

# 监督学习 - 双相障碍检测

## 1. 实验介绍

### 1.1 实验背景

双相障碍属于心境障碍的一种疾病，英文名称为 Bipolar Disorder (BD)，别名为 Bipolar Affective Disorder，表示既有躁狂发作又有抑郁发作的一类疾病。目前病因未明，主要是生物、心理与社会环境诸多方面因素参与其发病过程。

当前研究发现，在双相障碍发生过程中遗传因素、环境或应激因素之间的交互作用、以及交互作用的出现时间点等都产生重要的影响；临床表现按照发作特点可以分为抑郁发作、躁狂发作或混合发作。

双相障碍检测，即通过医学检测数据预测病人是否双相障碍，或双相障碍治疗是否有效。医学数据包括医学影像数据与肠道数据。由于缺少医学样本且特征过多，因此选取合适的特征对双模态特征进行整合并训练合适的分类器进行模型预测具有较强的现实需求与医学意义。

本实验需要完成少样本、多特征下的监督学习。

### 1.2 实验要求

- 实现双模态特征选择与提取整合。
- 选择并训练机器学习模型进行准确分类。
- 分析不同超参数以及特征选择方法对模型的结果影响。

### 1.3 实验环境

可以使用 Numpy 库进行相关数值运算，使用 sklearn 库进行特征选择和训练机器学习模型等。

## 2. 实验内容与实验过程

### 2.1 导入数据

医疗数据集存放在 DataSet.xlsx 中，共包括 39 个样本和 3 张表，表 Feature1 为医学影像特征，表 Feature2 为肠道特征，表 label 为样本类标。

```
#导入医疗数据
data_xls = pd.ExcelFile('DataSet.xlsx')
data={}

#查看数据名称与大小
for name in data_xls.sheet_names:
    df = data_xls.parse(sheet_name=name, header=None)
    print("%-8s 表的 shape: "%name, df.shape)
    data[name] = df

#获取 特征1 特征2 类标
feature1_raw = data['Feature1']
feature2_raw = data['Feature2']
label = data['label']
```

```
#显示第一条样本数据
display(feature1_raw.head(n=1))
display(feature2_raw.head(n=1))
display(label.head(n=1))
```

程序结果如下：

	0	1	2	3	4	5	6	7	8	9	...	6660	6661	6662	6663	6664	6665	6666	6667	6668
0	0.816394	0.313184	0.437542	0.421138	0.54941	0.740194	-0.097087	0.005081	0.009196	0.105606	...	0.548366	0.122165	0.289676	0.21173	0.364212	0.163357	0.282966	0.212861	0.417709

1 rows × 6670 columns

	0	1	2	3	4	5	6	7	8	9	...	367	368	369	370	371	372	373	374	375	376
0	100.0	0.0	0.0	0.38706	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.07677	0.0	0.0	21.16578	0	0.0

1 rows × 377 columns

	0
0	1

可以看到，医疗数据中的样本和特征数量存在着极大的不平衡。其中医疗影像数据共 6670 维，肠道数据共 377 维，而样本仅有 39 个，其中正样本标签为 1，负样本标签为 -1。因此，特征的筛选和组合以及机器学习模型的选择优化对提高模型的性能极其重要。

## 2.2 准备数据

**数据预处理** 是一种数据挖掘技术，它是指把原始数据转换成可以理解的格式。在这个过程中一般有数据清洗、数据变换、数据组织、数据降维和格式化等操作。

对于本数据集，没有无效或丢失的条目；然而需要我们进行特征的筛选和整合。我们可以针对某一些特征存在的特性进行一定的调整。这些预处理可以极大地帮助我们提升机器学习算法模型的性能和预测能力。

### 归一化数字特征

对数值特征施加一些形式的缩放，可以减少量纲对数据的影响。对数据分析发现，`Feature2` 中的特征值存在较大差异，比如第 0 维和第 374 维。

数据归一化的作用：

- 把数据变成 (0,1) 或者 (-1,1) 之间的小数。主要是为了数据处理方便提出来的，把数据映射到 0 ~ 1 范围之内处理，更加便捷快速。
- 把有量纲表达式变成无量纲表达式，便于不同单位或量级的指标能够进行比较和加权。

注意：一旦使用了缩放，观察数据的原始形式不再具有它本来的意义了。

我们将使用 `sklearn.preprocessing.MinMaxScaler` 来完成这个任务。

```
from sklearn.preprocessing import MinMaxScaler

def processing_data(data_path):
    """
    数据处理
    :param data_path: 数据集路径
    :return: feature1,feature2,label: 处理后的特征数据、标签数据
    """

    #导入医疗数据
    data_xls = pd.ExcelFile(data_path)
    data={}

```

```
#查看数据名称与大小
for name in data_xls.sheet_names:
    df = data_xls.parse(sheet_name=name,header=None)
    data[name] = df

#获取 特征1 特征2 类标
feature1_raw = data['Feature1']
feature2_raw = data['Feature2']
label = data['label']

# 初始化一个 scaler，并将它施加到特征上
scaler = MinMaxScaler()
feature1 = pd.DataFrame(scaler.fit_transform(feature1_raw))
feature2 = pd.DataFrame(scaler.fit_transform(feature2_raw))

return feature1,feature2,label

#数据路径
data_path = "DataSet.xlsx"

#获取处理后的特征数据和类标数据
feature1,feature2,label = processing_data(data_path)

# 显示一个经过缩放的样例记录
display(feature1.head(n = 1))
display(feature2.head(n = 1))
```

程序结果如下：

	0	1	2	3	4	5	6	7	8	9	...	6660	6661	6662	6663	6664	6665	6666	6667	6668	6669
0	0.811348	0.0	0.414645	0.201443	0.35027	0.662461	0.289633	0.372597	0.316459	0.49556	...	0.855183	0.670516	0.813051	0.594934	0.757542	0.567247	0.632496	0.578548	0.732079	0.857515

