# 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 11 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, 由于某些特征是连续型数据, 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 实验代码

## 1.1 感知机学习算法

Listing 1: 感知机原型算法实现

```python
class Perception():
    def __init__(self, dim) -> None:
        # 构造函数
        # dim:  特征维度
        # w:    权重
        # b:    偏置项

        self.dim = dim
        self.w = np.zeros([dim])
        self.b = 0

    def train(self, data_set, epoch, learning_rate):
        # 训练模型

        for i in range(epoch):
            for (x,y) in data_set:
                if self.predict(x)*y <= 0:
                    self.w = self.w + learning_rate*x*y
                    self.b = self.b + learning_rate*y
            # 计算准确率
            acc = self.accuracy(data_set)
            print('epoch: ',i+1, 'accuracy: ', acc)
            # 早停条件
            if acc == 1:
                break
        print('Trining complete.')

    def predict(self, x):
        # 预测

        return np.sign(np.dot(self.w,x) + self.b)
```

```
33      def accuracy(self, data_set):
34          # 计算精度
35
36          acc = 0
37          for (x,y) in data_set:
38              if self.predict(x)*y > 0:
39                  acc += 1
40          return acc/len(data_set)
```

Listing 2: 感知机对偶算法实现

```
1   class PerceptionDual():
2       def __init__(self, data) -> None:
3           # 构造函数
4           # data:    训练数据
5
6           self.data = data
7           self.N = len(data)
8           self.x = np.array([xx for (xx,_) in data])
9           self.y = np.array([yy for (_,yy) in data])
10          self.alpha = np.zeros([self.N])
11          self.b = 0
12
13      def train(self, epoch, learning_rate):
14          # 训练模型
15
16          for i in range(epoch):
17              for j in range(len(self.data)):
18                  if self.predict(self.data[j][0])*self.data[j
                        ][1] <= 0:
19                      self.alpha[j] = self.alpha[j] +
                            learning_rate
20                      self.b = self.b + learning_rate*self.data[j
                            ][1]
21              # 计算准确率
22              acc = self.accuracy(self.data)
23              print('epoch: ',i+1, 'accuracy: ', acc)
24              # 早停条件
```

```
25          if acc == 1:
26              break
27      print('Trining complete.')
28
29  def predict(self, x):
30      # 预测
31      return np.sign(np.dot(np.dot(self.x,x),self.alpha*self.
            y) + self.b)
32
33  def accuracy(self, data_set):
34      # 计算精度
35
36      acc = 0
37      for (x,y) in data_set:
38          if self.predict(x)*y > 0:
39              acc += 1
40      return acc/len(data_set)
```

