



東南大學  
SOUTHEAST UNIVERSITY

## 模式识别实验报告

专业: 人工智能

学号: 58122204

年级: 大二

姓名: 谢兴

签名:

时间:

# 目录

<b>1</b>	<b>实验一 KNN Classification</b>	<b>4</b>
1.1	问题描述 . . . . .	4
1.2	概述 . . . . .	4
1.3	任务说明 . . . . .	4
1.4	实现步骤与流程 . . . . .	4
1.4.1	实验思路 . . . . .	4
1.4.2	数学模型 . . . . .	5
1.4.3	关键难点 . . . . .	6
1.4.4	算法描述 . . . . .	6
1.4.5	马氏距离梯度计算公式推导 . . . . .	8
1.5	实验结果与分析 . . . . .	9
1.5.1	数据集的部分可视化分析 . . . . .	9
1.5.2	实验结果的分析 . . . . .	9
1.5.3	手动实现算法的评价 . . . . .	11
1.6	MindSpore 学习使用心得体会 . . . . .	14
1.7	代码附录（数据加载可视化展示部分，具体见 knn.ipynb 文件） . . . . .	15
1.7.1	knn.ipynb . . . . .	15
1.7.2	knn_mindspore.ipynb . . . . .	22
<b>2</b>	<b>实验二 Naïve Bayes Classification</b>	<b>24</b>
2.1	问题描述 . . . . .	24
2.1.1	概述 . . . . .	24
2.1.2	任务说明 . . . . .	24
2.2	实现步骤与流程 . . . . .	25

2.2.1	实验思路 . . . . .	25
2.2.2	数学模型 . . . . .	25
2.2.3	关键难点 . . . . .	26
2.2.4	算法描述 . . . . .	27
2.3	实验结果与分析 . . . . .	28
2.3.1	手动实现朴素贝叶斯算法 . . . . .	28
2.3.2	mindSpore 实现朴素贝叶斯算法 . . . . .	28
2.3.3	手动实现算法的评价 . . . . .	29
2.4	MindSpore 学习使用心得体会 . . . . .	29
2.5	代码附录 . . . . .	29
2.5.1	naive_bayes.ipynb . . . . .	29
2.5.2	naive_bayes_mindspore.ipynb . . . . .	32
<b>3</b>	<b>实验三 Neural Network Image Classification</b>	<b>36</b>
3.1	问题描述 . . . . .	36
3.2	概述 . . . . .	36
3.3	任务说明 . . . . .	36
3.4	实现步骤与流程 . . . . .	36
3.4.1	实验环境 . . . . .	36
3.4.2	实验思路 . . . . .	36
3.4.3	数学模型 . . . . .	38
3.4.4	关键难点 . . . . .	40
3.4.5	算法描述 . . . . .	40
3.5	实验结果与分析 . . . . .	40
3.5.1	实验结果展示 . . . . .	40
3.5.2	进一步探究 . . . . .	40

3.5.3	手动实现算法的评价	42
3.5.4	手动实现 BP 神经网络算法和基于 Mindspore 实现神经网络算法的对比	42
3.6	MindSpore 学习使用心得体会	47
3.7	代码附录	48
3.7.1	cifar-10-scratch.ipynb	48
3.7.2	bp_mindspore.ipynb	53
<b>4</b>	<b>心得体会</b>	<b>59</b>

# 1 实验一 KNN Classification

## 1.1 问题描述

## 1.2 概述

利用 KNN 算法，对 Iris 鸢尾花数据集中的测试集进行分类。

## 1.3 任务说明

1. 利用欧式距离作为 KNN 算法的度量函数，对测试集进行分类。实验报告中，要求在验证集上分析近邻数  $k$  对 KNN 算法分类精度的影响。
2. 利用马氏距离作为 KNN 算法的度量函数，对测试集进行分类。
3. 基于 MindSpore 平台提供的官方模型库，对相同的数据集进行训练，并与自己独立实现的算法对比结果（包括但不限于准确率、算法迭代收敛次数等指标），并分析结果中出现差异的可能原因，给出使用 MindSpore 的心得和建议。
- 4.（加分项）使用 MindSpore 平台提供的相似任务数据集（例如，其他的分类任务数据集）测试自己独立实现的算法并与 MindSpore 平台上的官方实现算法进行对比，并进一步分析差异及其成因。

## 1.4 实现步骤与流程

### 1.4.1 实验思路

1. 导入必要的库，包括 `numpy`, `pandas`, `matplotlib`, `plotly` 和 `seaborn`;
2. 数据加载和基本信息显示;
  - (a) 从 `data/train.csv` 文件中加载训练数据集
  - (b) 显示数据集的前几行数据

(c) 显示数据集的描述性统计信息

(d) 显示数据集的基本信息，包括数据类型和缺失值情况

3. 数据可视化；

(a) 使用 `seaborn` 绘制数据集的特征两两关系图，并按标签着色

(b) 使用 `plotly` 绘制数据分布的饼图

(c) 分别绘制每个特征（萼片长度、萼片宽度、花瓣长度、花瓣宽度）的箱线图和直方图

4. 实现基于 Euclidean 距离和基于 Mahalanobis 距离的 KNN 算法；

5. 数据处理和预测；

(a) 加载测试数据集并进行必要的类型转换和缺失值检查

(b) 使用预训练模型对测试数据进行预测，并将预测结果保存到 CSV 文件中

6. 最后对比欧式距离和马氏距离两种度量方式的分类效果和差异

(a) 比较多个预测结果文件 `task1_test_prediction.csv` 和 `task2_test_prediction.csv` 之间的差异，找出不同的行和列，并打印出不同值的位置

### 1.4.2 数学模型

KNN 算法的数学模型如下：

给定一个测试样本  $x$ ，KNN 算法通过计算  $x$  与训练集中所有样本之间的距离（常用欧氏距离），选择距离最近的  $k$  个样本，然后通过多数投票法决定  $x$  的类别。欧氏距离的计算公式为：

$$d(x_1, x_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

马氏距离的计算公式为：

$$d(x_1, x_2) = \sqrt{(x_1 - x_2)^T \Sigma^{-1} (x_1 - x_2)}$$

其中， $\Sigma$  为协方差矩阵。

### 1.4.3 关键难点

1. 如何高效地计算欧氏距离。
2. 如何在较大的数据集上进行快速的邻居搜索。

### 1.4.4 算法描述

基于 Euclidean 距离和 Mahalanobis 距离的 KNN 算法的伪代码分别见算法 1 和算法 2。

---

**Algorithm 1** K-Nearest Neighbors Based on Euclidean Distance

---

- 1: 初始化 KNN 分类器，邻居数为  $k$
  - 2: **procedure** 拟合 ( $X_{\text{train}}, y_{\text{train}}$ )
  - 3:     存储训练数据和标签
  - 4: **end procedure**
  - 5: **procedure** 预测 ( $X$ )
  - 6:     **for** 每个测试数据  $x$  **do**
  - 7:         计算  $x$  与所有训练样本之间的欧氏距离
  - 8:         对距离进行排序，选择最近的  $k$  个邻居
  - 9:         对这  $k$  个邻居的标签进行多数投票
  - 10:        将多数投票结果赋予  $x$
  - 11:     **end for**
  - 12:     **return** 预测的标签
  - 13: **end procedure**
-

---

**Algorithm 2** K-Nearest Neighbors Based on Mahalanobis Distance

---

```
1: 初始化 KNN 分类器, 邻居数为  $k$ , 矩阵  $A$  的维度为  $e$ , 学习率为  $\eta$ , 最大迭代次数为  $max\_iter$ 
2: procedure 拟合 ( $X\_train, y\_train$ )
3:   存储训练数据和标签
4:   初始化矩阵  $A$  为随机值
5:   for 迭代次数  $iteration = 1, 2, \dots, max\_iter$  do
6:     初始化梯度矩阵  $\nabla A$  为零
7:     for 每个训练样本  $x_i$  do
8:       获取与  $x_i$  同类的样本索引  $same\_class\_indices$ 
9:       for 每个同类样本  $x_j$  do
10:        if  $i == j$  then
11:          跳过
12:        end if
13:        计算  $p_{ij}$  值
14:        计算样本差异  $diff = x_i - x_j$ 
15:        更新梯度  $\nabla A += 2 \cdot p_{ij} \cdot (A \cdot diff) \cdot diff^T$ 
16:      end for
17:    end for
18:    按学习率更新矩阵  $A$ :  $A = A - \eta \cdot \nabla A / n$ 
19:  end for
20: end procedure
21: procedure 预测 ( $X$ )
22:   for 每个测试数据  $x$  do
23:     计算  $x$  与所有训练样本之间的马氏距离
24:     对距离进行排序, 选择最近的  $k$  个邻居
25:     对这  $k$  个邻居的标签进行多数投票
26:     将多数投票结果赋予  $x$ 
27:   end for
28:   return 预测的标签
29: end procedure
```

---



### 1.4.5 马氏距离梯度计算公式推导

假设我们有训练数据集  $\{(x_i, y_i)\}_{i=1}^n$ ，其中  $x_i \in \mathbb{R}^d$  为样本特征， $y_i$  为样本类别。为了优化马氏距离下的 KNN 算法，我们需要学习一个矩阵  $A \in \mathbb{R}^{e \times d}$ ，使得同类样本之间的距离最小化。马氏距离的计算公式为：

$$d_M(x_i, x_j) = \sqrt{(x_i - x_j)^\top A^\top A (x_i - x_j)} \quad (1)$$

为了优化矩阵  $A$ ，我们使用如下的目标函数：

$$\mathcal{L} = \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} p_{ij} \cdot d_M^2(x_i, x_j) \quad (2)$$

其中， $\mathcal{N}(i)$  表示与  $x_i$  同类的样本索引集合， $p_{ij}$  为权重，定义为：

$$p_{ij} = \frac{\exp(-d_M^2(x_i, x_j))}{\sum_{k \in \mathcal{N}(i)} \exp(-d_M^2(x_i, x_k))} \quad (3)$$