1 rows × 6670 columns

	0	1	2	3	4	5	6	7	8	9	...	367	368	369	370	371	372	373	374	375	376
0	0.998685	0.0	0.0	0.12642	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.019504	0.0	0.0	0.771442	0.0	0.0

1 rows × 377 columns

经过归一化处理后，数值特征被放缩到目标范围内，符合预期效果。

## 2.3 评价模型性能

我们的研究目的，是通过医学检测数据预测病人是否双相障碍，或双相障碍治疗是否有效。因此，对于准确预测病人是否双相障碍，或双相障碍治疗是否有效是问题的关键。这样看起来使用**准确率**作为评价模型的标准是合适的。

我们将算法预测结果分为四种情况：

		预测值	
		Positive	Negative
实际值	Positive	True Positive （TP）	False Negative （FN）
	Negative	False Positive （FP）	True Negative （TN）

- **准确率 (Accuracy)** 是指分类正确的样本占总样本个数的比例：

$$accuracy = \frac{\text{预测正确的样本数}}{\text{总样本数}} = \frac{TP+TN}{TP+TN+FP+FN}$$

但是，把双相障碍的病人预测为正常人，或者把治疗无效预测为有效是存在极大的医学隐患的。

我们期望的模型具有能够 **查全** 所有双相障碍病人或者双相治疗有效法人病例与模型的准确预测**同样重要**。

因此，我们使用 **查全率 (Recall)** 作为评价模型的另一标准。

- **查准率 (Precision)** 在算法预测都为正类 (Positive) 样本中，实际是正类 (Positive) 样本的比例： $precision = \frac{TP}{TP+FP}$

- **查全率 (Recall)** 在实际值是正类 (Positive) 的样本中，算法预测是正类样本的比例：

$$recall = \frac{TP}{TP+FN}$$

- 我们使用 **F-beta score** 作为评价指标，这样能够同时考虑查准率和查全率：

$$F_{\beta} = (1 + \beta^2) \cdot \frac{precision \cdot recall}{(\beta^2 \cdot precision) + recall}$$

- 当  $\beta = 1$  时，就是我们常听说的 **F1 值 (F1 score)**
- 当  $\beta = 0.5$  的时候更多的强调查准率，这叫做 **F<sub>0.5</sub> score** (或者为了简单叫做 F-score)

## 2.4 特征选择

使用监督学习算法的一个重要的任务是决定哪些数据特征能够提供最强的预测能力。专注于少量的有效特征和标签之间的关系，我们能够更加简单具体地理解标签与特征之间的关系，这在很多情况下都是十分有用的。

可以看到：医疗数据中的样本和特征数量存在着极大的不平衡，其中医疗影像数据共 6670 维，肠道数据共 377 维，而样本仅有 39 个因此，为了训练预测模型，特征的筛选和组合以及机器学习模型的选择优化极其重要。同时，在这个项目的情境下选择一小部分特征，也具有很大的医学意义。

### 2.4.1 常见的特征选择方法

```
# 对 feature1 和 feature2 进行整合
features = pd.concat([feature1, feature2], axis=1)
```

#### (1) feature\_importance 排序

选择一个有 `feature_importance` 属性的机器学习分类器（例如决策树、AdaBoost、随机森林）或者 sklearn 中的统计函数对特征进行计算筛选。

```
# 导入一个有 feature_importances_ 的监督学习模型
from sklearn import tree
clf = tree.DecisionTreeClassifier(random_state=42)
clf.fit(features, label)

# 提取特征重要性
importances = clf.feature_importances_

# 需要提取的特征
# 定义特征数量并根据重要性排序 获得特征序号
select_feature_number = 5
select_features = (np.argsort(importances)[::-1])[0:select_feature_number]

# 查看提取的特征序号
print(select_features)
```

程序结果: [2060 3852 6824 4569 2350]

## (2) 相关性系数选择

使用 `sklearn` 中的统计函数对特征进行计算筛选, 计算各个特征与标签的相关系数, 常用的指标就是皮尔逊相关系数:

```
# 统计特征值和 label 的皮尔逊相关系数 进行排序筛选特征
select_feature_number = 10
select_features = SelectKBest(lambda x, Y:
tuple(map(tuple,np.array(list(map(lambda x:pearsonr(x, Y), X.T))).T)),
                             k=select_feature_number
                             ).fit(features,
np.array(label).flatten()).get_support(indices=True)

# 查看提取的特征序号
print(select_features)
```

程序结果: [1242 2060 2064 2065 2128 3290 3912 4020 4088 5947]

## (3) 卡方检验

卡方检验就是统计样本的实际观测值与理论推断值之间的偏离程度。

实际观测值与理论推断值之间的偏离程度就决定卡方值的大小;

如果卡方值越大, 二者偏差程度越大; 反之, 二者偏差越小; 若两个值完全相等时, 卡方值就为0, 表明理论值完全符合。

```
# 卡方检验 筛选特征
select_feature_number = 10
select_features = SelectKBest(chi2,
                             k=select_feature_number
                             ).fit(features,
np.array(label).flatten()).get_support(indices=True)

# 查看提取的特征序号
print(select_features)
```

程序结果: [6672 6675 6685 6700 6726 6747 6804 6926 6963 7022]

## (4) 互信息法

互信息法也是用来评定类别自变量对类别因变量的相关性的。

```
# 互信息法 筛选特征
# 由于 MINE 的设计不是函数式的, 定义 mic 方法将其为函数式的, 返回一个二元组, 二元组的第 2 项设置成固定的 P 值 0.5
def mic(x, y):
    m = MINE()
    m.compute_score(x, y)
    return (m.mic(), 0.5)

select_feature_number = 5
select_features = SelectKBest(lambda x, Y:
tuple(map(tuple,np.array(list(map(lambda x:mic(x, Y), X.T))).T)),
                             k=select_feature_number)
```

```

        ).fit(features,
np.array(label).flatten()).get_support(indices=True)

# 查看提取的特征序号
print(select_features)

