

SVM机器学习实验——171250574杨逸存

实验简介

- 本实验使用python语言及sklearn、numpy、pandas等相关库实现了SVN分类算法在辨别垃圾邮件上的应用问题。

算法/代码详解

SVM算法原理

SVM算法是建立在感知机分类的基础上的。感知机意图寻找一个超平面 $\Pi: w^T x + b = 0$ 来区分两个类，但是SVM意图寻找一块最有的平面，即**让决策面离正负样本点的间隔都尽可能大**。对于一个数据集到超平面的距离表示为：

$$d(D, \Pi) = \min d((x, y), \Pi), \text{ 其中 } d((x, y), \Pi) = 1 / |w| \cdot y \cdot (w^T x + b)$$

则优化目标为：

$$\max d(D, \Pi), \text{ w.l.o.g. } \min (|w|^2 / 2) \text{ s.t. } y_i(w^T x_i + b) \geq 1$$

对于非线性可分的数据集，我们需要放松限制条件，同时增加损失函数的惩罚度：

$$\min [|w|^2 / 2 + C \cdot \sum_i \xi_i] \text{ s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i$$

经过数学转化：

$$\min L(D) = \min [|w|^2 / 2 + C \cdot \sum_i \text{ReLU}(1 - y_i(w^T x_i + b))] \text{ s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i$$

对改进的损失函数进行梯度下降，先求导：

$$L(D) = \min [|w|^2 / 2 + C \cdot \sum_i \text{ReLU}(1 - y_i(w^T x_i + b))] \text{ i.e. } L(x, y) = \min [|w|^2 / 2 + C \cdot \text{ReLU}(1 - y(w^T x + b))]$$

$$\text{当 } y(w^T x + b) \geq 1 \text{ 时, } \partial L(x, y) / \partial w = w, \partial L(x, y) / \partial b = 0$$

$$\text{当 } y(w^T x + b) < 1 \text{ 时, } \partial L(x, y) / \partial w = w - Cyx, \partial L(x, y) / \partial b = -Cy$$

i.e. 梯度下降表示为： $w \leftarrow w(1 - \eta)$ ， η 为学习率

当有 $y(w^T x + b) < 1$ 时，选出某个被错分类的样本实例 (x, y) ，再进行如下操作：

$$w \leftarrow w + \eta Cyx$$

$$b \leftarrow b + \eta Cy$$

最终完成收敛。

核心代码讲解

```

1 mail_matrix, mail_labels = read_files(mail_dir_path)
2 # 分割训练测试集，按照8: 2的比例分割训练集和测试集
3 mail_train, mail_test, mail_train_label, mail_test_label \
4     = train_test_split(mail_matrix, mail_labels, test_size=0.2,
5                           random_state=1)
6 # 使用TF-IDF将文本特征提取为向量
7 # 关键词document frequency最大阈值设为0.6
8 # 对关键词出现次数进行实际计数，而非0/1
9 count_vec = CountVectorizer(stop_words='english', max_df=0.6,
10                             decode_error='ignore', binary=False)
11 count_train = count_vec.fit_transform(mail_train)
12 tfidfTransformer = TfidfTransformer()
13 tfidf_train = tfidfTransformer.fit_transform(count_train)
14 # 训练模型
15 svm_model = LinearSVC()
16 svm_model.fit(tfidf_train, mail_train_label)

```

对测试集做同样处理后，进行predict操作，最终混淆矩阵和度量值代码如下：

```

1 '''
2 混淆矩阵：      预测正例(1)      预测反例(0)
3 真实正例(1)      TP              FN
4 真实反例(0)      FP              TN
5 设矩阵的labels字段为[1, 0]定义顺序
6 '''
7 confusion_m = pd.DataFrame(confusion_matrix(mail_test_label,
8                                               predict_label, labels=[1, 0]),
9                             index=['actual_ham', 'actual_spam'],
10                             columns=['predicted_ham',
11                                       'predicted_spam'])
12 print(confusion_m)
13 print('精确率Precision=TP/(TP+FP)，预测为正类的实例中真正预测正确的比例：', precision_score(mail_test_label, predict_label, pos_label=1))
14 print('准确率Accuracy=(TP+TN)/(TP+FN+FP+TN)，总体预测结果正确率：', accuracy_score(mail_test_label, predict_label))
15 print('召回率Recall=TP/(TP+FN)，指对于所有正类实例的预测正确率：', recall_score(mail_test_label, predict_label, pos_label=1))

```

实验结果分析

混淆矩阵及三个度量值：精确率、准确率、召回率显示如下

```

      predicted_ham  predicted_spam
actual_ham      3299           60
actual_spam      16      3369
精确率Precision=TP/(TP+FP)，预测为正类的实例中真正预测正确的比例： 0.9951734539969834
准确率Accuracy=(TP+TN)/(TP+FN+FP+TN)，总体预测结果正确率： 0.9887307236061684
召回率Recall=TP/(TP+FN)，指对于所有正类实例的预测正确率： 0.982137540934802

Process finished with exit code 0

```

binary设为True后实验结果，各方面度量值都有所下降，因此关键词出现的频数也对模型训练有所影响。

```

        predicted_ham predicted_spam
actual_ham      3293         66
actual_spam       22      3363
精确率Precision=TP/(TP+FP), 预测为正类的实例中真正预测正确的比例:  0.9933634992458522
准确率Accuracy=(TP+TN)/(TP+FN+FP+TN), 总体预测结果正确率:  0.9869513641755635
召回率Recall=TP/(TP+FN), 指对于所有正类实例的预测正确率:  0.9803512950282822

Process finished with exit code 0
```