实验三　朴素贝叶斯分类器

- 安装机器学习库scikit-learn

- 学会调用sklearn中的朴素贝叶斯分类器，对20newsgroups数据集进行分类

- 完成并提交实验报告 <span style="color:red">（4月9日上课前交给班长）</span>

- **结果示例**：分类准确率达到80%以上



```
naive_bayes ×
D:\Software\anaconda3\envs\matplot\python.exe
Overall accuracy:   0.8385214007782101


Process finished with exit code 0
```

- 下载学习通里：<span style="color:red">资料/实验课ppt/20news-bydate_py3.pkz</span>

- 新建一个python文件，输入以下代码：

```python
from sklearn.datasets import fetch_20newsgroups
news = fetch_20newsgroups(data_home='./', subset='all')
print(news.data[0])
```

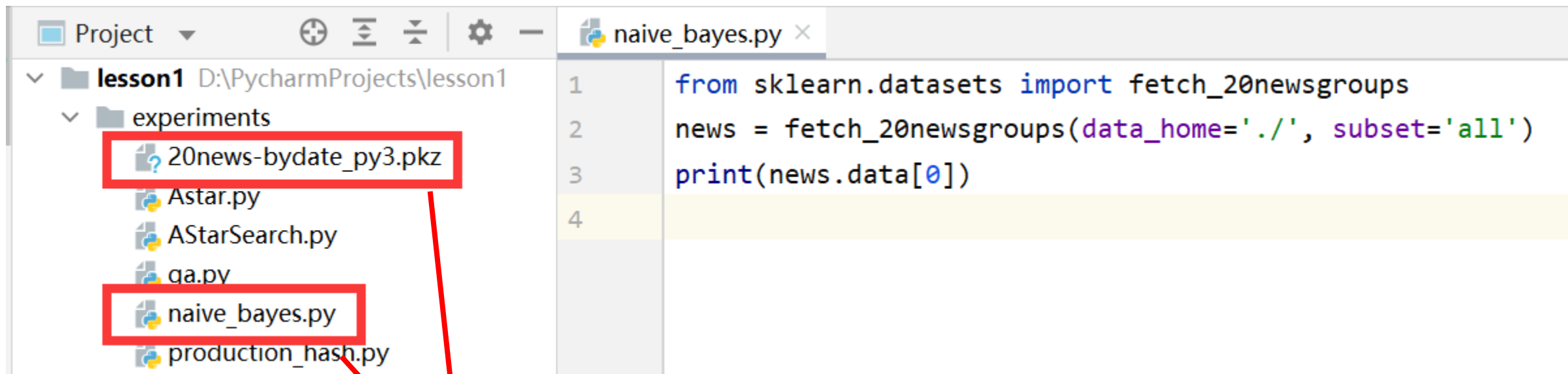- 将<span style="color:red">20news-bydate_py3.pkz</span>放在python文件的同目录下

- 运行上述python代码

- 例：



```python
from sklearn.datasets import fetch_20newsgroups
news = fetch_20newsgroups(data_home='./', subset='all')
print(news.data[0])
```

相同目录下！

- 成功示例：

✓ **20newsgroups数据集简介**

  ✓ 包含来自20个不同新闻组的文本数据。每个新闻组都包含多篇新闻文档，总共有18,846篇文档。

  ✓ 该数据集的文本数据涵盖了多个主题，包括科技、政治、体育、娱乐等。每个文档都被分配了一个特定的标签，表示其所属的新闻组类别。

$$P(X \mid A, B, C) = \frac{P(X)P(A \mid X)P(B \mid X, A)P(C \mid X, A, B)}{P(A, B, C)}$$

$$= \alpha \cdot P(X)P(A \mid X)P(B \mid X, A)P(C \mid X, A, B)$$

**朴素：** 假设特征$A$，$B$，$C$之间两两独立

**原条件概率公式可转化为：**

$$\alpha \cdot P(X)P(A \mid X)P(B \mid X, A)P(C \mid X, A, B)$$

$$= \alpha \cdot P(X)P(A \mid X)P(B \mid X)P(C \mid X)$$

- 导入需要使用的方法：

```python
# 朴素贝叶斯分类器
from sklearn.naive_bayes import MultinomialNB

# 文本向量化
from sklearn.feature_extraction.text import CountVectorizer

# 数据集分割
from sklearn.model_selection import train_test_split

# 计算分类准确度
from sklearn.metrics import accuracy_score
```