首先，我们对  $d_M^2(x_i, x_j)$  进行展开：

$$d_M^2(x_i, x_j) = (x_i - x_j)^\top A^\top A (x_i - x_j) \quad (4)$$

为了计算梯度  $\nabla_A \mathcal{L}$ ，我们需要对  $\mathcal{L}$  关于  $A$  求导：

$$\mathcal{L} = \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} p_{ij} (x_i - x_j)^\top A^\top A (x_i - x_j) \quad (5)$$

对  $A$  求导时，需要使用链式法则：

$$\frac{\partial \mathcal{L}}{\partial A} = \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} \left( \frac{\partial p_{ij}}{\partial A} (x_i - x_j)^\top A^\top A (x_i - x_j) + p_{ij} \frac{\partial ((x_i - x_j)^\top A^\top A (x_i - x_j))}{\partial A} \right) \quad (6)$$

首先计算  $p_{ij}$  对  $A$  的导数。由于  $p_{ij}$  包含在指数函数内，我们得到：

$$\frac{\partial p_{ij}}{\partial A} = p_{ij} \left( - \sum_{k \in \mathcal{N}(i)} p_{ik} \cdot 2(x_i - x_k)^\top A^\top \cdot (x_i - x_k) + 2(x_i - x_j)^\top A^\top \cdot (x_i - x_j) \right) \quad (7)$$

然后计算  $(x_i - x_j)^\top A^\top A(x_i - x_j)$  对  $A$  的导数：

$$\frac{\partial((x_i - x_j)^\top A^\top A(x_i - x_j))}{\partial A} = 2A(x_i - x_j)(x_i - x_j)^\top \quad (8)$$

将以上结果代入梯度公式中，我们得到：

$$\begin{aligned} \nabla_A \mathcal{L} = & \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} \left[ p_{ij} \left( - \sum_{k \in \mathcal{N}(i)} p_{ik} \cdot 2(x_i - x_k)^\top A^\top \cdot (x_i - x_k) \right. \right. \\ & \left. \left. + 2(x_i - x_j)^\top A^\top \cdot (x_i - x_j) \right) + 2p_{ij} A(x_i - x_j)(x_i - x_j)^\top \right] \end{aligned} \quad (9)$$

整理后得到最终的梯度公式：

$$\nabla_A \mathcal{L} = 2 \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} p_{ij} \left[ A(x_i - x_j)(x_i - x_j)^\top - \sum_{k \in \mathcal{N}(i)} p_{ik} A(x_i - x_k)(x_i - x_k)^\top \right] \quad (10)$$

## 1.5 实验结果与分析

### 1.5.1 数据集的部分可视化分析

1. `train.csv` 文件中训练数据的 `pairplot` 图如图 1 所示。
2. 训练数据的分布情况如图 2 所示，可以看出这三类鸢尾花的数据分布比例是不完全一致的，但三类的数据量大致相同。

### 1.5.2 实验结果的分析

1. 对于基于欧氏距离的 KNN 算法，当  $k = 3$  时，测试集的准确率为 93.33%；
2. 对于基于马氏距离的 KNN 算法，当  $k = 3$  时，测试集的准确率为 93.33%；



图 1: train.csv 训练数据的 pairplot 图

Data Distribution



图 2: train.csv 训练数据分布比例

3. 将  $k$  从 1 到 30 遍历，分析出基于欧氏距离的 KNN 算法对 Iris 数据集进行分类的准确率随  $k$  的变化情况，如图 3 所示。

4. 进一步扩大  $k$  的取值范围，将  $k$  从 1 遍历到 100，分析出基于欧氏距离的 KNN 算法对 Iris 数据集进行分类的准确率随  $k$  的变化情况，如图 4 所示。

从上述实验结果可以分析出，基于欧式距离的最佳  $k$  值为 5 或 27 或 29，此时的准确率最高，为 100%，说明  $k = 5$ ， $k = 27$ ， $k = 29$  时能完全正确的将 Iris 数据集中三类鸢尾花进行分类。而图 3 和图 4 中的准确率曲线同时也说明了 KNN 算法中的  $k$  值选择对分类准确率的影响：

1. 当  $k$  值较小时，模型容易受到噪声的影响，导致过拟合；
2. 当  $k$  值较大时，模型容易受到样本不均衡的影响，导致欠拟合。

而  $k$  值过大就相当于是对所有样本进行投票，根据本实验模型（基于欧氏距离）所得的结果发现， $k \geq 93$  时，准确率降到最低的 26.67%。

同时，通过图 3 和图 4，我们还发现，准确率并不是  $k$  变化就随着变化的，而是几近分段变化的，这从侧面说明了本实验 Iris 数据集的离散属性，即不同类别的鸢尾花在特征空间中的分布是不均匀的，同时也说明了 Iris 数据集三种鸢尾花的部分属性具有聚集性，这也与图 1 中 pairplot 的小图中三种鸢尾花的属性相分离契合。

### 1.5.3 手动实现算法的评价

传统的 KNN 方法的不足之处主要包括：

1. 分类速度慢：最近邻分类器是基于实例学习的懒惰学习方法，因为它是根据所给训练样本构造的分类器，是将所有训练样本首先存储起来，当要进行分类时，就临时进行计算处理。需要计算待分样本与训练样本库中每一个样本的相似度，才能求得与其最近的  $K$  个样本。对于高维样本或样本集规模较大的情况，其时间和空间复杂度较高，时间代价为  $O(mn)$ ，其中  $m$  为向量空间模型空间特征维数， $n$  为训练样本集大小。

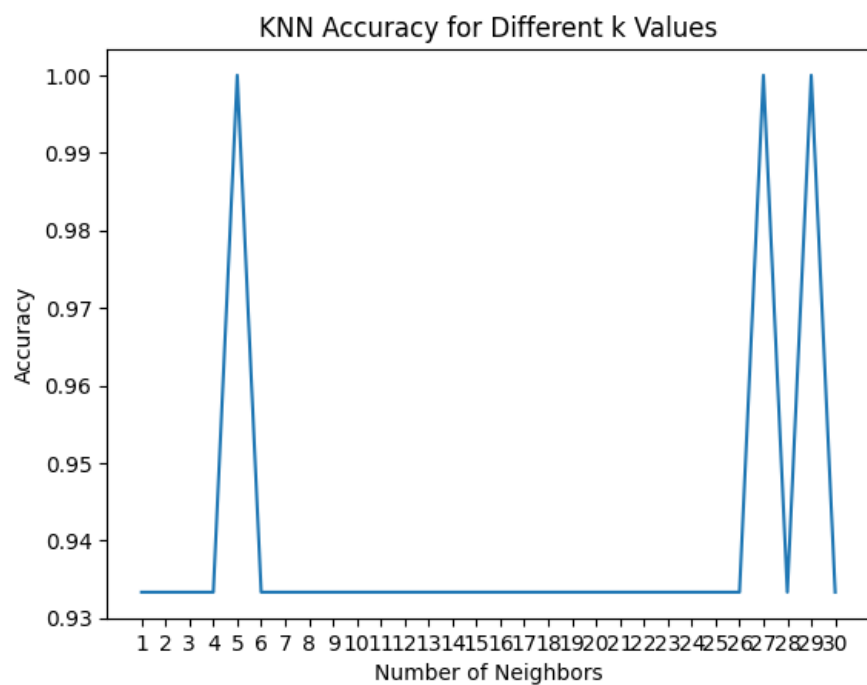


图 3: 分类准确率随  $k$  的变化情况

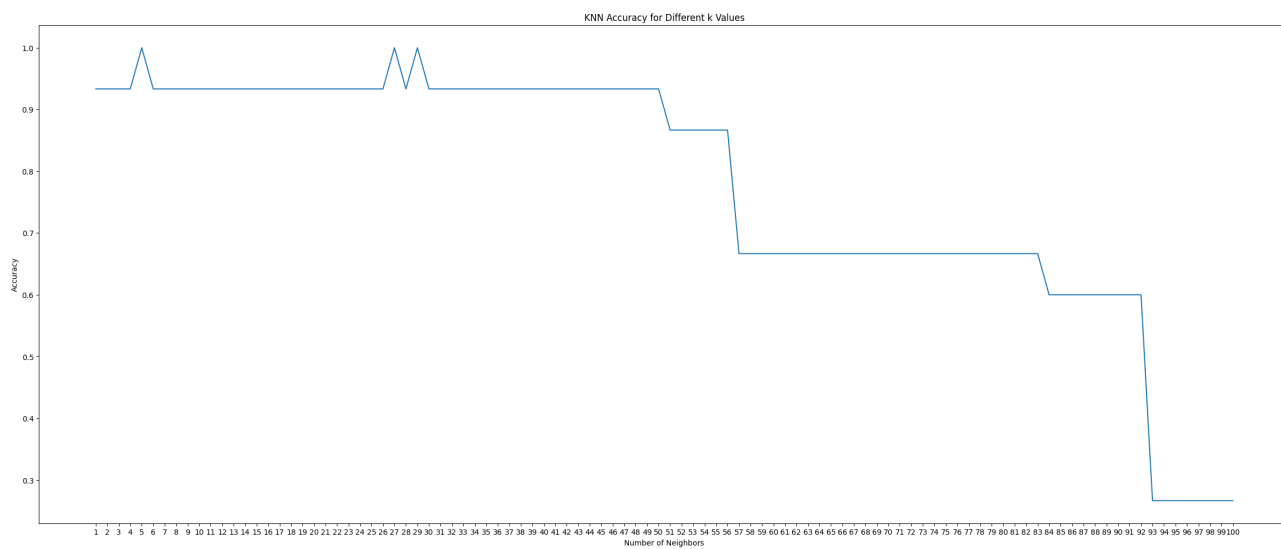


图 4: 分类准确率随  $k$  的变化情况

2. 样本库容量依赖性较强：对 KNN 算法在实际应用中的限制较大：有不少类别无法提供足够的训练样本，使得 KNN 算法所需要的相对均匀的特征空间条件无法得到满足，使得识别的误差较大。
3. 特征作用相同：与决策树归纳方法和神经网络方法相比，传统最近邻分类器认为每个属性的作用都是相同的（赋予相同权重）。样本的距离是根据样本的所有特征（属性）计算的。在这些特征中，有些特征与分类是强相关的，有些特征与分类是弱相关的，还有一些特征（可能是大部分）与分类不相关。这样，如果在计算相似度的时候，按所有特征作用相同来计算样本相似度就会误导分类过程。
4. K 值的确定：KNN 算法必须指定 K 值，K 值选择不当则分类精度不能保证。

