

原

Python sklearn数据分析中常用方法

2017年08月01日 15:50:15

青盏

阅读数：3175

标签：

python

数据分析

更多

一、数据处理

随机划分训练集和测试集：

```
1 from sklearn.model_selection import train_test_split
2
3 X_all = data_train.drop(['Survived', 'PassengerId'], axis=1) # 只包含特征集，不包含预测目标
4 y_all = data_train['Survived'] # 只包含预测目标
5
6 num_test = 0.20 # 测试集占据比例，如果是整数的话就是样本的数量
7 # 注意返回值: (X_train,y_train) 训练集的特征和Label || (X_test,y_test) 训练集的特征和Label
8 X_train, X_test, y_train, y_test = train_test_split(X_all, y_all, test_size=num_test, random_state=23)
9 # random_state 参数表示随机种子，如果为0或不填，每次随机产生的随机数组不同。

1 from sklearn.model_selection import StratifiedShuffleSplit
2 sss = StratifiedShuffleSplit(n_splits=10, test_size=0.1, random_state=0)
3 # sss 对象用于划分数据集
4 X = train[0::, 1::]
5 # X 为特征集
6 y = train[0::, 0]
7 # y 为Label 集
8 for train_index, test_index in sss.split(X, y):
9     X_train, X_test = X[train_index], X[test_index]
10    y_train, y_test = y[train_index], y[test_index]
```

文本特征提取：

sklearn.feature_extraction.text 文本相关的特征抽取

```
1 text.CountVectorizer: 将文本转换为每个词出现的个数的向量
2 text.TfidfVectorizer: 将文本转换为tfidf值的向量
3 text.HashingVectorizer: 文本的特征哈希

1 # CountVectorizer
2 ar1 = '今天 今天 天气 不错 我们 愉快 玩耍'
3 ar2 = '今天 锻炼 舒服 天气 一般'
4 ar3 = '天气 糟糕'
5 text = [ar1,ar2,ar3]
6 from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer
7 ct = CountVectorizer()
8 print(ct.fit_transform(text).todense())
9 print(ct.vocabulary_)
```

```
[[0 1 2 1 1 1 1 0 0 0]
 [1 0 1 1 0 0 0 0 1 1]
 [0 0 0 1 0 0 0 1 0 0]]
{'今天': 2, '天气': 3, '不错': 1, '我们': 5, '愉快': 4, '玩耍': 6, '锻炼': 1, '舒服': 1, '一般': 1, '糟糕': 1}
```

文本特征向量化，其实就是将所有文本中出现的单词组成一个词典，这个词典可以作为一个向量，对每个样例中出现的次数进行统计，从而每个样例度量。如上图。但如果只是统计词频是不够的，因为常用的语言中有些单词出现频率特别高，但是却没有啥意义。如“the、to”。因此我们需要降低这类单

TF-IDF思想是一个词语在一篇文章中出现次数越多，同时在所有文档中出现次数越少，越能够代表该文章。

```
1 # TfidfVectorizer
2 from sklearn.feature_extraction.text import TfidfTransformer
3 from sklearn.feature_extraction.text import CountVectorizer
4 transformer = TfidfTransformer()
```

```

5 tfidf = transformer.fit_transform(ct.fit_transform(text))
6 print(tfidf.todense())

1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf2 = TfidfVectorizer()
3 re = tfidf2.fit_transform(text)
4 print(re.todense())

```

```

[[ 0.         0.3874216  0.58928822  0.22881744  0.3874216  0.3874216
  0.3874216  0.         0.         0.         ]
 [ 0.50461134  0.         0.38376993  0.29803159  0.         0.         0.
  0.         0.50461134  0.50461134]
 [ 0.         0.         0.         0.50854232  0.         0.         0.
  0.861037  0.         0.         ]]

```

如果词库很大时，生成的词向量维度过大，可以使用hash方法对其进行降维.hash后的词向量就无法解释其意义。

```

1 from sklearn.feature_extraction.text import HashingVectorizer
2 vectorizer2=HashingVectorizer(n_features = 6,norm = None)
3 print(vectorizer2.fit_transform(text).todense())

```

```

[[ 0. -1. -1.  1.  0. -2.]
 [ 0.  1.  1.  1.  0.  0.]
 [ 0.  0.  0.  1. -1.  0.]]