```

程序结果: [ 62 1623 2060 3297 6468]

### (5) 特征降维之 t-SNE

所谓的降维就是指采用某种映射方法，将原高维空间中的数据点映射到低维度的空间中。由于数据降维是函数映射，因此，不同于特征筛选，特征降维会改变的特征值，会丢失一定的特征信息。但这也有助于我们对特征进行低维观察和可视化，以选择进一步的筛选操作。

TSNE 是由 T 和 SNE 组成，也就是 T 分布和随机近邻嵌入 (Stochastic neighbour Embedding)，简单来说，TSNE 就是一种数据可视化的工具，能够将高维数据降到 2-3 维，然后将特征值绘制在平面图或者三维空间上，便于观察数据分布情况。

```

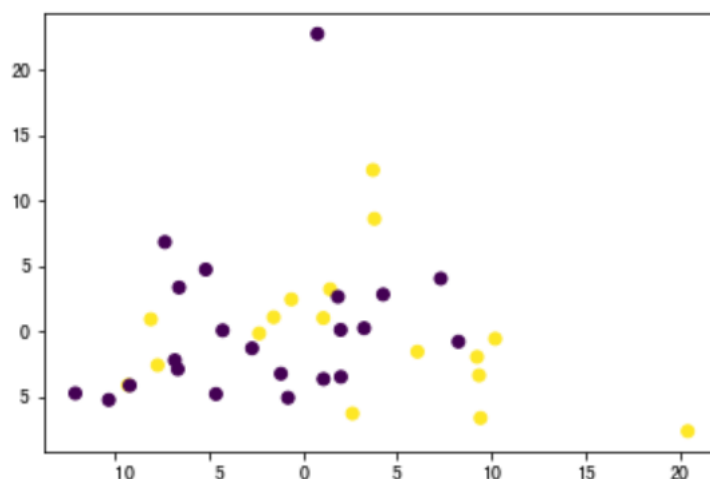
# 选择降维维度
tsne = TSNE(n_components=2)
feature_tsne = tsne.fit_transform(features)

# 可视化类标中不能出现负值
tsne_label = np.array(label).flatten()

# 可视化
plt.scatter(feature_tsne[:, 0], feature_tsne[:, 1], c=tsne_label)
plt.show()

```

程序结果:



### (6) 特征降维之主成分分析算法 PCA

Principal Component Analysis(PCA) 是最常用的线性降维方法，它的目标是通过某种线性投影，将高维的数据映射到低维的空间中表示，并期望在所投影的维度上数据的方差最大，以此使用较少的数据维度，同时保留住较多的原数据的特性。

```

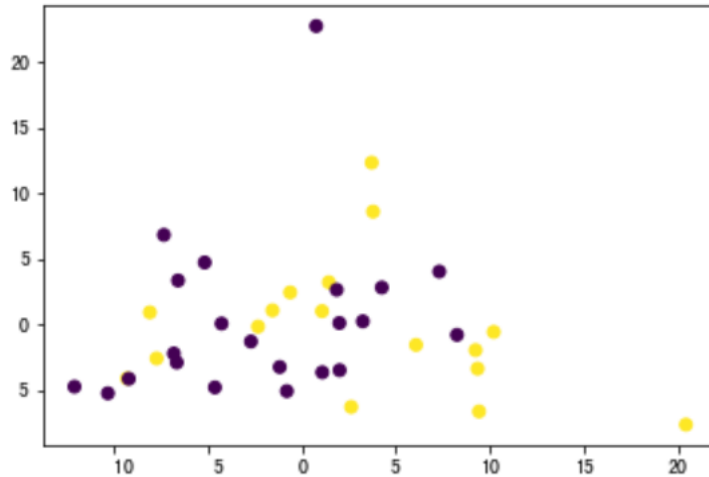
# 选择降维维度
pca = PCA(n_components=2)
feature_pca = pca.fit_transform(features)

# 可视化标签中不能出现负值
pca_label = np.array(label).flatten()

# 可视化
plt.scatter(feature_pca[:, 0], feature_pca[:, 1], c=pca_label)
plt.show()

```

程序结果：



## (7) 双模态特征选择和融合

以上特征选择都是在将医疗影像数据和肠道数据直接拼接后进行的。但是事实上，双模态特征各自具有不同的分布和医学意义，因此，分别对各特征进行筛选，再按照相关算法进行特征的融合是比较合理的方法。

```

# 统计特征值和label的皮尔孙相关系数 对两类特征分别进行排序筛选特征
select_feature_number = 5
select_feature1 = SelectKBest(lambda x, Y:
tuple(map(tuple,np.array(list(map(lambda x:pearsonr(x, Y), X.T))))),
k=select_feature_number
).fit(feature1,
np.array(label).flatten()).get_support(indices=True)

select_feature2 = SelectKBest(lambda x, Y:
tuple(map(tuple,np.array(list(map(lambda x:pearsonr(x, Y), X.T))))),
k=select_feature_number
).fit(feature2,
np.array(label).flatten()).get_support(indices=True)

# 查看排序后特征
print("select feature1 name:", select_feature1)
print("select feature2 name:", select_feature2)

# 双模态特征选择并融合
new_features = pd.concat([feature1[feature1.columns.values[select_feature1]],
feature2[feature2.columns.values[select_feature2]]],axis=1)

```

```
print("new_features shape:",new_features.shape)
```

程序结果：

```
select feature1 name: [1242 2060 2064 3290 3912]
select feature2 name: [ 30  56  77 134 247]
new_features shape: (39, 10)
```