KNN 的改进：对于 KNN 分类算法的改进方法主要可以分为加快分类速度、对训练样本库的维护、相似度的距离公式优化和 K 值确定四种类型。

#### 1. 加快 KNN 算法的分类速度

就学习而言，懒惰学习方法比积极学习方法要快，就计算量而言，它要比积极学习方法慢许多，因为懒惰学习方法在进行分类时，需要进行大量的计算。针对这一问题，到目前为止绝大多数解决方法都是基于减少样本量和加快搜索 K 个最近邻速度两个方面考虑的：

##### (a) 浓缩训练样本

当训练样本集中样本数量较大时，为了减小计算开销，可以对训练样本集进行编辑处理，即从原始训练样本集中选择最优的参考子集进行 K 最近邻寻找，从而减少训练样本的存储量和提高计算效率。这类方法主要包括 Condensing 算法、Wilson 的 Editing 算法和 Devijver 的 MultiEdit 算法，Kuncheva 使用遗传算法在这方面也进行了一些研究。

##### (b) 加快 K 个最近邻的搜索速度

这类方法是通过快速搜索算法，在较短时间内找到待分类样本的  $K$  个最近邻。在具体进行搜索时，不要使用盲目的搜索方法，而是要采用一定的方法加快搜索速度或减小搜索范围，例如可以构造交叉索引表，利用匹配成功与否的历史来修改样本库的结构，使用样本和概念来构造层次或网络来组织训练样本。

## 2. 相似度的距离公式的优化

为了改变传统 KNN 算法中特征作用相同的缺陷，可在相似度的距离公式中给特征赋予不同权重，例如在欧氏距离公式中给不同特征赋予不同权重。特征的权重一般根据各个特征在分类中的作用设定，可根据特征在整个训练样本库中的所起的作用大小来确定权重，也可根据在训练样本的局部样本（靠近待测试样本的样本集合）中的分类作用确定权重。

## 3. 对训练样本库的维护

对训练样本库进行维护以满足 KNN 算法的需要，包括对训练样本库中的样本进行添加或删除。对样本库的维护并不是简单的增加删除样本，而是可采用适当的办法来保证空间的大小，如符合某种条件的样本可以加入数据库中，同时可以对数据库库中已有符合某种条件的样本进行删除。从而保证训练样本库中的样本提供 KNN 算法所需要的相对均匀的特征空间。

## 4. $K$ 值选择 $K$ 的选择原则一般为：

- (a)  $K$  的选择往往通过大量独立的测试数据、多个模型来验证最佳的选择；
- (b)  $K$  值一般事先确定，也可以使用动态的，例如采用固定的距离指标，只对小于该指标的样本进行分析。

# 1.6 MindSpore 学习使用心得体会

本实验中使用了 mindspore 内置的 numpy，这个库与常规的 numpy 有一定区别，但是使用起来也是比较方便的。

在实验过程中，原本直接使用 `numpy` 在加载数据部分是将数据转换成 `numpy` 数组，而 `mindspore.numpy` 则是转换成 `mindspore.Tensor`，此外我还注意到，`mindspore` 的数据加载和处理部分与 `numpy` 有一定的区别，需要注意数据类型的转换，如 `asnumpy()`，在操作过程中注意数据转换成 `int32` 类型等。

## 1.7 代码附录（数据加载可视化展示部分，具体见 `knn.ipynb` 文件）

### 1.7.1 `knn.ipynb`

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import plotly.express as px
5  import seaborn as sns
6
7  iris_train = pd.read_csv("data/train.csv")
8
9  class KNN:
10     def __init__(self, n_neighbors=5):
11         self.n_neighbors = n_neighbors
12
13     def euclidean_distance(self, x1, x2):
14         return np.linalg.norm(x1 - x2)
15
16     def fit(self, X_train, y_train):
17         self.X_train = X_train
18         self.y_train = y_train
19
20     def predict(self, X):
21         # Create empty array to store the predictions
22         predictions = []
23         # Loop over X examples
24         for x in X:
25             # Get prediction using the prediction helper function
26             prediction = self._predict(x)
27             # Append the prediction to the predictions list
28             predictions.append(prediction)
29         return np.array(predictions)
30
```



```

31
32     def _predict(self, x):
33         # Create empty array to store distances
34         distances = []
35         # Loop over all training examples and compute the distance between x
and all the training examples
36         for x_train in self.X_train:
37             distance = self.euclidean_distance(x, x_train)
38             distances.append(distance)
39         distances = np.array(distances)
40
41         # Sort by ascendingly distance and return indices of the given n
neighbours
42         n_neighbors_idx = np.argsort(distances)[: self.n_neighbors]
43
44         # Get labels of n-neighbour indexes
45         labels = self.y_train[n_neighbors_idx]
46         labels = list(labels)
47         # Get the most frequent class in the array
48         most_occurring_value = max(labels, key=labels.count)
49         return most_occurring_value
50
51 class KNN2:
52     def __init__(self, n_neighbors=5, e=2, learning_rate=0.01, max_iter=100)
:
53         self.n_neighbors = n_neighbors
54         self.e = e
55         self.learning_rate = learning_rate
56         self.max_iter = max_iter
57         self.X_train = None
58         self.y_train = None
59         self.A = None
60
61     def fit(self, X_train, y_train):
62         self.X_train = X_train
63         self.y_train = y_train
64         n_samples, n_features = X_train.shape
65
66         # Initialize matrix A randomly
67         self.A = np.random.rand(self.e, n_features)
68

```

```

69     # Gradient descent to learn A
70     for iteration in range(self.max_iter):
71         gradient = np.zeros_like(self.A)
72         for i in range(n_samples):
73             xi = X_train[i]
74             yi = y_train[i]
75             same_class_indices = np.where(y_train == yi)[0]
76             for j in same_class_indices:
77                 if i == j:
78                     continue
79                 xj = X_train[j]
80                 pij = self._compute_pij(xi, xj, i, same_class_indices)
81                 diff = xi - xj
82                 gradient += 2 * pij * np.outer(self.A @ diff, diff)
83
84         self.A -= self.learning_rate * gradient / n_samples
85
86     def _compute_pij(self, xi, xj, i, same_class_indices):
87         numerator = np.exp(-self.mahalanobis_distance(xi, xj) ** 2)
88         denominator = 0
89         for k in same_class_indices:
90             if k == i:
91                 continue
92             xk = self.X_train[k]
93             denominator += np.exp(-self.mahalanobis_distance(xi, xk) ** 2)
94         return numerator / denominator
95
96     def mahalanobis_distance(self, x1, x2):
97         diff = x1 - x2
98         return np.sqrt(np.dot(np.dot(diff, self.A.T @ self.A), diff))
99
100    def predict(self, X):
101        # Create empty array to store the predictions
102        predictions = []
103        # Loop over X examples
104        for x in X:
105            # Get prediction using the prediction helper function
106            prediction = self._predict(x)
107            # Append the prediction to the predictions list
108            predictions.append(prediction)
109        return np.array(predictions)

```

```

110
111     def _predict(self, x):
112         # Create empty array to store distances
113         distances = []
114         # Loop over all training examples and compute the distance between x
115         # and all the training examples
116         for x_train in self.X_train:
117             distance = self.mahalanobis_distance(x, x_train)
118             distances.append(distance)
119         distances = np.array(distances)
120
121         # Sort by ascendingly distance and return indices of the given n
122         # neighbours
123         n_neighbors_idx = np.argsort(distances)[:self.n_neighbors]
124
125         # Get labels of n-neighbour indexes
126         labels = self.y_train[n_neighbors_idx]
127         labels = list(labels)
128         # Get the most frequent class in the array
129         most_occurring_value = max(labels, key=labels.count)
130         return most_occurring_value
131
132 # 读取训练集
133 train_data = pd.read_csv('data/train.csv')
134 X_train = train_data[['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal
135 Width']]
136 y_train = train_data['label']
137
138 # 读取验证集
139 test_data = pd.read_csv('data/val.csv')
140 X_test = test_data[['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal
141 Width']]
142 y_test = test_data['label']
143
144 # 检查数据类型并转换为浮点型
145 X_train = X_train.astype(float)
146 X_test = X_test.astype(float)
147
148 # 确保数据为NumPy数组
149 X_train = X_train.values
150 X_test = X_test.values

```

```

147
148 # 检查y_train和y_test是否为整数型
149 y_train = y_train.astype(int)
150 y_test = y_test.astype(int)
151
152 # 确保y_train和y_test为NumPy数组
153 y_train = y_train.values
154 y_test = y_test.values
155
156 model = KNN(3)
157 model.fit(X_train, y_train)
158 model2 = KNN2(3)
159 model2.fit(X_train, y_train)
160
161 def compute_accuracy(y_true, y_pred):
162     """
163     Computes the accuracy of a classification model.
164
165     Parameters:
166     y_true (numpy array): A numpy array of true labels for each data point.
167     y_pred (numpy array): A numpy array of predicted labels for each data
168     point.
169
170     Returns:
171     float: The accuracy of the model, expressed as a percentage.
172     """
173     y_true = y_true.flatten()
174     total_samples = len(y_true)
175     correct_predictions = np.sum(y_true == y_pred)
176     return (correct_predictions / total_samples)
177
178 predictions = model.predict(X_test)
179 accuracy = compute_accuracy(y_test, predictions)
180 print(f" our model got accuracy score of : {accuracy}")
181
182 # 初始化列表
183 a_index = list(range(1, 31))
184 accuracies = []
185
186 # 测试不同的邻居数

```

