



统计机器学习实验报告

实验名称：多种模型的葡萄酒品种分类任务

姓名：	王恒
学院：	数学与统计学院
专业：	计算数学
学号：	220220934161

2023 年 6 月 12 日

摘要

本文中,我们实现了朴素贝叶斯、逻辑斯谛回归、线性支持向量机、基于高斯核函数的非线性支持向量机、多层非线性感知机(神经网络)模型及其相应的学习算法,并用以上模型对葡萄酒品种进行分类,六种模型分别得到了 95.83%, 97.22%, 94.44%, 95.83%, 97.22% 的准确率.

具体的,我们将数据集按照 6 : 4 的比例划分为训练集和测试集,并进行归一化. 其中,由于该数据的所有特征都是连续型数据,我们在实现朴素贝叶斯模型时,将所有特征归一化到 $[-1, 1]$, 并将其划分为固定的小区间,模型中的概率均为特征值所处区间的概率. 对于其他四种回归模型,我们在输入数据后先采用了批量归一化技术,自动学习归一化参数,相应的,均采用小批量随机梯度下降算法在训练集上进行优化,在测试集上测试准确率. 我们将本实验报告的所有内容开源:

<https://github.com/WANGH950/Statistical-Machine-Learning/tree/main/2ND>.

关键词: 朴素贝叶斯, 逻辑斯谛回归, 支持向量机, 神经网络, 葡萄酒品种分类

目录

1	开始实验	1
2	实验结果分析	2
2.1	朴素贝叶斯模型结果分析	2
2.2	Logistic 回归模型结果分析	3
2.3	支持向量机模型结果分析	4
2.4	神经网络模型结果分析	5
A	模型及算法源码	7
A.1	朴素贝叶斯模型及其训练算法	7
A.2	Logistic 模型及相应学习算法	9
A.3	支持向量机	10
A.4	全连接神经网络	13

1 开始实验

首先加载需要用到的包

```
1 import torch
2 import torch.nn as nn
3 from collections import OrderedDict
4 import copy
5 import time
6 from sklearn import datasets
7 import matplotlib.pyplot as plt
8 torch.manual_seed(100)
```

我们的所有实验都基于 PyTorch 神经网络框架实现, 该框架不仅包含了数值计算包 `numpy` 中的几乎所有工具, 还提供了编写神经网络的基本工具, 起到了类似于建房子时提供” 砖头” 和” 水泥” 的作用. 本次实验中所有需要计算梯度的地方都通过 PyTorch 中提供的自动微分算法实现, 需要使用的优化算法也由 PyTorch 框架提供. 我们将随机数种子设置为 100, 以保证实验结果可以复现.

我们使用 `sklearn` 中的 `datasets` 自动下载并加载数据, 之后将其转化为 `torch.tensor` 类型的数据, 以便后续使用 `torch` 的工具接口进行处理.

```
1 # 加载数据
2 data = datasets.load_wine()
3 data_features = torch.tensor(data.data)
4 data_labels = torch.tensor(data.target)[: ,None]
5 data_set = torch.cat([data_features, data_labels], dim=1)
```

我们使用 `torch` 中提供的划分数据集的工具, 将数据按照 6:4 随机划分为训练集和测试集.

```
1 length = data_set.shape[0]
2 alpha = 0.6
3 train_len = int(length*alpha)
4 test_len = length - train_len
5 train_data, test_data = torch.utils.data.random_split(data_set,
    , [train_len, test_len])
6 train_data = torch.tensor([item.numpy() for item in train_data
    ])
7 test_data = torch.tensor([item.numpy() for item in test_data])
```

2 实验结果分析

算法源码参见附录A, 本项目报告及实验结果开源: <https://github.com/WANGH950/Statistical-Machine-Learning/tree/main/2ND>.

2.1 朴素贝叶斯模型结果分析

葡萄酒品种分类数据集中的 13 个特征都为连续型数据, 因此我们在模型训练时保存了训练集特征的最小值和最大值, 并将所有特征值归一化到 $[-1, 1]$, 确保它们的尺度相同.