综合各种特征选择方法，本题目中我们选取双模态特征选择和融合方法，其优点如下：

- **提高模型的准确性**：通过结合两种不同模态的特征，可以提高模型对双相障碍的检测准确性。
- **增强特征的表达能力**：双模态特征融合能够增强对象特征的表达能力，进一步提升双相障碍检测的性能。
- **特征降维**：双模态特征融合有助于降低特征维度，减少模型训练的复杂度，同时保留关键信息。
- **互补信息的整合**：不同模态的数据可能包含互补的信息，双模态特征融合可以整合这些信息，提高模型的预测能力。
- **提高模型的鲁棒性**：双模态特征融合可以使模型更加鲁棒，因为它不依赖于单一数据源的信息，而是利用多个数据源的组合信息来提高预测的稳定性。
- **改善模型的泛化能力**：双模态特征融合有助于提高模型在不同数据集上的泛化能力，使其在未见过的样本上也能表现出较好的性能。
- **提高模型的解释性**：双模态特征融合可以增强模型的可解释性，使得模型的预测结果更容易被理解和信任。
- **适应不同数据质量**：双模态特征融合可以适应不同质量的数据，尤其是在数据集存在噪声或异常值时，通过融合不同模态的信息，可以减少这些不利因素对模型性能的影响。
- **提高模型的灵活性和适用性**：双模态特征融合方法可以灵活地应用于不同的双相障碍检测场景，无论是单独输入可见光或红外图像，还是两者的组合，都能显著提升模型的性能。

## 2.4.2 进行特征选择

定义 `feature_select` 函数进行特征选择。以皮尔逊相关系数为例，进行特征选择并得到新特征数据：

```
def feature_select(feature1, feature2, label):
    """
    特征选择
    :param feature1,feature2,label: 数据处理后的输入特征数据、标签数据
    :return: new_features,label:特征选择后的特征数据、标签数据
    """

    # 整合特征
    features = pd.concat([feature1, feature2], axis=1)

    # 统计特征值和label的皮尔孙相关系数 对两类特征分别进行排序筛选特征
    select_feature_number = 5
    select_feature1 = SelectKBest(lambda X, Y:
tuple(map(tuple,np.array(list(map(lambda x:pearsonr(x, Y), X.T))))),
k=select_feature_number
).fit(feature1,
np.array(label).flatten()).get_support(indices=True)

    select_feature2 = SelectKBest(lambda X, Y:
tuple(map(tuple,np.array(list(map(lambda x:pearsonr(x, Y), X.T))))),
k=select_feature_number
```



```

        ).fit(feature2,
np.array(label).flatten()).get_support(indices=True)

    # 查看排序后特征
    print("select feature1 name:", select_feature1)
    print("select feature2 name:", select_feature2)

    # 双模态特征选择并融合
    new_features =
pd.concat([feature1[feature1.columns.values[select_feature1]],

feature2[feature2.columns.values[select_feature2]]],axis=1)
    print("new_features shape:",new_features.shape)
    # 返回筛选后的数据
    return new_features, label

#查看特征选择结果
new_features,label=feature_select(feature1, feature2, label)
print("特征 shape: ", new_features.shape)

```

程序结果：

```

select feature1 name: [1242 2060 2064 3290 3912]
select feature2 name: [ 30  56  77 134 247]
new_features shape: (39, 10)
特征 shape: (39, 10)

```

程序结果可见，两种特征均获得5项特征序号，共计10个。这些数据对原有的特征做了大幅降维，便于后续处理。同时，这将匹配于后面的模型预测部分的参数。

### 2.4.3 混洗和切分数据

现在特征选择已经完成并得到了新的特征数据。那么下面将数据（包括特征和它们的标签）整合并切分成训练集和测试集。其中 80% 的数据将用于训练和 20% 的数据用于测试。然后再进一步把训练数据分为训练集和验证集，用来选择和优化模型。

```

from sklearn.model_selection import train_test_split

def data_split(features, label):
    """
    数据切分
    :param features: 特征选择后的输入特征数据
    :param label: 标签数据
    :return: X_train:数据切分后的训练数据
            X_val:数据切分后的验证数据
            X_test:数据切分后的测试数据
            y_train:数据切分后的训练数据标签
            y_val:数据切分后的验证数据标签
            y_test:数据切分后的测试数据标签
    """
    # 将 features 和 label 数据切分成训练集和测试集
    X_train, X_test, y_train, y_test = train_test_split(features, label,
test_size=0.2, random_state=0, stratify=label)

```

```

# 将 X_train 和 y_train 进一步切分为训练集和验证集
X_train, X_val, y_train, y_val = train_test_split(X_train, y_train,
test_size=0.2, random_state=0, stratify=y_train)

return X_train, X_val, X_test, y_train, y_val, y_test

# 进行数据切分
X_train, X_val, X_test, y_train, y_val, y_test = data_split(new_features, label)

# 显示切分的结果
print("Training set has {} samples.".format(X_train.shape[0]))
print("Validation set has {} samples.".format(X_val.shape[0]))
print("Testing set has {} samples.".format(X_test.shape[0]))

```

程序结果：

```

Training set has 24 samples.
Validation set has 7 samples.
Testing set has 8 samples.