```

187 for k in a_index:
188     kcs = KNN(n_neighbors=k)
189     kcs.fit(X_train, y_train)
190     y_pred = kcs.predict(X_test)
191     accuracy = compute_accuracy(y_test, y_pred)
192     accuracies.append(accuracy)
193
194 # 转换为Pandas Series
195 a_series = pd.Series(accuracies, index=a_index)
196
197 # 绘制结果
198 plt.plot(a_index, a_series)
199 plt.xlabel('Number of Neighbors')
200 plt.ylabel('Accuracy')
201 plt.xticks(a_index)
202 plt.title('KNN Accuracy for Different k Values')
203 plt.show()
204
205 test_csv = pd.read_csv('data/test_data.csv')
206 test_pred = test_csv[['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal
    Width']]
207
208 # 检查数据类型并转换为浮点型
209 test_pred = test_pred.astype(float)
210
211 # 确保数据为NumPy数组
212
213 test_pred = test_pred.values
214
215 #对test_csv文件的每一项进行预测
216
217 predictions = model.predict(test_pred)
218 predictions = predictions.astype(int)
219 predictions = pd.DataFrame(predictions, columns=['label'])
220 predictions.to_csv('task1_test_prediction.csv', index=False)
221 predictions.head().style.background_gradient(sns.color_palette("YlOrBr",
    as_cmap=True))
222
223 test_csv2 = pd.read_csv('data/test_data.csv')
224 test_pred2 = test_csv[['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal
    Width']]

```

```

225
226 # 检查数据类型并转换为浮点型
227 test_pred2 = test_pred2.astype(float)
228
229 # 确保数据为NumPy数组
230
231 test_pred2 = test_pred2.values
232
233 #对test_csv文件的每一项进行预测
234
235 predictions2 = model2.predict(test_pred2)
236 predictions2 = predictions2.astype(int)
237 predictions2 = pd.DataFrame(predictions2, columns=['label'])
238 predictions2.to_csv('task2_test_prediction.csv', index=False)
239 iris_train.head().style.background_gradient(sns.color_palette("YlOrBr",
    as_cmap=True))
240
241 # 读取两个CSV文件
242 file1 = pd.read_csv('task1_test_prediction.csv')
243 file2 = pd.read_csv('task2_test_prediction.csv')
244
245 # 比较两个DataFrame并找出不同
246 comparison = file1 == file2
247
248 # 找出不同的行和列
249 differences = comparison[comparison == False]
250
251 # 打印出不同的值和位置
252 print("Differences found at these locations:")
253 print(differences)
254
255 # 显示不同的行
256 diff_rows = differences.dropna(how='all')
257 print("Rows with differences:")
258 print(diff_rows)
259
260 # 显示不同的列
261 diff_cols = differences.dropna(axis=1, how='all')
262 print("Columns with differences:")
263 print(diff_cols)

```

## 1.7.2 knn\_mindspore.ipynb

```
1  import mindspore.context as context
2  import mindspore.numpy as mnp
3  import mindspore.ops as ops
4  import pandas as pd
5  from mindspore import Tensor
6
7  # 设置MindSpore的运行环境
8  context.set_context(mode=context.GRAPH_MODE, device_target="CPU")
9
10 # 自定义数据加载函数
11 def load_data(file_path, has_label=True):
12     data = pd.read_csv(file_path)
13     if has_label:
14         X = data.iloc[:, :-1].values
15         y = data.iloc[:, -1].values
16         return X, y
17     else:
18         X = data.values
19         return X
20
21 # 加载数据
22 X_train, y_train = load_data('data/train.csv')
23 X_val, y_val = load_data('data/val.csv')
24 X_test = load_data('data/test_data.csv', has_label=False)
25
26 # 将数据转换为MindSpore张量
27 X_train = Tensor(X_train, mnp.float32)
28 y_train = Tensor(y_train, mnp.int32)
29 X_val = Tensor(X_val, mnp.float32)
30 y_val = Tensor(y_val, mnp.int32)
31 X_test = Tensor(X_test, mnp.float32)
32
33 # 定义KNN算法
34 class KNN:
35     def __init__(self, k=3):
36         self.k = k
37
38     def fit(self, X, y):
39         self.X_train = X
40         self.y_train = y
```

```

41
42     def predict(self, X):
43         distances = self.compute_distances(X)
44         return self.predict_labels(distances)
45
46     def compute_distances(self, X):
47         num_test = X.shape[0]
48         num_train = self.X_train.shape[0]
49         dists = mnp.zeros((num_test, num_train))
50         for i in range(num_test):
51             dists[i, :] = mnp.sqrt(mnp.sum((self.X_train - X[i, :])**2, axis
=1))
52         return dists
53
54     def predict_labels(self, dists):
55         num_test = dists.shape[0]
56         y_pred = mnp.zeros(num_test, dtype=mnp.int32)
57         for i in range(num_test):
58             closest_y = []
59             sorted_indices = ops.Sort(axis=0)(dists[i, :])[1]
60             closest_y = self.y_train[sorted_indices[:self.k]]
61             y_pred[i] = mnp.bincount(closest_y).argmax()
62         return y_pred
63
64 # 创建KNN实例并训练
65 knn = KNN(k=3)
66 knn.fit(X_train, y_train)
67
68 # 在验证集上进行预测
69 y_val_pred = knn.predict(X_val)
70 accuracy_val = mnp.mean((y_val_pred == y_val).astype(mnp.float32))
71 print("Validation Accuracy:", accuracy_val.asnumpy())
72
73 # 在测试集上进行预测
74 y_test_pred = knn.predict(X_test)
75 print("Test Predictions:", y_test_pred.asnumpy())
76
77 # 将预测结果保存到CSV文件
78 df_predictions = pd.DataFrame(y_test_pred.asnumpy(), columns=['label'])
79 df_predictions.to_csv('task3_test_prediction.csv', index=False)

```



## 2 实验二 Naïve Bayes Classification

### 2.1 问题描述

#### 2.1.1 概述

利用朴素贝叶斯算法，对 MNIST 数据集中的测试集进行分类。

#### 2.1.2 任务说明

1. 在课程学习中同学们已经学习了贝叶斯分类理论并掌握了其基本原理，即利用贝叶斯公式

$$p(\omega_j|x) = \frac{p(x|\omega_j)p(\omega_j)}{p(x)}$$

对  $p(\omega_j|x)$  作出预测。由于  $p(x)$  为一固定值，所以一般不在计算过程中求得  $p(x)$  的具体值。在实际运用中，为了方便计算，通常假设数据特征之间相互独立，即

$$p(x|\omega_j) = p(x_1|\omega_j) \cdot p(x_2|\omega_j) \cdots p(x_d|\omega_j), \quad x \in \mathbb{R}^d,$$

这便是著名的朴素贝叶斯算法。

2. MNIST 数据集本身以二进制形式保存，所以首先需要选择合适的编程语言编写读写二进制数据的程序完成对图片、标记信息的初步提取工作。读取了图片信息后，发现每个像素点的值在  $[0,1]$  区间内，这是图像压缩后的结果，所以可以先将像素值乘以 255 再取整，得到每一个点的灰度值。将图像二值化，得到可以用于分类的  $28 \times 28$  个特征向量以及对应的标签数据，之后便可以交由贝叶斯分类器进行学习。
3. 基于 MindSpore 平台提供的官方模型库，对相同的数据集进行训练，并与自己独立实现的算法对比结果（包括但不限于准确率、算法迭代收敛次数等指标），并分析结果中出现差异的可能原因，给出使用 MindSpore 的心得和建议。

4. (加分项) 使用 MindSpore 平台提供的相似任务数据集 (例如, 其他的分类任务数据集) 测试自己独立实现的算法并与 MindSpore 平台上的官方实现算法进行对比, 并进一步分析差异及其成因。

## 2.2 实现步骤与流程

### 2.2.1 实验思路

1. 读取数据集的图片和标签信息;
2. 对图片信息进行预处理, 包括归一化、二值化和将图像展开成一维向量;
3. 实现朴素贝叶斯算法, 包括拟合和预测两个步骤;
4. 使用预训练模型对测试数据进行预测, 计算准确率;
5. 可视化部分模型对测试数据的预测结果。

### 2.2.2 数学模型

朴素贝叶斯学习步骤如下。先计算类先验概率分布:

$$P(Y = c_k) = \frac{1}{N} \sum_{i=1}^N I(\hat{y}_i = c_k), \quad k = 1, 2, \dots, K \quad (11)$$

其中  $c_k$  表示第  $k$  个类别,  $y_i$  表示第  $i$  个样本的类标记。类先验概率分布可以通过极大似然估计得到。

然后计算类条件概率分布:

$$P(X = x|Y = c_k) = P(X^{(1)} = x^{(1)}, \dots, X^{(n)} = x^{(n)}|Y = c_k), \quad k = 1, 2, \dots, K \quad (12)$$

直接对  $P(X = x|Y = c_k)$  进行估计不太可行, 因为参数量太大。但是朴素贝叶斯的一

个最重要的假设就是条件独立性假设，即：

$$P(X = x|Y = c_k) = P(X^{(1)} = x^{(1)}, \dots, X^{(n)} = x^{(n)}|Y = c_k) = \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k) \quad (13)$$

有了条件独立性假设之后，便可基于极大似然估计计算类条件概率。

类先验概率分布和类条件概率分布都计算得到之后，基于贝叶斯公式即可以计算类后验概率：

$$P(Y = c_k|X = x) = \frac{P(X = x|Y = c_k)P(Y = c_k)}{\sum_k P(X = x|Y = c_k)P(Y = c_k)} \quad (14)$$

代入类条件计算公式，有：

$$P(Y = c_k|X = x) = \frac{\prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)P(Y = c_k)}{\sum_k \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)P(Y = c_k)} \quad (15)$$

基于上述公式即可以学习一个朴素贝叶斯模型。给定新的数据样本时，计算其最大后验概率即可：

$$\hat{y} = \arg \max_{c_k} \frac{\prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)P(Y = c_k)}{\sum_k \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)P(Y = c_k)} \quad (16)$$

其中，分母对于所有的  $\hat{y}$  都是一样的，所以上述式可进一步简化为：

$$\hat{y} = \arg \max_{c_k} \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)P(Y = c_k) \quad (17)$$

方程 17 即为朴素贝叶斯算法的预测公式。

### 2.2.3 关键难点

朴素贝叶斯算法的难点在于如何高效地计算类条件概率分布。由于朴素贝叶斯算法的条件独立性假设，可以将类条件概率分布分解为各个特征的条件概率分布的乘积。这样可以大大减少计算量。

## 2.2.4 算法描述

朴素贝叶斯算法实现的伪代码如算法 3 所示。

---

**Algorithm 3** 朴素贝叶斯分类器

---