```

二、模型选择

```

1 # machine Learning
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.svm import SVC, LinearSVC
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn.naive_bayes import GaussianNB
7 from sklearn.linear_model import Perceptron
8 from sklearn.linear_model import SGDClassifier
9 from sklearn.tree import DecisionTreeClassifier

```

逻辑回归:

```

1 logreg = LogisticRegression()
2 logreg.fit(X_train, Y_train)
3 Y_pred = logreg.predict(X_test)
4 acc_log = round(logreg.score(X_train, Y_train) * 100, 2)
5 acc_log

1 # 查看特征系数
2 coeff_df = pd.DataFrame(train_df.columns.delete(0))
3 coeff_df.columns = ['Feature']
4 coeff_df["Correlation"] = pd.Series(logreg.coef_[0])
5
6 coeff_df.sort_values(by='Correlation', ascending=False)

```

SVC支持向量机:

```

1 svc = SVC()
2 svc.fit(X_train, Y_train)
3 Y_pred = svc.predict(X_test)
4 acc_svc = round(svc.score(X_train, Y_train) * 100, 2)

```

```
1  ## Linear SVC
2  # linear_svc = LinearSVC()
3  # linear_svc.fit(X_train, Y_train)
4  # Y_pred = linear_svc.predict(X_test)
5  # acc_linear_svc = round(linear_svc.score(X_train, Y_train) * 100, 2)
6  # acc_linear_svc
```

K近邻学习KNN:

```
1  # knn = KNeighborsClassifier(n_neighbors = 3)
2  # knn.fit(X_train, Y_train)
3  # Y_pred = knn.predict(X_test)
4  # acc_knn = round(knn.score(X_train, Y_train) * 100, 2)
5  # acc_knn
```

朴素贝叶斯分类器:

```
1  # gaussian = GaussianNB()
2  # gaussian.fit(X_train, Y_train)
3  # Y_pred = gaussian.predict(X_test)
4  # acc_gaussian = round(gaussian.score(X_train, Y_train) * 100, 2)
5  # acc_gaussian
```

感知机:

```
1  # perceptron = Perceptron()
2  # perceptron.fit(X_train, Y_train)
3  # Y_pred = perceptron.predict(X_test)
4  # acc_perceptron = round(perceptron.score(X_train, Y_train) * 100, 2)
5  # acc_perceptron
```

随机梯度下降法:

```
1  # sgd = SGDClassifier()
2  # sgd.fit(X_train, Y_train)
3  # Y_pred = sgd.predict(X_test)
4  # acc_sgd = round(sgd.score(X_train, Y_train) * 100, 2)
5  # acc_sgd
```

决策树:

```
1  ## Decision Tree
2  # decision_tree = DecisionTreeClassifier()
3  # decision_tree.fit(X_train, Y_train)
4  # Y_pred = decision_tree.predict(X_test)
5  # acc_decision_tree = round(decision_tree.score(X_train, Y_train) * 100, 2)
6  # acc_decision_tree
```

随机森林:

```
1  ## Random Forest
2  # random_forest = RandomForestClassifier(n_estimators=100)
3  # random_forest.fit(X_train, Y_train)
4  # Y_pred = random_forest.predict(X_test)
5  # random_forest.score(X_train, Y_train)
6  # acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)
7  # acc_random_forest
```

```
1  # 基于准确率搜索最佳参数的随机森林
2  from sklearn.ensemble import RandomForestClassifier
3  from sklearn.metrics import make_scorer, accuracy_score
4  from sklearn.model_selection import GridSearchCV
5
6  # Choose the type of classifier.
7  clf = RandomForestClassifier()
8
```

```

9  # Choose some parameter combinations to try
10 parameters = {'n_estimators': [4, 6, 9],
11               'max_features': ['log2', 'sqrt', 'auto'],
12               'criterion': ['entropy', 'gini'],
13               'max_depth': [2, 3, 5, 10],
14               'min_samples_split': [2, 3, 5],
15               'min_samples_leaf': [1, 5, 8]
16               }
17
18 # Type of scoring used to compare parameter combinations
19 acc_scorer = make_scorer(accuracy_score)
20
21 # Run the grid search
22 grid_obj = GridSearchCV(clf, parameters, scoring=acc_scorer)
23 grid_obj = grid_obj.fit(X_train, y_train)
24
25 # Set the clf to the best combination of parameters
26 clf = grid_obj.best_estimator_
27
28 # Fit the best algorithm to the data.
29 clf.fit(X_train, y_train)