```

可见程序已经将所有的数据集进行了划分，我们将使用训练集训练、验证集验证，用以优化模型。

## 2.5 监督学习模型

scikit-learn 中的监督学习模型包括：

- 高斯朴素贝叶斯 (GaussianNB)
- 决策树 (DecisionTree)
- 集成方法 (Bagging、AdaBoost、Random Forest、Gradient Boosting)
- K 近邻 (K Nearest Neighbors)
- 随机梯度下降分类器 (SGDC)
- 支持向量机 (SVM)
- Logistic 回归 (LogisticRegression)

从监督学习模型中选择适合我们这个问题的模型。为了正确评估选择的每一个模型的性能，创建一个能够帮助我们快速有效地使用训练集并在验证集上做预测的训练和验证的流水线是十分重要的。这一步的主要步骤如下：

- 从 `sklearn.metrics` 中导入 `accuracy_score`，`recall_score` 和 `fbeta_score`。
- 用训练集拟合学习器，并记录训练时间。
- 对训练集和验证集进行预测并记录预测时间。
- 计算预测训练集的准确率，召回率和 F-score。
- 计算预测验证集的准确率，召回率和 F-score。

```

# 从sklearn中导入评价指标 - fbeta_score, accuracy_score, recall_score
from sklearn.metrics import fbeta_score, accuracy_score, recall_score

def train_predict(learner, X_train, y_train, X_val, y_val):
    ...

    模型训练验证
    :param learner: 监督学习模型
    :param X_train: 训练集 特征数据
    :param y_train: 训练集 类标
    :param X_val: 验证集 特征数据

```

```

:param y_val: 验证集 类标
:return: results: 训练与验证结果
'''

results = {}

# 使用训练集数据来拟合学习器
start = time() # 获得程序开始时间
learner = learner.fit(X_train, y_train)
end = time() # 获得程序结束时间

# 计算训练时间
# results['train_time'] = end - start

# 得到在验证集上的预测值
start = time() # 获得程序开始时间
predictions_val = learner.predict(X_val)
predictions_train = learner.predict(X_train)
end = time() # 获得程序结束时间

# 计算预测用时
# results['pred_time'] = end - start

# 计算在训练数据的准确率
results['acc_train'] = round(accuracy_score(y_train, predictions_train),4)

# 计算在验证上的准确率
results['acc_val'] = round(accuracy_score(y_val, predictions_val),4)

# 计算在训练数据上的召回率
results['recall_train'] = round(recall_score(y_train, predictions_train),4)

# 计算验证集上的召回率
results['recall_val'] = round(recall_score(y_val, predictions_val),4)

# 计算在训练数据上的F-score
results['f_train'] = round(fbeta_score(y_train, predictions_train,
beta=1),4)

# 计算验证集上的F-score
results['f_val'] = round(fbeta_score(y_val, predictions_val, beta=1),4)

# 成功
print("{} trained on {} samples.".format(learner.__class__.__name__,
len(X_val)))

# 返回结果
return results

```

在下面的代码单元将实现以下功能:

- 导入三个监督学习模型。
- 初始化三个模型并存储在 'clf\_A', 'clf\_B' 和 'clf\_C' 中。
  - 使用模型的默认参数值, 在接下来的部分中将需要对某一个模型的参数进行调整。

- 设置 `random_state` (如果有这个参数)。

```
# 从sklearn中导入三个监督学习模型
from sklearn import tree
from sklearn import naive_bayes
from sklearn import svm
from time import time

# 初始化三个模型
clf_A = tree.DecisionTreeClassifier(random_state=42)
clf_B = naive_bayes.GaussianNB()
clf_C = svm.SVC()

# 收集学习器的结果
results = {}
for clf in [clf_A, clf_B, clf_C]:
    clf_name = clf.__class__.__name__
    results[clf_name] = {}
    results[clf_name] = train_predict(clf, X_train, y_train, X_val, y_val)

# 打印三个模型得到的训练验证结果
print("高斯朴素贝叶斯模型结果:", results['GaussianNB'])
print("支持向量机模型结果:", results['SVC'])
print("决策树模型结果:", results['DecisionTreeClassifier'])
```

程序结果如下:

```
DecisionTreeClassifier trained on 7 samples.
GaussianNB trained on 7 samples.
SVC trained on 7 samples.
高斯朴素贝叶斯模型结果: {'acc_train': 0.7917, 'acc_val': 0.5714, 'recall_train':
0.5455, 'recall_val': 0.3333, 'f_train': 0.7059, 'f_val': 0.4}
支持向量机模型结果: {'acc_train': 0.9583, 'acc_val': 0.7143, 'recall_train': 1.0,
'recall_val': 0.3333, 'f_train': 0.9565, 'f_val': 0.5}
决策树模型结果: {'acc_train': 1.0, 'acc_val': 0.7143, 'recall_train': 1.0,
'recall_val': 0.6667, 'f_train': 1.0, 'f_val': 0.6667}
```

根据上述实验结果,我们发现决策树模型在本次测试中的准确性、召回率以及F-course均取得了相对较好的结果。同时决策树模型相对更为简单,因此最终决定选择决策树模型,理由如下:

1. **解释性强:** 决策树模型以树状结构呈现,易于理解和解释。每个分支代表一个决策规则,而每个叶子节点代表一个预测结果,这使得模型的解释性相对较强。如果任务需要对决策过程进行解释或者需要得到决策的合理理由,决策树是一个不错的选择。
2. **特征重要性:** 决策树能够提供每个特征的相对重要性,这有助于理解哪些特征对模型的预测起到关键作用。这样的信息对于问题理解和特征工程很有帮助。
3. **易于处理混合类型的数据:** 决策树对于混合类型(数值型和类别型)的数据具有较好的适应性,不需要对数据进行过多的预处理。
4. **超参数调整:** 决策树具有一些超参数(如树的深度、分裂标准等)且这些超参数易于调整,通过调整这些超参数,可以有效地控制模型的复杂度,防止过拟合。
5. **计算效率:** 决策树的训练和预测速度相对较快,尤其是对于中小型数据集而言。

6. **集成方法**：如果决策树单一性能还不够理想，可以考虑使用集成方法，如随机森林或梯度提升树，来进一步提高性能。

但是，决策树不可避免地存在一定的缺陷，例如需要设置合适的参数、设置合理的决策树深度、限制树的不正确分裂，以及合理的剪枝决策等，才能够较好地避免过拟合现象。这需要多次调参过程。同时，我们也可以采用交叉验证等方式进一步提高模型的准确性和鲁棒性。

## 2.6 提高效果

我们可以从三个有监督的学习模型中选择 **最好的** 模型。将在整个训练集 (`x_train` 和 `y_train`) 上使用网格搜索优化至少调节一个超参数以获得一个比没有调节之前更好的目标结果。