```
1: 初始化:
2: 定义类别列表 classes
3: 定义类别先验概率 class_priors
4: 定义特征概率 feature_probs
5: procedure FIT(X, y)
6:   classes = np.unique(y)
7:   for 每个类别 cls in classes do
8:      $X_{cls} = X[y == cls]$ 
9:      $P(cls) = \frac{|X_{cls}|}{|X|}$ 
10:     $P(x_i|cls) = \frac{\sum X_{cls,i} + 1}{|X_{cls}| + 2}$ 
11:   end for
12: end procedure
13: procedure PREDICT(X)
14:   for 每个样本 x in X do
15:     for 每个类别 cls in classes do
16:        $\log P(cls|x) = \log P(cls) + \sum (\log P(x_i|cls) \times x_i + \log(1 - P(x_i|cls)) \times (1 - x_i))$ 
17:     end for
18:     选择最大后验概率对应的类别
19:   end for
20: end procedure
```

---

## 2.3 实验结果与分析

### 2.3.1 手动实现朴素贝叶斯算法

利用课程所提供的 MNIST 数据集，我们手动实现了朴素贝叶斯算法，并在测试集上进行了分类。可以看出实验的准确率为 **84.27%**。

为便于观察手动实现朴素贝叶斯算法模型的性能，我们在部分验证集上进行了测试，测试结果如图 5 所示。可以看出在多数图片的预测中，手动实现的朴素贝叶斯算法还是能够将 MNIST 数据集中的数字进行正确分类的。

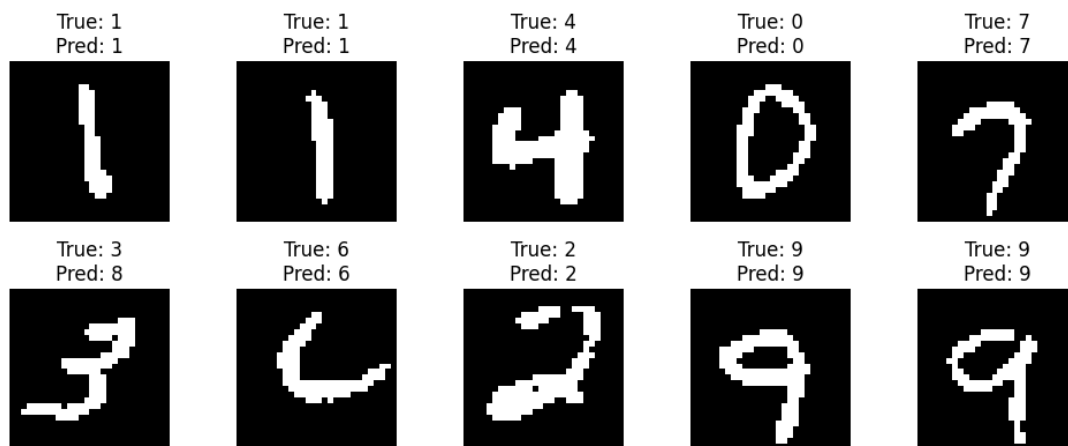


图 5: 手动实现 Naïve Bayes 测试结果

### 2.3.2 mindSpore 实现朴素贝叶斯算法

我们使用了 `mindspore.numpy` 来代替原有的 `numpy` 实现了朴素贝叶斯算法。

使用 `mindspore` 实现的朴素贝叶斯算法在 MNIST 数据集上的分类准确率为 **84.27%**。准确率和手动实现的朴素贝叶斯算法相同，说明 `mindspore` 实现的朴素贝叶斯算法是正确的。

但是我发现，`mindspore.numpy` 实现要比 `numpy` 花费的时间要长，这可能是因为：

1. MindSpore 的计算图模式（`GRAPH_MODE`）可以优化性能，但需要一些额外的编译时间。相比之下，使用 `NumPy` 的代码是动态执行的，没有编译阶段。

2. MindSpore 中的 Tensor 操作在某些情况下可能比 NumPy 更慢，特别是在小规模数据和简单计算时。因为 MindSpore 的设计初衷是用于大规模的深度学习任务，优化了大规模数据的处理，而不是小规模简单任务。
3. 代码中有很多从 NumPy 到 mindspore.Tensor 的类型转换，这些操作本身也会带来一些开销。类型转换会导致额外的计算时间和内存开销。

### 2.3.3 手动实现算法的评价

实验结果表明,手动实现的朴素贝叶斯算法在 MNIST 数据集上的分类准确率为 **84.27%**。这个准确率表明多数的 MNIST 数据集中的数字都能够被正确分类。该结果初步表明手动实现朴素贝叶斯算法的性能是可以接受的。

未来可以采用半朴素贝叶斯算法，适当考虑一部分属性间的相互依赖信息，从而既不需要进行完全联合概率计算，又不至于彻底忽略了比较强的相互依赖关系。因为朴素贝叶斯算法的假设是：假设属性之间相互独立，但实际上属性之间是有关联的，所以半朴素贝叶斯算法是对朴素贝叶斯算法的一种改进。

## 2.4 MindSpore 学习使用心得体会

在这个实验中，我们使用了 mindspore.numpy 来代替原有的 numpy 实现了朴素贝叶斯算法。我感觉 mindspore 的使用还是不如直接用 numpy 方便，因为 mindspore.numpy 中没有 frombuffer 这个属性，所以在读取 MNIST 数据集时，我是使用 numpy 来读取数据，然后再转换成 mindspore.Tensor。但是这个示例为我深入理解 mindspore 的操作提供了一个很好的机会。

## 2.5 代码附录

### 2.5.1 naive\_bayes.ipynb

```
1 import numpy as np
```

```

2  import struct
3  import gzip
4  import matplotlib.pyplot as plt
5
6  def read_images(file_path):
7      with gzip.open(file_path, 'rb') as f:
8          # 读取文件头信息：魔数和图片数量
9          magic, num_images = struct.unpack(">II", f.read(8))
10         # 读取图片的行数和列数
11         num_rows, num_cols = struct.unpack(">II", f.read(8))
12         print(f"Magic number: {magic}, Number of images: {num_images}, Rows:
{num_rows}, Columns: {num_cols}")
13         # 读取图片数据
14         images = np.frombuffer(f.read(), dtype=np.uint8).reshape(num_images,
num_rows, num_cols)
15         return images
16
17  def read_labels(file_path):
18      with gzip.open(file_path, 'rb') as f:
19          # 读取文件头信息：魔数和标签数量
20          magic, num_labels = struct.unpack(">II", f.read(8))
21          print(f"Magic number: {magic}, Number of labels: {num_labels}")
22          # 读取标签数据
23          labels = np.frombuffer(f.read(), dtype=np.uint8)
24          return labels
25
26  # 设置数据集文件路径
27  train_images_path = 'data/train-images-idx3-ubyte.gz'
28  train_labels_path = 'data/train-labels-idx1-ubyte.gz'
29  test_images_path = 'data/t10k-images-idx3-ubyte.gz'
30  test_labels_path = 'data/t10k-labels-idx1-ubyte.gz'
31
32  # 读取数据集
33  X_train = read_images(train_images_path)
34  y_train = read_labels(train_labels_path)
35  X_test = read_images(test_images_path)
36  y_test = read_labels(test_labels_path)
37
38  print(f"Training data shape: {X_train.shape}")
39  print(f"Training labels shape: {y_train.shape}")
40  print(f"Test data shape: {X_test.shape}")

```

```

41 print(f"Test labels shape: {y_test.shape}")
42
43 # 将像素值归一化到[0, 1]范围
44 X_train = X_train / 255.0
45 X_test = X_test / 255.0
46
47 # 将图像二值化
48 X_train = (X_train > 0.5).astype(int)
49 X_test = (X_test > 0.5).astype(int)
50
51 # 将图像展开成一维向量
52 X_train = X_train.reshape(X_train.shape[0], -1)
53 X_test = X_test.reshape(X_test.shape[0], -1)
54
55 class NaiveBayes:
56     def __init__(self):
57         self.classes = None
58         self.class_priors = {}
59         self.feature_probs = {}
60
61     def fit(self, X, y):
62         self.classes = np.unique(y)
63         n_samples, n_features = X.shape
64
65         for cls in self.classes:
66             X_cls = X[y == cls]
67             self.class_priors[cls] = X_cls.shape[0] / n_samples
68             self.feature_probs[cls] = (X_cls.sum(axis=0) + 1) / (X_cls.shape
[0] + 2)
69
70     def predict(self, X):
71         posteriors = []
72         for x in X:
73             posterior_probs = {}
74             for cls in self.classes:
75                 prior = np.log(self.class_priors[cls])
76                 conditional = np.sum(np.log(self.feature_probs[cls]) * x +
np.log(1 - self.feature_probs[cls]) * (1 - x))
77                 posterior_probs[cls] = prior + conditional
78                 posteriors.append(max(posterior_probs, key=posterior_probs.get))
79         return np.array(posteriors)

```



```

80
81 # 实例化朴素贝叶斯分类器
82 nb = NaiveBayes()
83
84 # 对数据进行训练
85 nb.fit(X_train, y_train)
86
87 # 对测试数据进行预测
88 y_pred = nb.predict(X_test)
89 accuracy = np.mean(y_pred == y_test)
90 print(f"Accuracy: {accuracy:.4f}")
91
92 # 可视化部分预测结果
93 import random
94
95 # 随机选择10个测试样本
96 indices = random.sample(range(X_test.shape[0]), 10)
97 images = X_test[indices].reshape(-1, 28, 28)
98 true_labels = y_test[indices]
99 pred_labels = y_pred[indices]
100
101 plt.figure(figsize=(10, 10))
102 for i in range(10):
103     plt.subplot(5, 5, i + 1)
104     plt.imshow(images[i], cmap='gray')
105     plt.title(f"True: {true_labels[i]}\nPred: {pred_labels[i]}")
106     plt.axis('off')
107 plt.tight_layout()
108 plt.show()

```

## 2.5.2 naive\_bayes\_mindspore.ipynb