由于归一化参数仅通过训练集得到, 用这些参数对测试数据进行归一化后, 其可能会分布到 $[-1, 1]$ 以外, 因此我们将 $[-1, 1]$ 等距划分为 $S-2$ 个小区间, 并分别计算了各特征在 $(-\infty, -1), [1, +\infty), [x_k, x_{k+1}), x_k = -1 + \frac{2(k-1)}{S-2}, k = 1, 2, \dots, S-1$ 个小区间内的条件概率

$$P_{\lambda}(X^j \in S_{jk} | Y = c_k) = \frac{\sum_{i=1}^N I(x_j \in S_{jk}, y_i = c_k) + \lambda}{\sum_{i=1}^N I(y_i = c_k) + S\lambda} \quad (1)$$

和先验概率

$$P_{\lambda}(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k) + \lambda}{N + K\lambda}, \quad (2)$$

其中, $S_{jk} = [x_k, x_{k+1}), k = 1, \dots, S-1, S_{j0} = (-\infty, -1), S_{jS} = [1, +\infty), N$ 为训练样本数量, λ 为 Laplacian 平滑参数, K 为类别个数. 注意, 为方便起见, 我们将每个特征的区间划分数量设置为同一值 S .

Listing 1: 朴素贝叶斯模型训练算法部分源码

```

1  ...
2  for i in range(self.K):
3      labels_i = labels == self.full_labels[i]
4      self.pre_prob[i] = (labels_i.sum() + self.lamb) / (N + self
        .K*self.lamb)
5      for j in range(self.n):
6          self.cond_prob[i,j,0] = 1 / self.S
7          for k in range(self.S-1):
8              l = -1 + k * delta_x
9              r = 1 + delta_x
10             features_ij = features[labels_i[:,0],j]
11             features_ijk = features_ij[(features_ij>=l)*(
                features_ij<r)]

```

```

12         self.cond_prob[i,j,k+1] = (features_ijk.shape[0] +
13         self.lamb) / (labels_i.sum() + self.S*self.lamb)
...

```

我们实例化模型, 训练模型, 并在测试集上测试, 得到了 95.83% 的准确率.

2.2 Logistic 回归模型结果分析

我们采用多项逻辑斯谛回归模型

$$P(Y = k|x) = \frac{\exp(w_k \cdot x + b_k)}{1 + \sum_{i=1}^{K-1} \exp(w_i \cdot x + b_i)}, k = 1, 2, \dots, K-1$$

$$P(Y = K|x) = \frac{1}{1 + \sum_{i=1}^{K-1} \exp(w_i \cdot x + b_i)}.$$
(3)

以处理本文的三分类 ($K = 3$) 任务.

为学习模型参数, 我们实现了多类别的负对数似然损失函数

$$L(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(\pi_{\theta_k}(x_i)).$$
(4)

```

1  # 负对数似然损失函数
2  class NLLLoss(nn.Module):
3      def __init__(self) -> None:
4          super().__init__()
5
6      def forward(self, y_pre, y_rel):
7          assert y_pre.shape == y_rel.shape
8          loss = -torch.log(y_pre)*y_rel
9          return torch.sum(loss)