## 1.2　k-近邻算法

Listing 3: kd-树构造算法和搜索算法实现

```
1      class Node():
2  def __init__(self, value, data, label) -> None:
3      # 构造函数
4      # value:        节点的划分超平面参数
5      # data:         落在超平面上的数据点
6      # label:        落在超平面上的数据点对应的标签
7
8      self.value = value
9      self.data = data
10     self.label = label
11     self.left = None
12     self.right = None
13
14  def set_left(self, node):
15      # 设置左子节点
16
```

```
17            if node != None:
18                self.left = node
19
20        def set_right(self, node):
21            # 设置右子节点
22
23            if node != None:
24                self.right = node
25
26  class KDTree():
27        def __init__(self) -> None:
28            # 构造函数
29            # 用于存储KDTree，支持直接实例化对象时直接输入一个kd-树
30            self.root = None
31
32        def create(self, data, label, j = 0):
33            # 递归构造平衡KD树
34
35            num, k = data.shape
36            if num == 0:
37                return None
38            else:
39                l = j % k
40                ind_sorted = np.argsort(data[:,l])
41                ind_median = ind_sorted[num//2]
42                value_ = int(np.median(data[ind_median,l]))
43                data_ = data[data[:,l]==value_]
44                label_ = label[data[:,l]==value_]
45                node = Node(
46                    value=value_,
47                    data=data_,
48                    label=label_
49                )
50                node.set_left(
51                    self.create(
52                        data=data[data[:,l]<value_],
53                        label=label[data[:,l]<value_],
```

```
54                    j=j+1
55                )
56            )
57            node.set_right(
58                self.create(
59                    data=data[data[:,l]>value_],
60                    label=label[data[:,l]>value_],
61                    j=j+1
62                )
63            )
64            if j == 0:
65                self.root = node
66            else:
67                return node
68
69    def search(self, x, j = 0, node = None):
70        # 递归搜索KD树
71
72        if self.root == None:
73            print("You haven't created a KDTree yet.")
74            return None
75        if j == 0:
76            node = self.root
77        k = x.shape[0]
78        l = j % k
79        # 叶子节点停止条件
80        if self.is_leaf(node):
81            distance = np.linalg.norm(x-node.data,2,1)
82            index = np.argmin(distance)
83            return node.data[index], node.label[index]
84        else:
85            # 计算当前节点中的最近数据点
86            distance = np.linalg.norm(x-node.data,2,1)
87            min_distance = np.min(distance)
88            index = np.argmin(distance)
89            nearest = node.data[index]
90            label = node.label[index]
```

```
91              # 递归计算子节点的最近数据点, 并比较
92              if x[l] < node.value and node.left != None:
93                  nearest_, label_ = self.search(
94                      x = x,
95                      j = j+1,
96                      node = node.left
97                  )
98                  if np.linalg.norm(x-nearest_,2) < min_distance:
99                      nearest = nearest_
100                     label = label_
101             elif x[l] > node.value and node.right != None:
102                 nearest_, label_ = self.search(
103                     x = x,
104                     j = j+1,
105                     node = node.right
106                 )
107                 if np.linalg.norm(x-nearest_,2) < min_distance:
108                     nearest = nearest_
109                     label = label_
110             return nearest, label
111
112     def is_leaf(self, node: Node):
113         # 判断是否是叶子节点
114
115         if node.left != None or node.right != None:
116             return False
117         else:
118             return True
```
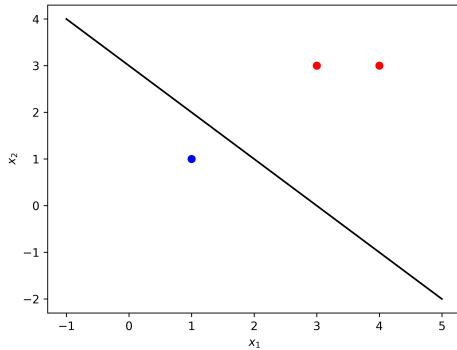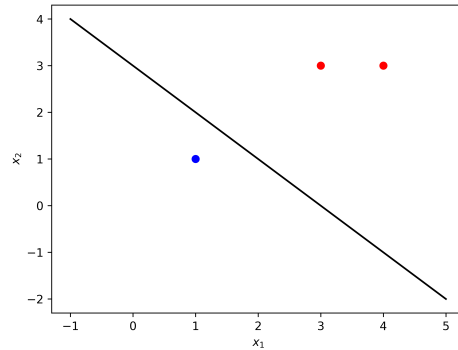
# 2　实验结果分析

　　MNIST 数据集通过 torchvision 加载, 所有算法基于 python 的 numpy 实现, 数据预处理的一部分地方借助 pytorch 实现. 细节请参考:
https://github.com/WANGH950/Statistical-Machine-Learning/tree/main/1ST

## 2.1　感知机学习算法结果分析



(a) 感知机原型算法训练结果　　　　　　(b) 感知机对偶算法训练结果

图 1: 感知机模型训练结果, 其中红点表示正类, 蓝点表示负类, 黑色实线表示感知机分割超平面.

　　以 0.01 的学习率, 感知机原型算法在题目给定的数据上, 经过 5 次训练后收敛, 得到结果. 以 0.01 的学习率, 感知机对偶算法在题目给定的数据上, 经过 5 次训练后收敛, 得到结果. (图 1)

Listing 4: MNIST 数据加载

```
import torchvision
data=torchvision.datasets.MNIST(
    root='MNIST',
    train=True,
    transform=torchvision.transforms.ToTensor(),
    download=True
)
```