```

1 import struct
2 import gzip
3 import matplotlib.pyplot as plt
4 import numpy as np # 用于从缓冲区读取数据
5 import mindspore.context as context
6 import mindspore.numpy as mnp
7 import mindspore.ops as ops
8 from mindspore import Tensor
9

```

```

10 # 设置MindSpore的运行环境
11 context.set_context(mode=context.GRAPH_MODE, device_target="CPU")
12
13 def read_images(file_path):
14     with gzip.open(file_path, 'rb') as f:
15         # 读取文件头信息：魔数和图片数量
16         magic, num_images = struct.unpack(">II", f.read(8))
17         # 读取图片的行数和列数
18         num_rows, num_cols = struct.unpack(">II", f.read(8))
19         print(f"Magic number: {magic}, Number of images: {num_images}, Rows:
20 {num_rows}, Columns: {num_cols}")
21         # 读取图片数据
22         buffer = f.read()
23         images = np.frombuffer(buffer, dtype=np.uint8).reshape(num_images,
24 num_rows, num_cols)
25         return Tensor(images, dtype=mnps.float32)
26
27 def read_labels(file_path):
28     with gzip.open(file_path, 'rb') as f:
29         # 读取文件头信息：魔数和标签数量
30         magic, num_labels = struct.unpack(">II", f.read(8))
31         print(f"Magic number: {magic}, Number of labels: {num_labels}")
32         # 读取标签数据
33         buffer = f.read()
34         labels = np.frombuffer(buffer, dtype=np.uint8)
35         return Tensor(labels, dtype=mnps.int32)
36
37 # 设置数据集文件路径
38 train_images_path = 'data/train-images-idx3-ubyte.gz'
39 train_labels_path = 'data/train-labels-idx1-ubyte.gz'
40 test_images_path = 'data/t10k-images-idx3-ubyte.gz'
41 test_labels_path = 'data/t10k-labels-idx1-ubyte.gz'
42
43 # 读取数据集
44 X_train = read_images(train_images_path)
45 y_train = read_labels(train_labels_path)
46 X_test = read_images(test_images_path)
47 y_test = read_labels(test_labels_path)
48
49 print(f"Training data shape: {X_train.shape}")
50 print(f"Training labels shape: {y_train.shape}")

```

```

49 print(f"Test data shape: {X_test.shape}")
50 print(f"Test labels shape: {y_test.shape}")
51
52 # 将数据转换为mindspore.numpy数组
53 X_train = mnp.array(X_train) / 255.0
54 X_test = mnp.array(X_test) / 255.0
55 y_train = mnp.array(y_train)
56 y_test = mnp.array(y_test)
57
58 # 将图像二值化
59 X_train = (X_train > 0.5).astype(mnp.float32)
60 X_test = (X_test > 0.5).astype(mnp.float32)
61
62 # 将图像展开成一维向量
63 X_train = X_train.reshape(X_train.shape[0], -1)
64 X_test = X_test.reshape(X_test.shape[0], -1)
65
66 class NaiveBayes:
67     def __init__(self):
68         self.classes = None
69         self.class_priors = {}
70         self.feature_probs = {}
71
72     def fit(self, X, y):
73         self.classes = mnp.unique(y)
74         n_samples, n_features = X.shape
75
76         for cls in self.classes:
77             indices = mnp.where(y == cls, mnp.ones_like(y), mnp.zeros_like(y)
78             )).astype(bool)
79             X_cls = X[indices]
80             self.class_priors[int(cls.asnumpy())] = X_cls.shape[0] /
n_samples
81             self.feature_probs[int(cls.asnumpy())] = (X_cls.sum(axis=0) + 1)
82             / (X_cls.shape[0] + 2)
83
84     def predict(self, X):
85         posteriors = []
86         for x in X:
87             posterior_probs = {}
88             for cls in self.classes:

```

```

87         cls_int = int(cls.asnumpy())
88         prior = mnp.log(Tensor(self.class_priors[cls_int], dtype=mnp
        .float32))
89         conditional = mnp.sum(mnp.log(Tensor(self.feature_probs[
        cls_int], dtype=mnp.float32)) * x + mnp.log(1 - Tensor(self.feature_probs[
        cls_int], dtype=mnp.float32)) * (1 - x))
90         posterior_probs[cls_int] = prior + conditional
91         posteriors.append(max(posterior_probs, key=posterior_probs.get))
92         return mnp.array(posteriors)
93
94 # 实例化朴素贝叶斯分类器
95 nb = NaiveBayes()
96
97 # 对数据进行训练
98 nb.fit(X_train, y_train)
99
100 # 对测试数据进行预测
101 y_pred = nb.predict(X_test)
102 accuracy = mnp.mean((y_pred == y_test).astype(mnp.float32))
103 print(f"Accuracy: {accuracy.asnumpy():.4f}")
104
105 # 可视化部分预测结果
106 import random
107
108 # 随机选择10个测试样本
109 indices = random.sample(range(X_test.shape[0]), 10)
110 images = X_test[indices].reshape(-1, 28, 28)
111 true_labels = y_test[indices]
112 pred_labels = y_pred[indices]
113
114 plt.figure(figsize=(10, 10))
115 for i in range(10):
116     plt.subplot(5, 5, i + 1)
117     plt.imshow(images[i].asnumpy(), cmap='gray')
118     plt.title(f"True: {true_labels[i].asnumpy()}\nPred: {pred_labels[i].
        asnumpy()}")
119     plt.axis('off')
120 plt.tight_layout()
121 plt.show()

```

## 3 实验三 Neural Network Image Classification

### 3.1 问题描述

### 3.2 概述

利用神经网络算法，对 CIFAR 数据集中的测试集进行分类。

### 3.3 任务说明

1. 基于神经网络模型及 BP 算法，根据训练集中的数据对你设计的神经网络模型进行训练，随后对给定的打乱的测试集中的数据进行分类。
2. 基于 MindSpore 平台提供的官方模型库，对相同的数据集进行训练，并与自己独立实现的算法对比结果（包括但不限于准确率、算法迭代收敛次数等指标），并分析结果中出现差异的可能原因。
- 3.（加分项）使用 MindSpore 平台提供的相似任务数据集（例如，其他的分类任务数据集）测试自己独立实现的算法并与 MindSpore 平台上的官方实现算法进行对比，并进一步分析差异及其成因。

### 3.4 实现步骤与流程

#### 3.4.1 实验环境

实验环境见表 1。

#### 3.4.2 实验思路

1. 数据加载与预处理:
  - (a) 从 CIFAR-10 数据集中加载图像数据和标签。

表 1: Experiment Environment

Items	Version
CPU	Intel Core i5-1135G7
RAM	16 GB
Python	3.11.5
Operating system	Windows11

(b) 将图像数据从原始格式转换为适用于神经网络的格式 ( $32 \times 32 \times 3 \rightarrow 3072 \times 1$ )。

(c) 将图像数据进行归一化处理，将像素值缩放到  $[0, 1]$  范围内。

## 2. 数据划分:

(a) 将加载的训练数据划分为训练集和验证集，以便在训练过程中进行模型评估。

## 3. 定义全连接神经网络结构:

(a) 输入层：包含 3072 个神经元，对应每张图像的 3072 个像素值。

(b) 隐藏层：包含 100 个神经元，使用 ReLU 激活函数。

(c) 输出层：包含 10 个神经元，对应 10 个类别，使用 Softmax 激活函数。

## 4. 定义前向传播与反向传播:

(a) 前向传播：计算输入数据通过网络后的输出。

(b) 反向传播：根据预测结果和实际标签计算损失，并更新网络权重。

## 5. 模型训练:

(a) 使用小批量梯度下降优化网络权重。

(b) 在每个 epoch 结束后，计算并记录训练损失、验证损失和验证准确率。

#### 6. 模型评估:

(a) 在测试集上评估模型性能，计算并输出测试集准确率。

(b) 绘制训练过程中损失和准确率的变化曲线。

#### 7. 混淆矩阵:

(a) 计算混淆矩阵，分析分类结果的具体表现。

(b) 绘制混淆矩阵，直观展示模型在各个类别上的分类效果。

#### 8. 预测结果展示

(a) 展示模型在测试集部分图片上的分类结果。

本实验的评估指标和超参数如表 2 所示。

表 2: 实验评估指标和超参数

参数	值
learning rate	0.01
epochs	100
loss function	Cross Entropy
performance	accuracy
batch size	64
num hiddens	128

### 3.4.3 数学模型

BP 神经网络的数学模型如图 6 所示。

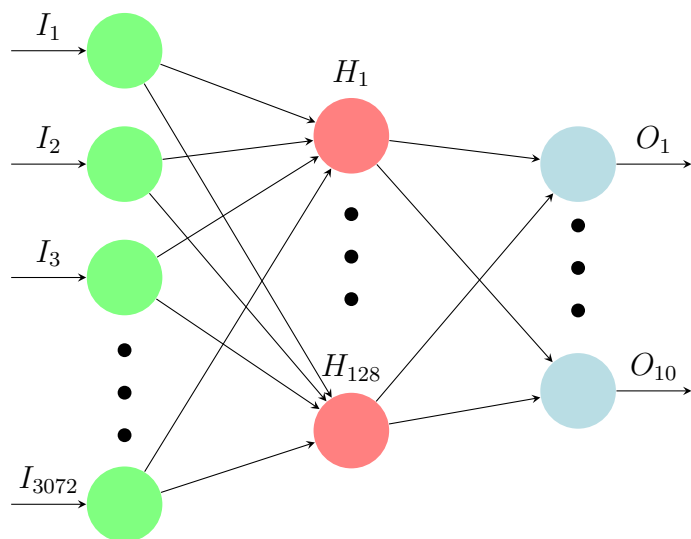


图 6: 全连接神经网络结构示意图

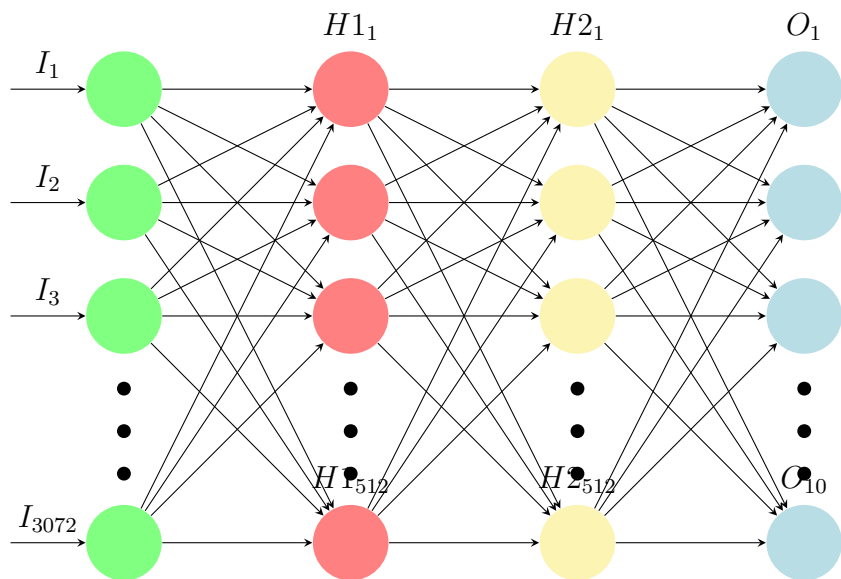


图 7: 基于 MindSpore 平台实现的多隐藏层神经网络结构示意图



#### 3.4.4 关键难点

在训练过程中实时评估模型的性能，并根据评估结果调整模型的超参数，如学习率、批次大小、隐藏层神经元数量等。需要平衡模型的复杂度和泛化能力，避免过拟合或欠拟合。

#### 3.4.5 算法描述

使用 BP 神经网络训练算法，如算法 4 所示。

### 3.5 实验结果与分析

#### 3.5.1 实验结果展示

1. 准确率曲线本实验训练集和验证集上的准确率随 epoch 的变化曲线如图 8 所示。从图中可以看出，随着 epoch 的增加，验证集的准确率逐渐提高，最终收敛到 **50%**。
2. 损失曲线本实验验证集上的交叉熵损失随 epoch 的变化曲线如图 9 所示。从图中可以看出，随着 epoch 的增加，训练集和验证集的交叉熵损失逐渐降低，最终收敛到一个稳定值。其中训练集交叉熵损失稳定在 **1.0** 左右，验证集交叉熵损失稳定在 **1.4** 左右。
3. 混淆矩阵本实验的混淆矩阵如图 10 所示。从图中可以看出，模型在 CIFAR-10 数据集上的分类效果较好，大部分类别的分类准确率较高。
4. 预测结果本实验的预测结果如图 11 所示。从图中可以看出，模型在测试集上的分类效果较好，大部分图片的分类结果正确。

#### 3.5.2 进一步探究

虽然上述基于全连接神经网络的分类器在 CIFAR-10 数据集上取得了一定的分类效果，但是其准确率仍然较低，仅为 50% 左右。结合机器学习课程所学知识，可以尝试以下方法进一步提高分类器的性能：

---

**Algorithm 4** 全连接神经网络训练算法

---