```

遍历模型方法：

```

1  import matplotlib.pyplot as plt
2  import seaborn as sns
3
4  from sklearn.model_selection import StratifiedShuffleSplit
5  from sklearn.metrics import accuracy_score, log_loss
6  from sklearn.neighbors import KNeighborsClassifier
7  from sklearn.svm import SVC
8  from sklearn.tree import DecisionTreeClassifier
9  from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
10 from sklearn.naive_bayes import GaussianNB
11 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis
12 from sklearn.linear_model import LogisticRegression
13
14 classifiers = [
15     KNeighborsClassifier(3),
16     SVC(probability=True),
17     DecisionTreeClassifier(),
18     RandomForestClassifier(),
19     AdaBoostClassifier(),
20     GradientBoostingClassifier(),
21     GaussianNB(),
22     LinearDiscriminantAnalysis(),
23     QuadraticDiscriminantAnalysis(),
24     LogisticRegression()]
25
26 log_cols = ["Classifier", "Accuracy"]
27 log = pd.DataFrame(columns=log_cols)
28
29 sss = StratifiedShuffleSplit(n_splits=10, test_size=0.1, random_state=0)
30 # sss 对象用于划分数据集
31 X = train[0::, 1::]
32 # X为特征集
33 y = train[0::, 0]
34 # y为Label集
35
36 acc_dict = {}
37
38 for train_index, test_index in sss.split(X, y):
39     X_train, X_test = X[train_index], X[test_index]
40     y_train, y_test = y[train_index], y[test_index]
41
42     for clf in classifiers:
43         name = clf.__class__.__name__
44         clf.fit(X_train, y_train)
45         train_predictions = clf.predict(X_test)
46         acc = accuracy_score(y_test, train_predictions)

```

```
47         if name in acc_dict:
48             acc_dict[name] += acc
49         else:
50             acc_dict[name] = acc
51
52     for clf in acc_dict:
53         acc_dict[clf] = acc_dict[clf] / 10.0
54         # 计算平均准确率
55     log_entry = pd.DataFrame([[clf, acc_dict[clf]]], columns=log_cols)
56     log = log.append(log_entry)
57
58     plt.xlabel('Accuracy')
59     plt.title('Classifier Accuracy')
60
61     sns.set_color_codes("muted")
62     sns.barplot(x='Accuracy', y='Classifier', data=log, color="b")
63     # 画条形图分析
```

叠加多层 (2) 模型教程

三、模型评估

使用k折交叉验证法:

```
1  from sklearn.cross_validation import KFold
2
3  def run_kfold(clf):
4      kf = KFold(891, n_folds=10)
5      outcomes = []
6      fold = 0
7      for train_index, test_index in kf:
8          fold += 1
9          X_train, X_test = X_all.values[train_index], X_all.values[test_index]
10         y_train, y_test = y_all.values[train_index], y_all.values[test_index]
11         clf.fit(X_train, y_train)
12         predictions = clf.predict(X_test)
13         accuracy = accuracy_score(y_test, predictions)
14         outcomes.append(accuracy)
15         print("Fold {0} accuracy: {1}".format(fold, accuracy))
16     mean_outcome = np.mean(outcomes)
17     print("Mean Accuracy: {0}".format(mean_outcome))
18     run_kfold(clf)
```

四、其他

保存模型:

```
1  Pickle:
2  >>> import pickle
3  >>> s = pickle.dumps(clf)
4  >>> clf2 = pickle.loads(s)
5  >>> clf2.predict(X[0:1])
6  joblib:
7  >>> from sklearn.externals import joblib
8  >>> joblib.dump(clf, 'filename.pkl')
9  >>> clf = joblib.load('filename.pkl')
```

Python全栈学完需要多少钱?

零基础学爬虫, 你要掌握学习那些技能? 需要学多久?



想对作者说点什么

SVC算法预测汽车是否值得购买 - 梁某

汽车是否值得购买 这是一个关于汽车测评的数据集, 类别变量为汽车的测评: (unacc, ACC, good, vgood) 分...

👁 135

来自: 梁某