```

我们使用小批量梯度下降方法训练数据, 每次在训练数据中均匀随机采样 batch 条数据进行训练. 我们还在训练模型时将数据的标签进行 one-hot 编码, 以使用该方法.

```

1  index = torch.randint(0,N,[batch_size]) # 随机选取batch条数据
2  data_i = data_set[index,:]
3  x_i = data_i[:, :-1].to(torch.float32)
4  y_i = nn.functional.one_hot(data_i[:, -1].to(torch.int64),
    num_classes=model.class_num)

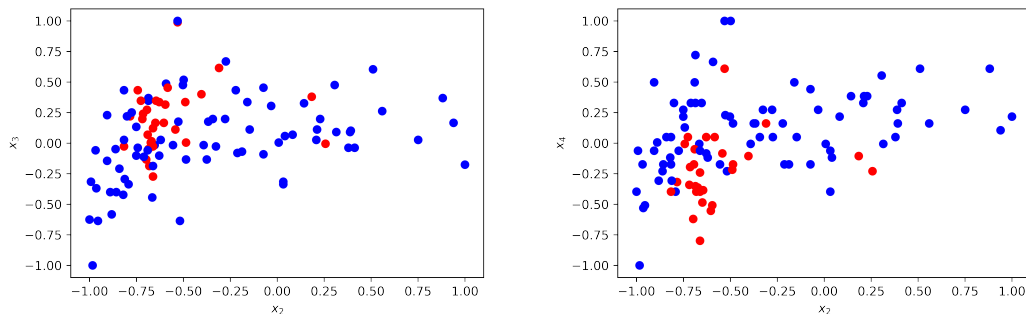
```

在数据经过线性层以前,我们先对其进行了批量归一化,自动学习相应的归一化参数.

我们实例化模型,选择批处理大小为 20,在 $1e-3$ 的学习率下在训练数据上迭代了 2,000 次,得到了最终的模型.我们使用训练好的模型在测试集上进行测试,得到了 97.22% 的准确率.该非线性模型效果略好于朴素贝叶斯模型.

2.3 支持向量机模型结果分析

我们分别实现了线性和非线性支持向量机模型.如图1,简单观察数据,我们可以得到第二个特征和第三个特征以及第二个特征和第四个特征的非线性可分关系.因此,针对非线性支持向量机模型,我们简单的考虑高斯核函数得到的特征



(a) 第二个特征和第三个特征关于标签 0 和非 0 绘制的散点图 (b) 第二个特征和第四个特征关于标签 0 和非 0 绘制的散点图

图 1: 非线性可分数据现象图像.

和原始特征结合,再经过线性分类器对模型进行分类.

我们实现了合页损失函数

$$L(\theta) = \sum_{i=1}^N [1 - y_i(w \cdot x_i + b)]_+ + \lambda \|w\|^2, \quad (5)$$

以用于优化模型,其中 $\theta = \{w, b\}$.

```

1 # 合页损失函数
2 class HingeLoss(nn.Module):
3     def __init__(self, lamb) -> None:
4         super(HingeLoss, self).__init__()
5         self.lamb = lamb
6
7     def forward(self, res_pre_linear_values, res_rel, parameters):

```

8

```

return torch.sum(torch.relu(1 - res_rel*
    res_pre_linear_values)) + self.lamb*torch.norm(
    parameters)**2

```

由于支持向量机是二分类模型, 我们将原始三分类问题分解为两个二分类问题, 并针对这两类二分类问题各训练了一个模型.

对于线性模型, 我们选取批处理大小 20, 正则化参数 $\lambda = 0.1$, 以 $1e-3$ 的学习率在训练数据上迭代了 2,000 次得到了最终的模型. 对于非线性模型, 我们选择核函数参数 $N_s = 20$, 其余参数与线性模型相同. 我们使用训练好的线性模型和非线性模型, 在测试数据上分别得到了 94.44%, 95.83% 的准确率. 结果表明, 非线性模型的效果略好于线性模型, 但都不如逻辑斯谛回归模型.

2.4 神经网络模型结果分析

多层感知机, 又名全连接前馈神经网络是现在深度学习模型中最基础也是最重要的模型之一. 其通过交替堆叠线性函数和非线性函数 (激活函数) 而达到对任意函数的万能逼近的作用. 全连接神经网络的结构如图2.

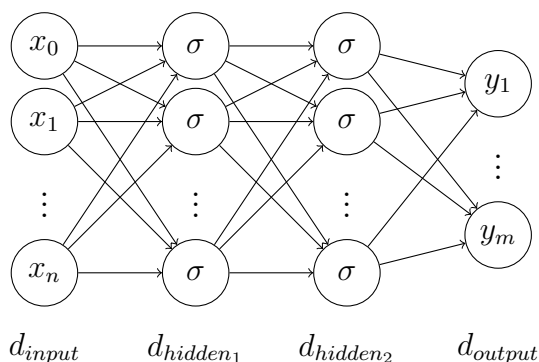


图 2: 两层隐藏层的全连接前馈神经网络结构, 其中带箭头的实线表示线性连接, 每条实线都对应一个权重, σ 是非线性函数, 为网络提供非线性因素.