```
1: 初始化网络的权重  $\mathbf{W1}$ ,  $\mathbf{W2}$  和偏置  $\mathbf{b1}$ ,  $\mathbf{b2}$  为随机值
2: procedure 训练 ( $X_{\text{train}}$ ,  $Y_{\text{train}}$ ,  $X_{\text{val}}$ ,  $Y_{\text{val}}$ , batch_size, learning_rate, epochs)
3:   for 每个 epoch do
4:     for 每个 mini-batch do
5:       前向传播:
6:         设输入  $\mathbf{X}$ 
7:         计算隐藏层的输出和激活值:
8:            $\mathbf{Z1} = \mathbf{XW1} + \mathbf{b1}$ 
9:            $\mathbf{A1} = \max(0, \mathbf{Z1})$  ▷ ReLU 激活函数
10:        计算输出层的输出和激活值:
11:           $\mathbf{Z2} = \mathbf{A1W2} + \mathbf{b2}$ 
12:           $\mathbf{A2} = \text{softmax}(\mathbf{Z2})$ 
13:        计算损失:
14:          使用交叉熵损失函数计算损失  $L$ 
15:        反向传播:
16:          计算输出层的误差  $\delta^{(2)}$ 
17:          计算隐藏层的误差  $\delta^{(1)}$ 
18:          计算梯度:
19:             $\nabla_{\mathbf{W2}} = \mathbf{A1}^T \delta^{(2)}$ 
20:             $\nabla_{\mathbf{b2}} = \sum \delta^{(2)}$ 
21:             $\nabla_{\mathbf{W1}} = \mathbf{X}^T \delta^{(1)}$ 
22:             $\nabla_{\mathbf{b1}} = \sum \delta^{(1)}$ 
23:          更新权重和偏置:
24:            使用学习率  $\eta$  更新  $\mathbf{W1}$ ,  $\mathbf{b1}$ ,  $\mathbf{W2}$ ,  $\mathbf{b2}$ 
25:        end for
26:      计算训练集和验证集的损失和准确率
27:    end for
28: end procedure
```

---

- 尝试使用更复杂的神经网络结构，如 CNN、RNN、Transformers 等
- 尝试使用更高级的优化算法，如 Adam、RMSprop 等
- 尝试使用更复杂的数据增强技术，如旋转、平移、缩放等
