
8     print(size, '->', score)
9 mp.figure('Learning Curve', facecolor='lightgray')
10 mp.title('Learning Curve', fontsize=20)
11 mp.xlabel('train_size', fontsize=14)
12 mp.ylabel('F1 Score', fontsize=14)
13 mp.tick_params(labelsize=10)
14 mp.grid(linestyle=':')
15 mp.plot(train_sizes, test_means, 'o-', c='dodgerblue', label='Training')
16 mp.legend()
17 mp.show()
```

案例：预测工人工资收入。

读取adult.txt，针对不同形式的特征选择不同类型的编码器，训练模型，预测工人工资收入。

1. 自定义标签编码器，若为数字字符串，则使用该编码器，保留特征数字值的意义。

```
1 class DigitEncoder():
2
3     def fit_transform(self, y):
4         return y.astype(int)
5
6     def transform(self, y):
7         return y.astype(int)
8
9     def inverse_transform(self, y):
10        return y.astype(str)
```

1. 读取文件，整理样本数据，对样本矩阵中的每一列进行标签编码。

```
1 num_less, num_more, max_each = 0, 0, 7500
2 data = []
3
4 txt = np.loadtxt('../data/adult.txt', dtype='U20', delimiter=', ')
5 for row in txt:
6     if(' ?' in row):
7         continue
8     elif(str(row[-1]) == '<=50K'):
9         num_less += 1
10        data.append(row)
11    elif(str(row[-1]) == '>50K'):
12        num_more += 1
13        data.append(row)
14
15 data = np.array(data).T
16 encoders, x = [], []
```

```

17 for row in range(len(data)):
18     if str(data[row, 0]).isdigit():
19         encoder = DigitEncoder()
20     else:
21         encoder = sp.LabelEncoder()
22     if row < len(data) - 1:
23         x.append(encoder.fit_transform(data[row]))
24     else:
25         y = encoder.fit_transform(data[row])
26     encoders.append(encoder)

```

1. 划分训练集与测试集，基于朴素贝叶斯分类算法构建学习模型，输出交叉验证分数，验证测试集。

```

1 x = np.array(x).T
2 train_x, test_x, train_y, test_y = ms.train_test_split(
3     x, y, test_size=0.25, random_state=5)
4 model = nb.GaussianNB()
5 print(ms.cross_val_score(
6     model, x, y, cv=10, scoring='f1_weighted').mean())
7 model.fit(train_x, train_y)
8 pred_test_y = model.predict(test_x)
9 print((pred_test_y == test_y).sum() / pred_test_y.size)

```

1. 模拟样本数据，预测收入级别。

```

1 data = [['39', 'State-gov', '77516', 'Bachelors',
2         '13', 'Never-married', 'Adm-clerical', 'Not-in-family',
3         'white', 'Male', '2174', '0', '40', 'United-States']]
4 data = np.array(data).T
5 x = []
6 for row in range(len(data)):
7     encoder = encoders[row]
8     x.append(encoder.transform(data[row]))
9 x = np.array(x).T
10 pred_y = model.predict(x)
11 print(encoders[-1].inverse_transform(pred_y))

```

支持向量机(SVM)

支持向量机原理

1. 寻求最优分类边界

正确：对大部分样本可以正确地划分类别。

泛化：最大化支持向量间距。

公平：与支持向量等距。

简单：线性，直线或平面，分割超平面。

2. 基于核函数的升维变换

通过名为核函数的特征变换，增加新的特征，使得低维度空间中的线性不可分问题变为高维度空间中的线性可分问题。

线性核函数：linear，不通过核函数进行维度提升，仅在原始维度空间中寻求线性分类边界。

基于线性核函数的SVM分类相关API：


```

1 model = svm.SVC(kernel='linear')
2 model.fit(train_x, train_y)

```

案例：对simple2.txt中的数据进行分类。

```

1 import numpy as np
2 import sklearn.model_selection as ms
3 import sklearn.svm as svm
4 import sklearn.metrics as sm
5 import matplotlib.pyplot as mp
6 x, y = [], []
7 data = np.loadtxt('../data/multiple2.txt', delimiter=',', dtype='f8')
8 x = data[:, :-1]
9 y = data[:, -1]
10 train_x, test_x, train_y, test_y = \
11     ms.train_test_split(x, y, test_size=0.25, random_state=5)
12 # 基于线性核函数的支持向量机分类器
13 model = svm.SVC(kernel='linear')
14 model.fit(train_x, train_y)
15 n = 500
16 l, r = x[:, 0].min() - 1, x[:, 0].max() + 1
17 b, t = x[:, 1].min() - 1, x[:, 1].max() + 1
18 grid_x = np.meshgrid(np.linspace(l, r, n),
19                     np.linspace(b, t, n))
20 flat_x = np.column_stack((grid_x[0].ravel(), grid_x[1].ravel()))
21 flat_y = model.predict(flat_x)
22 grid_y = flat_y.reshape(grid_x[0].shape)
23 pred_test_y = model.predict(test_x)
24 cr = sm.classification_report(test_y, pred_test_y)
25 print(cr)
26 mp.figure('SVM Linear Classification', facecolor='lightgray')
27 mp.title('SVM Linear Classification', fontsize=20)
28 mp.xlabel('x', fontsize=14)
29 mp.ylabel('y', fontsize=14)
30 mp.tick_params(labelsize=10)
31 mp.pcolormesh(grid_x[0], grid_x[1], grid_y, cmap='gray')
32 mp.scatter(test_x[:, 0], test_x[:, 1], c=test_y, cmap='brg', s=80)
33 mp.show()

```

多项式核函数：poly · 通过多项式函数增加原始样本特征的高次方幂

$$y = x_1 + x_2$$

$$y = x_1^2 + 2x_1x_2 + x_2^2$$

$$y = x_1^3 + 3x_1^2x_2 + 3x_1x_2^2 + x_2^3$$

案例 · 基于多项式核函数训练sample2.txt中的样本数据。

```

1 # 基于线性核函数的支持向量机分类器
2 model = svm.SVC(kernel='poly', degree=3)
3 model.fit(train_x, train_y)

```

径向基核函数：rbf · 通过高斯分布函数增加原始样本特征的分布概率

案例 · 基于径向基核函数训练sample2.txt中的样本数据。

```
1 # 基于径向基核函数的支持向量机分类器
2 # C：正则强度
3 # gamma：正态分布曲线的标准差
4 model = svm.SVC(kernel='rbf', C=600, gamma=0.01)
5 model.fit(train_x, train_y)
```