我们搭建了以 Sigmoid Linear Unit(SiLU)

$$SiLU(x) = \frac{x \exp(x)}{1 + \exp(x)} \quad (6)$$

为激活函数的神经网络模型, 并在输出层使用 Softmax 函数

$$Softmax(y_i) = \frac{\exp(y_i)}{\sum_{i=1}^K \exp(y_i)} \quad (7)$$

归一化.

我们实例化了一个 3 层隐藏层, 每个隐藏层维度为 30 的全连接前馈神经网络. 该模型以 $1e-3$ 的学习率, 20 的批处理大小, 在训练数据上经过 3,000 次迭代, 得到了最终的模型. 模型在测试数据上得到了 97.22% 的准确率, 精度与逻辑斯谛回归模型相似, 要优于其他模型.

A 模型及算法源码

A.1 朴素贝叶斯模型及其训练算法

Listing 2: 朴素贝叶斯模型及其训练算法实现

```
1 # 定义模型
2 class NaiveBayes(nn.Module):
3     def __init__(self, n, full_labels, S, lamb) -> None:
4         super(NaiveBayes, self).__init__()
5         # 归一化参数
6         self.max = None
7         self.min = None
8         self.n = n # 特征数量
9         self.full_labels = full_labels # 所有标签
10        self.K = len(full_labels) # 标签数量
11        self.lamb = lamb # 贝叶斯估计参数 lambda
12        self.S = S # 每个特征分划区间数, 这里默认都为 S
13        self.cond_prob = torch.zeros([self.K, self.n, S]) # 条件
            概率
14        self.pre_prob = torch.zeros([self.K]) # 先验概率
15
16    def forward(self, features):
17        B, n = features.shape
18        assert n == self.n
19        post_prob = torch.ones([B, 1]) * self.pre_prob
20        # 归一化
21        features = (features - self.min) / (self.max - self.min
            ) * 2 - 1
22        delta_x = 2 / (self.S - 2)
23        for i in range(B):
24            for j in range(self.K):
25                for k in range(n):
26                    # (-\infty, -1) 和 (1, \infty) 的概率
27                    if features[i, k] < -1: post_prob[i, j] *=
                        self.cond_prob[j, k, 0]
28                    elif features[i, k] >= 1: post_prob[i, j] *=
                        self.cond_prob[j, k, -1]
```

```

29         else:
30             for h in range(self.S-2):
31                 l = -1 + h * delta_x
32                 r = l + delta_x
33                 if features[i,k] >= l and features[
34                     i,k] < r:
35                     post_prob[i,j] *= self.
36                         cond_prob[j,k,h+1]
37                     break
38     return self.full_labels[torch.argmax(post_prob,dim=1)]
39
40 def fit(self, train_data):
41     # 计算先验概率
42     N,_ = train_data.shape
43     self.max = torch.max(train_data[:,:,:-1],dim=0).values
44     self.min = torch.min(train_data[:,:,:-1],dim=0).values
45     train_data[:,:,:-1] = (train_data[:,:,:-1] - self.min) / (
46         self.max - self.min) * 2 - 1
47     features = train_data[:,:,:-1] # 特征
48     labels = train_data[:,-1:].int() # 标签
49     delta_x = 2 / (self.S - 2)
50     for i in range(self.K):
51         labels_i = labels == self.full_labels[i]
52         self.pre_prob[i] = (labels_i.sum() + self.lamb) / (
53             N + self.K*self.lamb)
54         for j in range(self.n):
55             self.cond_prob[i,j,0] = 1 / self.S
56             for k in range(self.S-1):
57                 l = -1 + k * delta_x
58                 r = l + delta_x
59                 features_ij = features[labels_i[:],0],j]
60                 features_ijk = features_ij[(features_ij>=1)
61                     *(features_ij<r)]
62                 self.cond_prob[i,j,k+1] = (features_ijk.
63                     shape[0] + self.lamb) / (labels_i.sum()
64                     + self.S*self.lamb)
65     return self.pre_prob, self.cond_prob