Listing 5: MNIST 数据预处理

```
train_data = data.train_data
train_label = data.train_labels
# 转化为二分类
data_set = []
for i in range(train_data.shape[0]):
    if train_label[i] < 5:
        y = 1
    else:
```

```
 9          y = -1
10      data_set.append(
11          (train_data[i].reshape([28**2]).numpy()/255,y)
12      )
```

<div align="center">Listing 6: 使用对偶感知机算法处理 MNIST 数据</div>

```
1  model = PerceptionDual(data_set[:1000])
2  model.train(
3      epoch=30,
4      learning_rate=0.001
5  )
```

如 Listing 6 所示, 我们使用 1,000 条数据, 在 0.001 的学习率下训练了 30 次, 最后在训练数据上得到了 0.923 的准确率.

## 2.2　k-近邻算法结果分析

<div align="center">Listing 7: MNIST 数据预处理</div>

```
1  data_num,dimx,dimy = train_data.shape
2  data = train_data.reshape([data_num,dimx*dimy])+torch.randint
       (0,10,[data_num,dimx*dimy])
3  label = train_label.unsqueeze(-1)
```

如 Listing 7 所示, 由于手写数字是稀疏矩阵, 我们通过对其添加随机噪声以保证构造的 kd-树是好的 (不加噪声会导致数据都分布在一个节点上). 这里, 我们添加 $U(0, 10)$ 的均匀整数噪声.

<div align="center">Listing 8: 构造 kd-树</div>

```
1  model = KDTree()
2  tree = model.create(
3      data=data[:30000],
4      label=label[:30000]
5  )
```

如 Listing 8 所示, 我们使用前 30,000 条数据构造 kd-树. 我们选取第 40,000 条数据进行预测, 图 1 展示了预测结果和真实结果.

同时, 我们还使用后 30,000 条数据作为测试集进行测试 (listing 9), 得到了 93.223% 的准确率. 结果表明, 我们的算法实现准确.
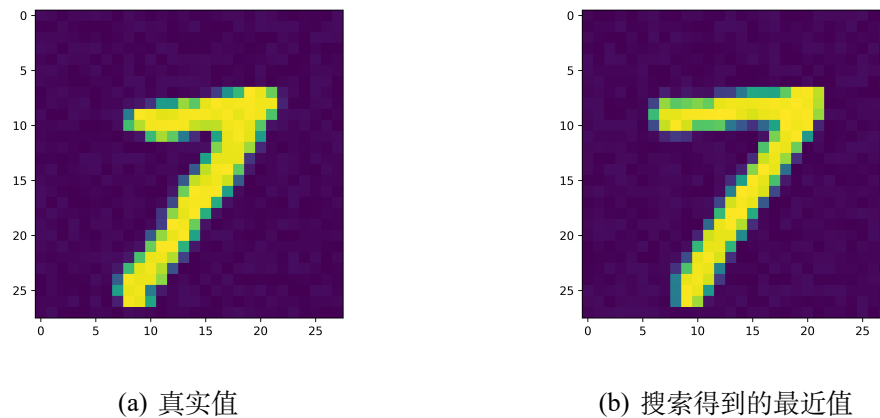
(a) 真实值

(b) 搜索得到的最近值

图 2: kd-树搜索结果.

Listing 9: 使用 kd-树分类测试数据

```
1  k = 30000
2  acc = 0
3  data_test = data[-k:]
4  label_test = label[-k:]
5  for i in range(k):
6      data_rel = data_test[i]
7      label_rel = label_test[i]
8      data_pre,label_pre = model.search(data_rel)
9      if label_pre == label_rel:
10         acc += 1
11 acc /= k
```

所有实验结果和运行效率以源码为准:

https://github.com/WANGH950/Statistical-Machine-Learning/tree/main/1ST

# A 模型及算法源码

## A.1 朴素贝叶斯模型及其训练算法

Listing 10: 朴素贝叶斯模型及其训练算法实现

```python
# 定义模型
class NaiveBayes(nn.Module):
    def __init__(self, n, full_labels, S, lamb) -> None:
        super(NaiveBayes,self).__init__()
        # 归一化参数
        self.max = None
        self.min = None
        self.n = n # 特征数量
        self.full_labels = full_labels # 所有标签
        self.K = len(full_labels) # 标签数量
        self.lamb = lamb # 贝叶斯估计参数lambda
        self.S = S # 每个特征分划区间数，这里默认都为S
        self.cond_prob = torch.zeros([self.K,self.n,S]) # 条件
            概率
        self.pre_prob = torch.zeros([self.K]) # 先验概率

    def forward(self, features):
        B,n = features.shape
        assert n == self.n
        post_prob = torch.ones([B,1])*self.pre_prob
        # 归一化
        features = (features - self.min) / (self.max - self.min
            ) * 2 - 1
        delta_x = 2 / (self.S - 2)
        for i in range(B):
            for j in range(self.K):
                for k in range(n):
                    # (-\infty,-1)和(1,\infty)的概率
                    if features[i,k] < -1: post_prob[i,j] *=
                        self.cond_prob[j,k,0]
                    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[
                                        i,k] < r:
34                                      post_prob[i,j] *= self.
                                            cond_prob[j,k,h+1]
35                                      break
36          return self.full_labels[torch.argmax(post_prob,dim=1)]
37
38      def fit(self, train_data):
39          # 计算先验概率
40          N,_ = train_data.shape
41          self.max = torch.max(train_data[:,:-1],dim=0).values
42          self.min = torch.min(train_data[:,:-1],dim=0).values
43          train_data[:,:-1] = (train_data[:,:-1] - self.min) / (
                self.max - self.min) * 2 - 1
44          features = train_data[:,:-1] # 特征
45          labels = train_data[:,-1:].int() # 标签
46          delta_x = 2 / (self.S - 2)
47          for i in range(self.K):
48              labels_i = labels == self.full_labels[i]
49              self.pre_prob[i] = (labels_i.sum() + self.lamb) / (
                    N + self.K*self.lamb)
50              for j in range(self.n):
51                  self.cond_prob[i,j,0] = 1 / self.S
52                  for k in range(self.S-1):
53                      l = -1 + k * delta_x
54                      r = l + delta_x
55                      features_ij = features[labels_i[:,0],j]
56                      features_ijk = features_ij[(features_ij>=l)
                            *(features_ij<r)]
57                      self.cond_prob[i,j,k+1] = (features_ijk.
                            shape[0] + self.lamb) / (labels_i.sum()
                            + self.S*self.lamb)
58          return self.pre_prob, self.cond_prob
```

## A.2 Logistic 模型及相应学习算法

Listing 11: 逻辑斯谛回归模型及损失函数和训练方法

```python
# 定义模型
class Logistic(nn.Module):
    def __init__(self,feature_num,class_num) -> None:
        super(Logistic,self).__init__()
        self.feature_num = feature_num
        self.class_num = class_num
        self.linear = nn.Sequential(
            nn.BatchNorm1d(feature_num), # 批量归一化
            nn.Linear(feature_num,class_num-1,bias=False)
        )

    def forward(self, x):
        B,d = x.shape
        assert d == self.feature_num and B > 0
        y = torch.cat([torch.exp(self.linear(x)),torch.ones([B
            ,1])],dim=1)
        y = y / torch.sum(y,dim=1,keepdim=True)
        return y

# 负对数似然损失函数
class NLLLoss(nn.Module):
    def __init__(self) -> None:
        super().__init__()

    def forward(self,y_pre,y_rel):
        assert y_pre.shape == y_rel.shape
        loss = -torch.log(y_pre)*y_rel
        return torch.sum(loss)

# 定义模型训练函数
def train(model, data_set, batch_size, epoch = 1000,
    learning_rate = 1e-3):
    N,_ = data_set.shape
    criterion = NLLLoss()
```

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

## A.3  支持向量机

Listing 12: 线性支持向量机

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

```python
3       def __init__(self,feature_num) -> None:
4           super(SVMLinear,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()
```

14

```python
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 13: 非线性支持向量机

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

## A.4　全连接神经网络

Listing 14: 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
   =2) -> None:
 3     super(MLP,self).__init__()
 4     self.feature_num = feature_num
 5     self.class_num = class_num
 6     self.hidden_dim = hidden_dim
 7     self.layer_num = layer_num
 8     self.model = nn.Sequential(
 9         OrderedDict(
10             [("input_layer",
11                 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)
```