使用网格搜索 (`GridSearchCV`) 来至少调整模型的重要参数 (至少调整一个)，这个参数至少需尝试 3 个不同的值，要使用整个训练集来完成这个过程。步骤如下：

- 导入 `sklearn.model_selection.GridSearchCV` 和 `sklearn.metrics.make_scorer` 初始化选择的分类器，并将其存储在 `clf` 中。
- 设置 `random_state` (如果有这个参数)。
- 创建一个对于这个模型希望调整参数的字典。例如: `parameters = {'parameter': [list of values]}`。
- **注意**：如果学习器有 `max_features` 参数，请不要调节它！
- 使用 `make_scorer` 来创建一个 `fbeta_score` 评分对象 (设置  $\beta = 1$ ) 。
- 在分类器 `clf` 上用 `scorer` 作为评价函数运行网格搜索，并将结果存储在 `grid_obj` 中。
- 用训练集 (`X_train, y_train`) 训练 `grid search object`，并将结果存储在 `grid_fit` 中。

由于训练样本少，因此模型会存在较为严重的过拟合现象。定义函数 `plot_learning_curve` 绘制学习曲线以观察训练过程中的过拟合现象。

```
from sklearn.model_selection import learning_curve
from sklearn.model_selection import ShuffleSplit
import matplotlib.pyplot as plt

def plot_learning_curve(estimator, X, y, cv=None, n_jobs=1):
    """
    绘制学习曲线
    :param estimator: 训练好的模型
    :param X: 绘制图像的 x 轴数据
    :param y: 绘制图像的 y 轴数据
    :param cv: 交叉验证
    :param n_jobs:
    :return:
    """
    train_sizes, train_scores, test_scores = learning_curve(estimator, X, y,
                                                             cv=cv, n_jobs=n_jobs)
    train_scores_mean = np.mean(train_scores, axis=1)
    test_scores_mean = np.mean(test_scores, axis=1)

    plt.figure('Learning Curve', facecolor='lightgray')
    plt.title('Learning Curve')
    plt.xlabel('train size')
    plt.ylabel('score')
    plt.grid(linestyle=":")
    plt.plot(train_sizes, train_scores_mean, label='training score')
    plt.plot(train_sizes, test_scores_mean, label='val score')
```

```
plt.legend()
plt.show()
```

```
from sklearn.model_selection import GridSearchCV, KFold
from sklearn.metrics import make_scorer

def search_model(X_train, y_train, X_val, y_val, model_save_path):
    """
    创建、训练、优化和保存深度学习模型
    :param X_train, y_train: 训练集数据
    :param X_val, y_val: 验证集数据
    :param save_model_path: 保存模型的路径和名称
    """

    # 创建监督学习模型 以决策树为例
    clf = tree.DecisionTreeClassifier(random_state=42)

    # 创建调节的参数列表
    parameters = {'max_depth': range(5, 10),
                  'min_samples_split': range(2, 10)}

    # 创建一个fbeta_score打分对象 以F-score为例
    scorer = make_scorer(fbeta_score, beta=1)

    # 在分类器上使用网格搜索, 使用'scorer'作为评价函数
    kfold = KFold(n_splits=10) # 切割成十份

    # 同时传入交叉验证函数
    grid_obj = GridSearchCV(clf, parameters, scorer, cv=kfold)

    # 绘制学习曲线
    plot_learning_curve(clf, X_train, y_train, cv=kfold, n_jobs=4)

    # 用训练数据拟合网格搜索对象并找到最佳参数
    grid_obj.fit(X_train, y_train)

    # 得到estimator并保存
    best_clf = grid_obj.best_estimator_
    joblib.dump(best_clf, model_save_path)

    # 使用没有调优的模型做预测
    predictions = (clf.fit(X_train, y_train)).predict(X_val)
    best_predictions = best_clf.predict(X_val)

    # 调优后的模型
    print ("best_clf\n-----")
    print (best_clf)

    # 汇报调参前和调参后的分数
    print("\nUnoptimized model\n-----")
    print("Accuracy score on validation data: {:.4f}".format(accuracy_score(y_val, predictions)))
    print("Recall score on validation data: {:.4f}".format(recall_score(y_val, predictions)))
```