- 定义news_predict方法，用于训练朴素贝叶斯分类器，并对测试集中的样本进行分类：

```python
def news_predict(train_sample, train_label, test_sample):
    vectorizer = CountVectorizer()
    X_train = vectorizer.fit_transform(train_sample)

    # 训练朴素贝叶斯分类器
    NB_classifier = MultinomialNB()
    NB_classifier.fit(X_train, train_label)

    X_test = vectorizer.transform(test_sample)

    # 预测测试样本的类别
    y_pred = NB_classifier.predict(X_test)
    return y_pred
```

**#** 文本向量化
```
vectorizer = CountVectorizer()
X_train = vectorizer.fit_transform(train_sample)
```

- 数据集 $Dataset = \{Training\ set,\ Test\ set\}$

- 一般训练集占80%，测试集占20%

- 样本 $\vec{x_i} = (a_i, b_i, c_i, ...)$，其中$a_i, b_i, c_i$均为特征(feature)
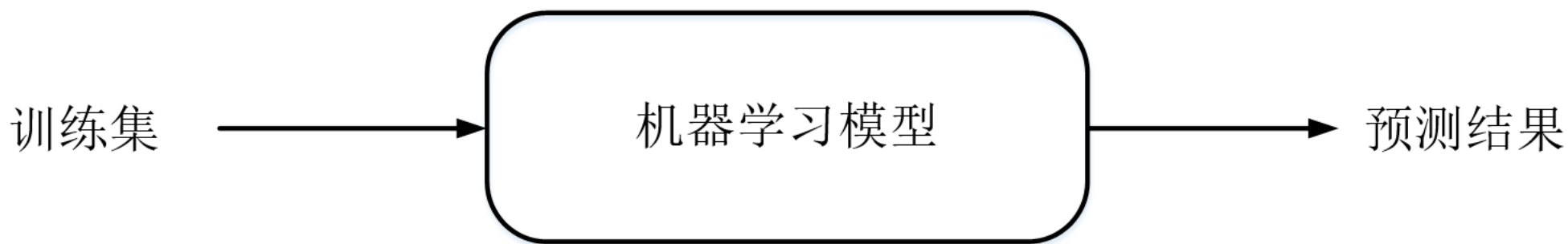
- 由于原始数据不是向量，因此要先对每个样本进行向量化

# 训练朴素贝叶斯分类器
```
NB_classifier = MultinomialNB()
NB_classifier.fit(X_train, train_label)
```

**训练阶段：**



- 设计一个损失函数 $loss = loss\_func($ 预测结果，真实结果 $)$
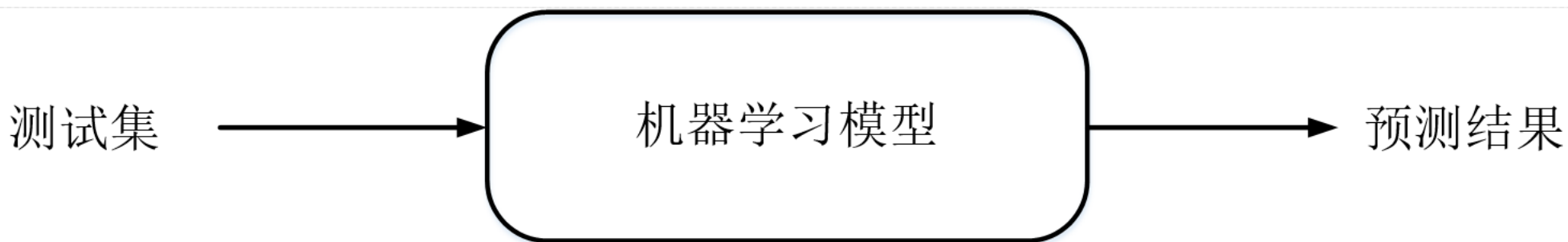- 通过优化算法降低loss，使得预测结果和真实结果尽可能接近
- 优化loss的过程 = 调整机器学习模型的参数

```
X_test = vectorizer.transform(test_sample)
# 预测测试样本的类别
y_pred = NB_classifier.predict(X_test)
```

**测试阶段：**

测试集 → 机器学习模型 → 预测结果

- 计算 $acc = Accuracy($预测结果，真实结果$)$，以评估模型性能

```python
# x表示新闻文本，y表示类别(标签)
x = news.data
y = news.target

# 把数据集分成测试集和训练集
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)

# 调用之前写好的news_predict方法
y_predict = news_predict(x_train, y_train, x_test)

# 计算总体分类准确度
acc = accuracy_score(y_test, y_predict)
print('Overall accuracy: ', acc)
```

谢 谢！