```

A.2 Logistic 模型及相应学习算法

Listing 3: 逻辑斯谛回归模型及损失函数和训练方法

```
1 # 定义模型
2 class Logistic(nn.Module):
3     def __init__(self, feature_num, class_num) -> None:
4         super(Logistic, self).__init__()
5         self.feature_num = feature_num
6         self.class_num = class_num
7         self.linear = nn.Sequential(
8             nn.BatchNorm1d(feature_num), # 批量归一化
9             nn.Linear(feature_num, class_num-1, bias=False)
10        )
11
12    def forward(self, x):
13        B, d = x.shape
14        assert d == self.feature_num and B > 0
15        y = torch.cat([torch.exp(self.linear(x)), torch.ones([B
16            , 1])], dim=1)
17        y = y / torch.sum(y, dim=1, keepdim=True)
18        return y
19
20 # 负对数似然损失函数
21 class NLLLoss(nn.Module):
22     def __init__(self) -> None:
23         super().__init__()
24
25     def forward(self, y_pre, y_rel):
26         assert y_pre.shape == y_rel.shape
27         loss = -torch.log(y_pre)*y_rel
28         return torch.sum(loss)
29
30 # 定义模型训练函数
31 def train(model, data_set, batch_size, epoch = 1000,
32         learning_rate = 1e-3):
33     N, _ = data_set.shape
34     criterion = NLLLoss()
```

```
33     optim = torch.optim.SGD(model.parameters(), learning_rate,
34                               momentum=0) # 动量设置为0
35
36     start = time.time()
37     loss_values = torch.zeros(epoch)
38     for i in range(epoch):
39         model.train()
40         optim.zero_grad()
41         index = torch.randint(0, N, [batch_size]) # 随机选取batch
42             条数据
43         data_i = data_set[index, :]
44         x_i = data_i[:, :-1].to(torch.float32)
45         y_i = nn.functional.one_hot(data_i[:, -1].to(torch.int64),
46                                     num_classes=model.class_num)
47         outputs = model(x_i)
48         loss = criterion(outputs, y_i)
49         loss.backward()
50         optim.step()
51
52         model.eval()
53         loss_values[i] = loss.item()
54
55         print('\r%5d/{}|{}{}|{:.2f}s   [Loss: %e]'.format(
56             epoch,
57             "#" * int((i+1)/epoch*50),
58             " " * (50 - int((i+1)/epoch*50)),
59             time.time() - start) %
60             (i+1,
61              loss_values[i]), end = ' ', flush=True)
61     print("\nTraining has been completed.")
62     return loss_values
```