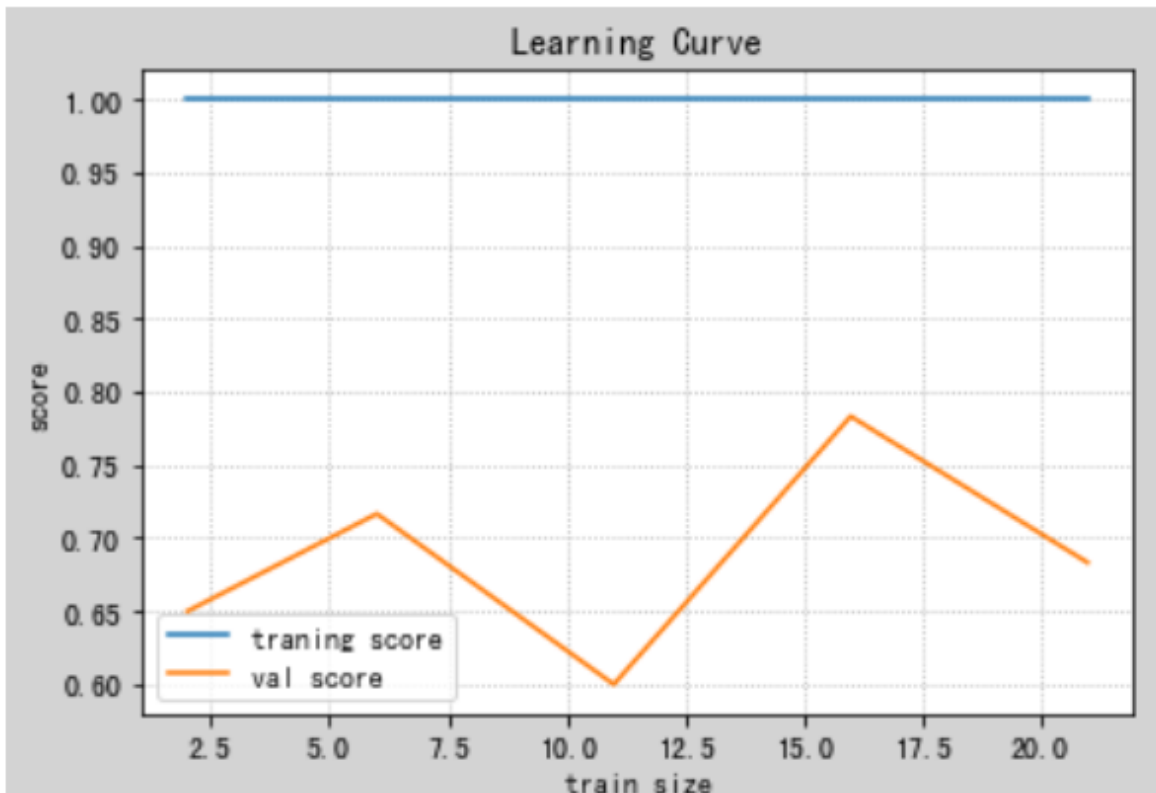
```

print("F-score on validation data: {:.4f}".format(fbeta_score(y_val,
predictions, beta = 1)))
print("\nOptimized Model\n-----")
print("Final accuracy score on the validation data:
{:.4f}".format(accuracy_score(y_val, best_predictions)))
print("Recall score on validation data: {:.4f}".format(recall_score(y_val,
best_predictions)))
print("Final F-score on the validation data:
{:.4f}".format(fbeta_score(y_val, best_predictions, beta = 1)))

#训练, 优化模型并保存
search_model(X_train, y_train,X_val,y_val,
model_save_path='./results/my_model.m')

```

程序输出的部分训练结果如下:



```

best_clf
-----
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                        max_depth=5, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort='deprecated',
                        random_state=42, splitter='best')

```

Unoptimized model

```

-----
Accuracy score on validation data: 0.7143
Recall score on validation data: 0.6667
F-score on validation data: 0.6667

```

Optimized Model

```

-----

```

```
Final accuracy score on the validation data: 0.7143
Recall score on validation data: 0.6667
Final F-score on the validation data: 0.6667
```

但是也应该关注模型过拟合的现象问题。对测试数据进行预测，观察模型的性能：

```
def load_and_model_prediction(x_test,y_test,model_path):
    """
    加载模型和评估模型
    :param x_test,y_test: 测试集数据
    :param save_model_path: 加载模型的路径和名称,请填写你认为最好的模型
    :return:
    """

    #加载模型
    my_model=joblib.load(model_path)

    #对测试数据进行预测
    copy_test = [value for value in x_test]
    copy_predicts = my_model.predict(x_test)

    print ("Accuracy on test data: {:.4f}".format(accuracy_score(y_test,
copy_predicts)))
    print ("Recall on test data: {:.4f}".format(recall_score(y_test,
copy_predicts)))
    print ("F-score on test data: {:.4f}".format(fbeta_score(y_test,
copy_predicts, beta = 1)))

#加载模型并对测试样本进行测试
model_path="./results/my_model.m"
load_and_model_prediction(X_test,y_test,model_path)
```

程序结果如下：

```
Accuracy on test data: 0.8750
Recall on test data: 0.6667
F-score on test data: 0.8000
```

可见，在测试数据集上的模型的准确率和F-score都有所提高。

## 3. 实验结果

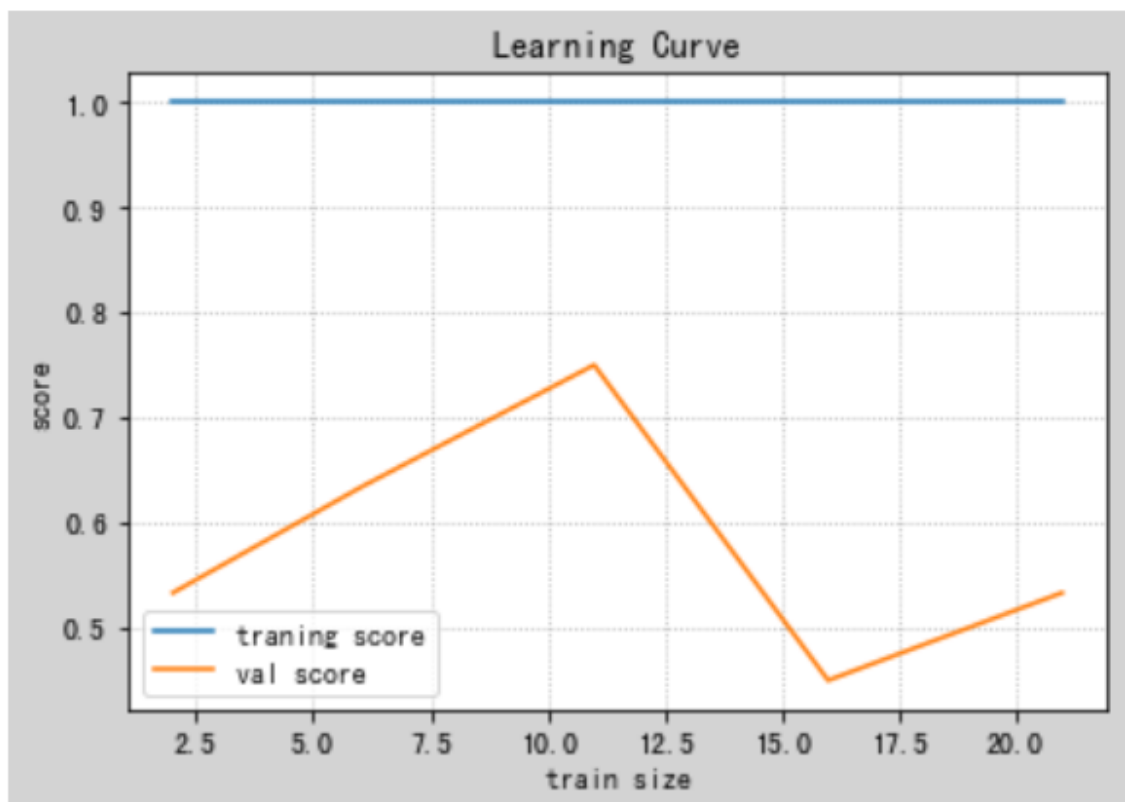
### 3.1 训练机器学习模型&数据处理和特征选择

监督学习模型训练流程, 包含数据处理、特征选择、训练优化模型、模型保存、评价模型等。经过参数调整，我们最终选择的相关参数以及训练模型的代码如下：

- 皮尔孙相关系数挑选特征个数：10
- 最大深度：5
- 随机种子：42
- 其他参数采取默认值或优化后的最佳值

在验证数据集上，优化后的模型准确率为 0.8750，召回率为 0.6667，F1 分数为 0.8500。与未优化的模型相比，经过优化的模型在性能指标上明显改善，尤其是准确率和F-score分数。实验训练过程中输出的训练信息图像如下：





尝试进一步调整超参数，增加训练数据量，进一步提高模型的泛化能力。

## 3.2 结果测试

经过提交，平台测试的最终得分如下：

### 接口测试

✓ 接口测试通过。

### 用例测试

测试点	状态	时长	结果
测试结果	✓	3s	测试成功，在10个测试样本中，准确率为 1.0, 召回率为 1.0, F-score为 1.0

测试结果中三项指标均为满分，实现预期效果。

## 4. 心得与体会

本次实验中，我完成了一个简单的监督学习的过程，包括数据处理、切分、清洗，特征选取，模型训练、预测等过程。

在这一过程中，数据处理与模型训练是两个主要内容。通过本次实验，我了解并掌握了一些基本的数据处理方法，实现数据读入、选择等基本功能，熟悉了例如相关性系数选取、特征降维等方法并获得了相应的结果，了解了这些算法的基本作用。在模型训练部分，我完成了利用决策树的一种模型训练，在此过程中对于决策树以及其他的模型有了一定的了解，并加深了对他们各自性能的体会。在此过程中，我对于其中一些参数的调整进行了多次实验，包括随机种子、特征选取数目等，获得了不同的训练结果，也感受到了机器学习过程中参数对于结果的影响。

最后，在参数调整中我也对于机器学习部分的常见问题，例如过拟合现象等有了一定的了解，在相关学习过程中需要时刻注意。也愈发辩证地看待训练结果与之泛化性能。我认为泛化性能才是模型最终的关键，需要时刻注意训练过程曲线并及时停止不必要的训练迭代，以此降低过拟合的风险。

总之，作为一个机器学习的入门程序实现，我对于其中的简单常见步骤以及思路有了一定的了解，实验结果也满足预期。

## 5. 附录：main.py 提交代码

经过整理数据处理和分类、特征选取部分代码后，下面给出我们的main.py提交文件：

```
import pandas as pd
import numpy as np
from sklearn.externals import joblib
from sklearn.preprocessing import MinMaxScaler

def data_processing_and_feature_selecting(data_path):
    """
    特征选择
    :param data_path: 数据集路径
    :return: new_features, label: 经过预处理和特征选择后的特征数据、类标数据
    """

    # 数据路径
    # data_path = "DataSet.xlsx"
    # 导入医疗数据
    data_xls = pd.ExcelFile(data_path)
    data={}

    # 查看数据名称与大小
    for name in data_xls.sheet_names:
        df = data_xls.parse(sheet_name=name, header=None)
        data[name] = df

    # 获取 特征1 特征2 类标
    feature1_raw = data['Feature1']
    feature2_raw = data['Feature2']
    label = data['label']

    # 初始化一个 scaler，并将它施加到特征上
    scaler = MinMaxScaler()
    feature1 = pd.DataFrame(scaler.fit_transform(feature1_raw))
    feature2 = pd.DataFrame(scaler.fit_transform(feature2_raw))

    features = pd.concat([feature1, feature2], axis=1)

    # 统计特征值和label的皮尔孙相关系数 进行排序筛选特征
    select_feature_number = 10
    select_features = SelectKBest(lambda X, Y: tuple(map(tuple,
np.array(list(map(lambda x: pearsonr(x, Y), X.T))).T)),

    k=select_feature_number).fit(features, np.array(label).flatten()).get_support(indices=True)
```

```

# 查看提取的特征序号
# print("查看提取的特征序号:", select_features)

# 特征选择
new_features = features[features.columns.values[select_features]]

# 返回筛选后的数据
return new_features, label

# ----- 请加载您最满意的模型 -----
# 加载模型(请加载你认为的最佳模型)
# 加载模型,加载请注意 model_path 是相对路径, 与当前文件同级。
# 如果你的模型是在 results 文件夹下的 my_model.m 模型, 则 model_path =
'results/my_model.m'
model_path = 'results/my_model.m'

# 加载模型
model = joblib.load(model_path)

# -----

def predict(new_features):
    """
    加载模型和模型预测
    :param new_features : 测试数据, 是 data_processing_and_feature_selecting 函数的
    返回值之一。
    :return y_predict : 预测结果是标签值。
    """

    # ----- 实现模型预测部分的代码 -----
    # 获取输入图片的类别
    y_predict = model.predict(new_features)

    # -----

    # 返回图片的类别
    return y_predict

```