A.3 支持向量机

Listing 4: 线性支持向量机

```
1 # 定义模型
2 class SVMLinear(nn.Module):
```

```
3     def __init__(self, feature_num) -> None:
4         super(SVMLLinear, self).__init__()
5         self.feature_num = feature_num
6         self.batch_norm = nn.BatchNorm1d(feature_num) # 特征批
           量归一化, 学习归一化参数
7         self.linear = nn.Linear(feature_num, 1)
8
9     def forward(self, x, signed = True):
10        B, d = x.shape
11        assert B > 0 and d == self.feature_num
12        normed = self.batch_norm(x)
13        outputs = self.linear(normed)
14        if signed:
15            return torch.sign(outputs)
16        else:
17            return outputs
18
19 # 合页损失函数
20 class HingeLoss(nn.Module):
21     def __init__(self, lamb) -> None:
22         super(HingeLoss, self).__init__()
23         self.lamb = lamb
24
25     def forward(self, res_pre_linear_values, res_rel, parameters):
26         return torch.sum(torch.relu(1 - res_rel *
           res_pre_linear_values)) + self.lamb * torch.norm(
           parameters) ** 2
27
28 # 定义训练函数
29 def train_svm(model, data_set, batch_size, lamb=1, epoch =
   1000, learning_rate = 1e-3):
30     N, _ = data_set.shape
31     criterion = HingeLoss(lamb)
32     optim = torch.optim.SGD(model.parameters(), learning_rate,
           momentum=0) # 动量设置为0
33
34     start = time.time()
```

```
35     loss_values = torch.zeros(epoch)
36     for i in range(epoch):
37         model.train()
38         optim.zero_grad()
39         index = torch.randint(0,N,[batch_size]) # 随机选取batch
           条数据
40         data_i = data_set[index,:]
41         x_i = data_i[:, :-1].to(torch.float32)
42         y_i = data_i[:, -1:].to(torch.int32)
43         outputs = model(x_i, signed=False)
44         loss = criterion(outputs, y_i, model.linear.weight)
45         loss.backward()
46         optim.step()
47
48         model.eval()
49         loss_values[i] = loss.item()
50
51         print('\r%5d/{}|{}{}|{: .2f}s   [Loss: %e]'.format(
52             epoch,
53             "#"*int((i+1)/epoch*50),
54             " "*(50-int((i+1)/epoch*50)),
55             time.time() - start) %
56             (i+1,
57             loss_values[i]), end = ' ', flush=True)
58     print("\nTraining has been completed.")
59     return loss_values
```

Listing 5: 非线性支持向量机

```
1 # 高斯核非线性SVM
2 class SVMNonLinear(nn.Module):
3     def __init__(self, feature_num, kernel_data) -> None:
4         super(SVMNonLinear, self).__init__()
5         self.feature_num = feature_num
6         self.kernel_features = kernel_data[:, :-1].to(torch.
           float32)
7         self.kernel_labels = kernel_data[:, -1].to(torch.int32)
8         self.batch_norm = nn.BatchNorm1d(feature_num)
```

```
9         self.linear = nn.Linear(feature_num+kernel_data.shape
10                                   [0],1)
11
12         self.kernel_params = nn.Parameter(torch.rand(1)+5,
13                                             requires_grad=True) # 核函数标准差
14
15     def forward(self, x, signed = True):
16         B,d = x.shape
17         assert B > 0 and d == self.feature_num
18         # 归一化特征
19         normed_data = self.batch_norm(x)
20         normed_kernel_data = self.batch_norm(self.
21                                             kernel_features)
22         kernel_features = torch.exp(-torch.mm(normed_data,
23                                             normed_kernel_data.T) / self.kernel_params**2) *
24         self.kernel_labels
25         outputs = self.linear(torch.cat([normed_data,
26                                         kernel_features],dim=1)) # 核方法特征和原特征结合
27         if signed:
28             return torch.sign(outputs)
29         else:
30             return outputs
```

A.4 全连接神经网络

Listing 6: 全连接神经网络模型实现

```
1 class MLP(nn.Module):
2     def __init__(self,feature_num,class_num,hidden_dim=20,layer_num
3                   =2) -> None:
4         super(MLP,self).__init__()
5         self.feature_num = feature_num
6         self.class_num = class_num
7         self.hidden_dim = hidden_dim
8         self.layer_num = layer_num
9         self.model = nn.Sequential(
10             OrderedDict(
11                 [("input_layer",
12                  nn.Sequential(
```



```
12         nn.BatchNorm1d(feature_num),
13         nn.Linear(feature_num,hidden_dim),
14         nn.Tanh()
15     ))] +
16     [("hidden_layer_"+str(i+1),
17      nn.Sequential(
18          nn.Linear(hidden_dim,hidden_dim),
19          nn.Tanh()
20      )) for i in range(layer_num-1)] +
21     [("output_layer",
22      nn.Sequential(
23          nn.Linear(hidden_dim,class_num),
24          nn.Softmax(dim=1)
25      ))]
26     )
27 )
28
29 # 前向传播
30 def forward(self,x):
31     return self.model(x)
```