- 尝试使用更复杂的模型评估指标，如 F1-score、ROC 曲线等
- 尝试使用更复杂的模型融合技术，如集成学习、模型融合等
- 尝试使用更复杂的超参数调优技术，如网格搜索、贝叶斯优化等
- 尝试使用更复杂的模型解释技术，如 LIME、SHAP 等

### 3.5.3 手动实现算法的评价

本实验实现了单隐藏层前馈神经网络对 CIFAR-10 数据集进行分类，取得了 50% 左右的准确率。与 mindspore 平台相比，手动实现的算法在准确率和收敛速度上均有所不足，可能的原因如下：

1. 模型复杂度：手动实现的算法只使用了单隐藏层的前馈神经网络，模型复杂度较低，难以捕捉数据集的复杂特征。
2. 优化算法：手动实现的算法使用了简单的小批量梯度下降优化算法，收敛速度较慢，难以达到较高的准确率。而 mindspore 中采用的是 SGD、Adam 等高级的优化算法。
3. 超参数调优：手动实现的算法中的超参数（学习率、批次大小、隐藏层神经元数量等）未经过充分调优，可能导致模型性能不佳。

### 3.5.4 手动实现 BP 神经网络算法和基于 Mindspore 实现神经网络算法的对比

1. 准确率对比：基于 Mindspore 实现的神经网络在 CIFAR-10 数据集上的准确率约为 53%，如图 12 所示，高于手动实现的算法的准确率 50%。

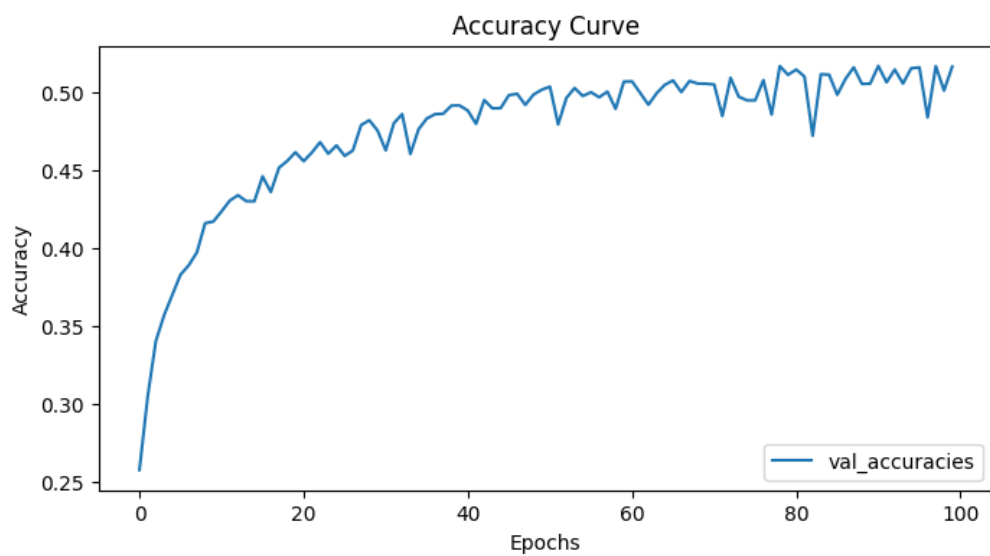


图 8: Accuracy 随 epoch 的变化曲线

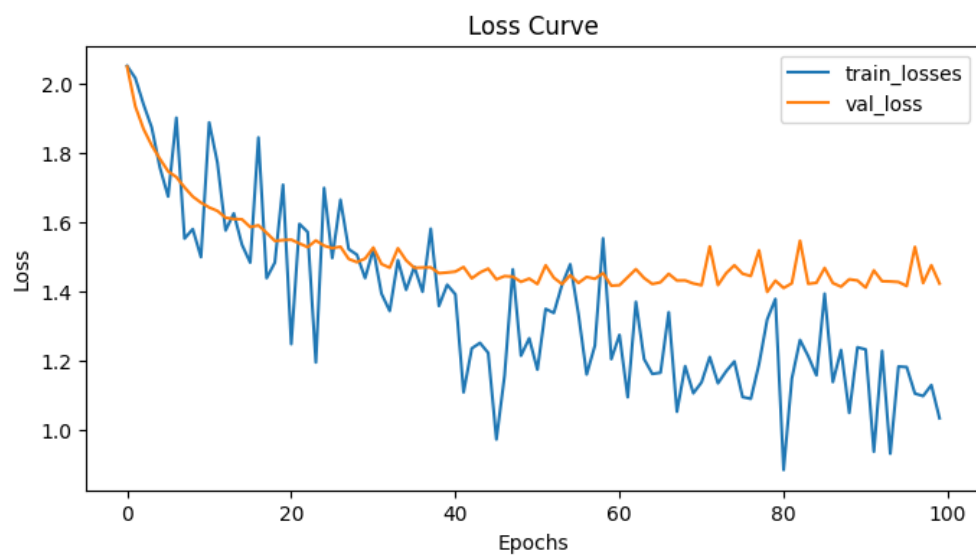


图 9: Loss 随 epoch 的变化曲线

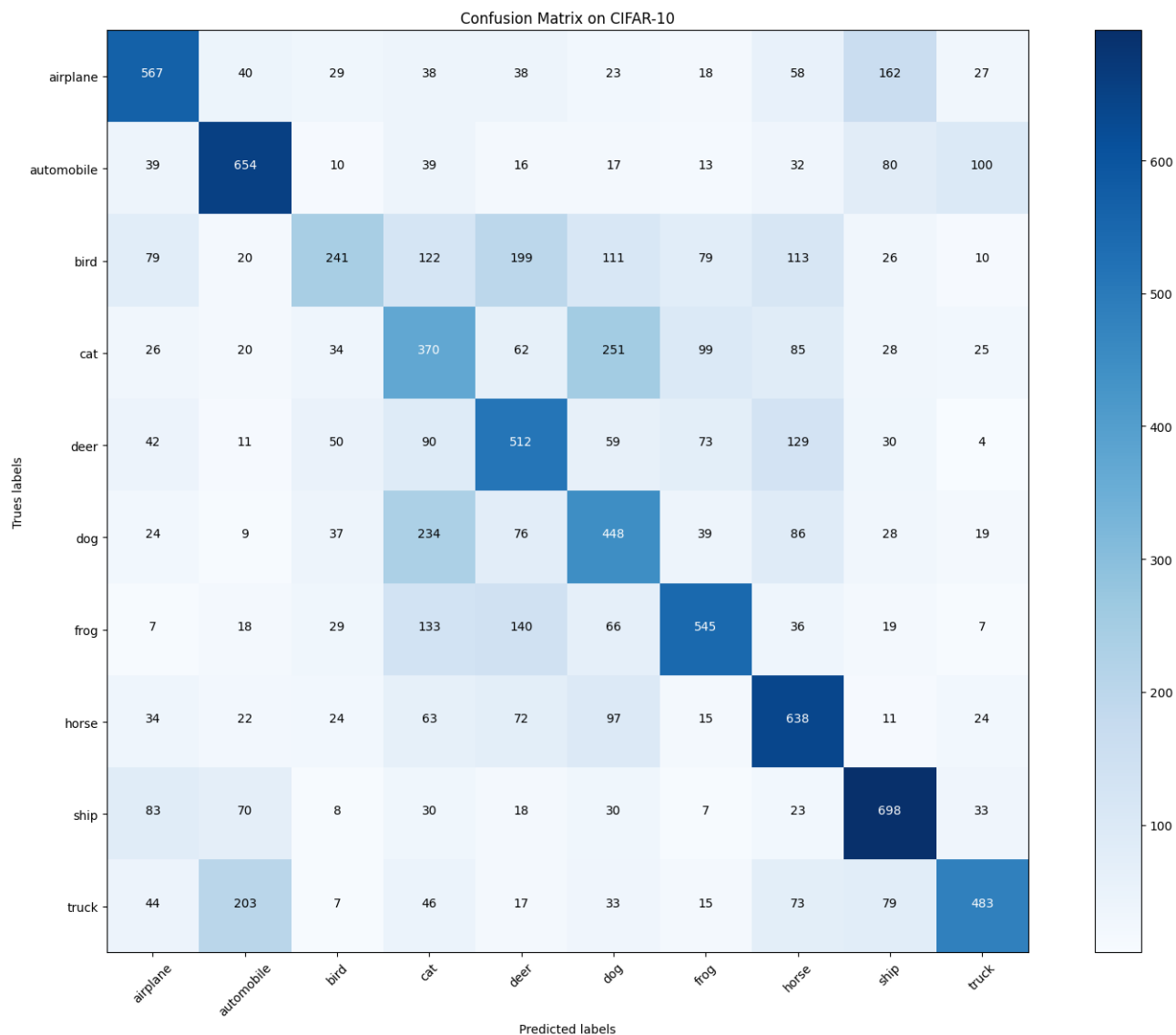


图 10: Confusion Matrix on CIFAR-10

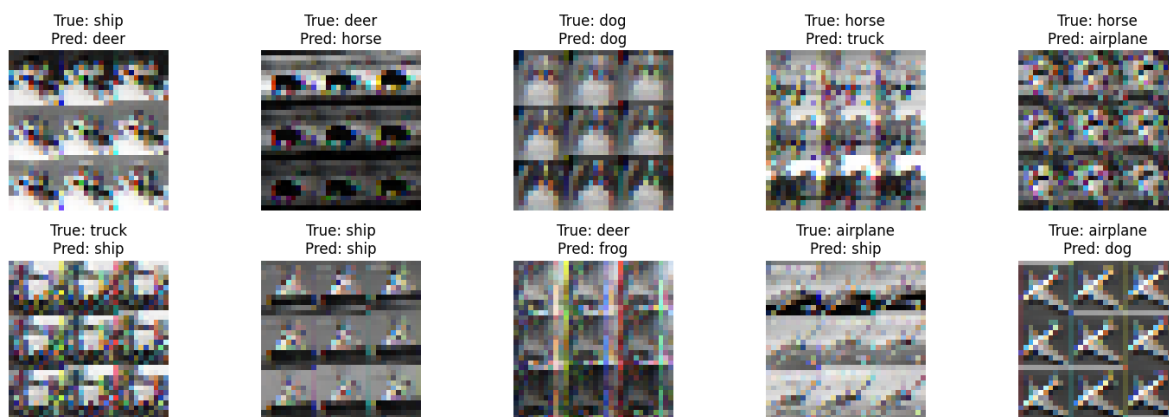


图 11: 随机测试数据展示分类效果

2. 收敛速度对比：基于 Mindspore 实现的神经网络在 CIFAR-10 数据集上的收敛速度较快，仅需 **10** 个 epoch 即可收敛，而手动实现的算法需要 **100** 个 epoch 才能收敛。
3. 模型复杂度对比：基于 Mindspore 实现的神经网络使用了多隐藏层的神经网络结构，模型复杂度更高，能够更好地捕捉数据集的复杂特征。
4. 优化算法对比：基于 Mindspore 实现的神经网络使用了 SGD、Adam 等高级的优化算法，收敛速度更快，性能更好。
5. 超参数调优对比：基于 Mindspore 实现的神经网络中的超参数经过充分调优，模型性能更佳。

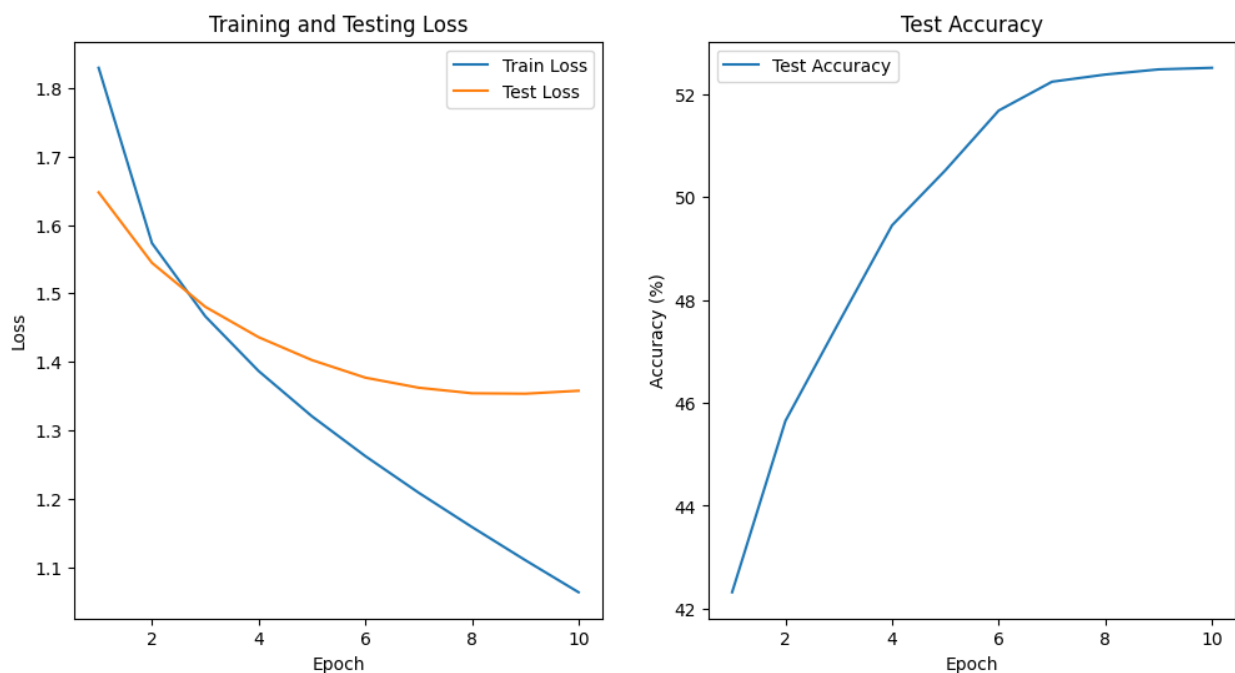


图 12: 基于 Mindspore 实现的神经网络在 CIFAR-10 数据集上的 Loss 和 Accuracy 曲线

出现上述差异可能的原因如下：

1. 模型复杂度：基于 Mindspore 实现的神经网络使用了多隐藏层的神经网络结构，模型复杂度更高，能够更好地捕捉数据集的复杂特征。

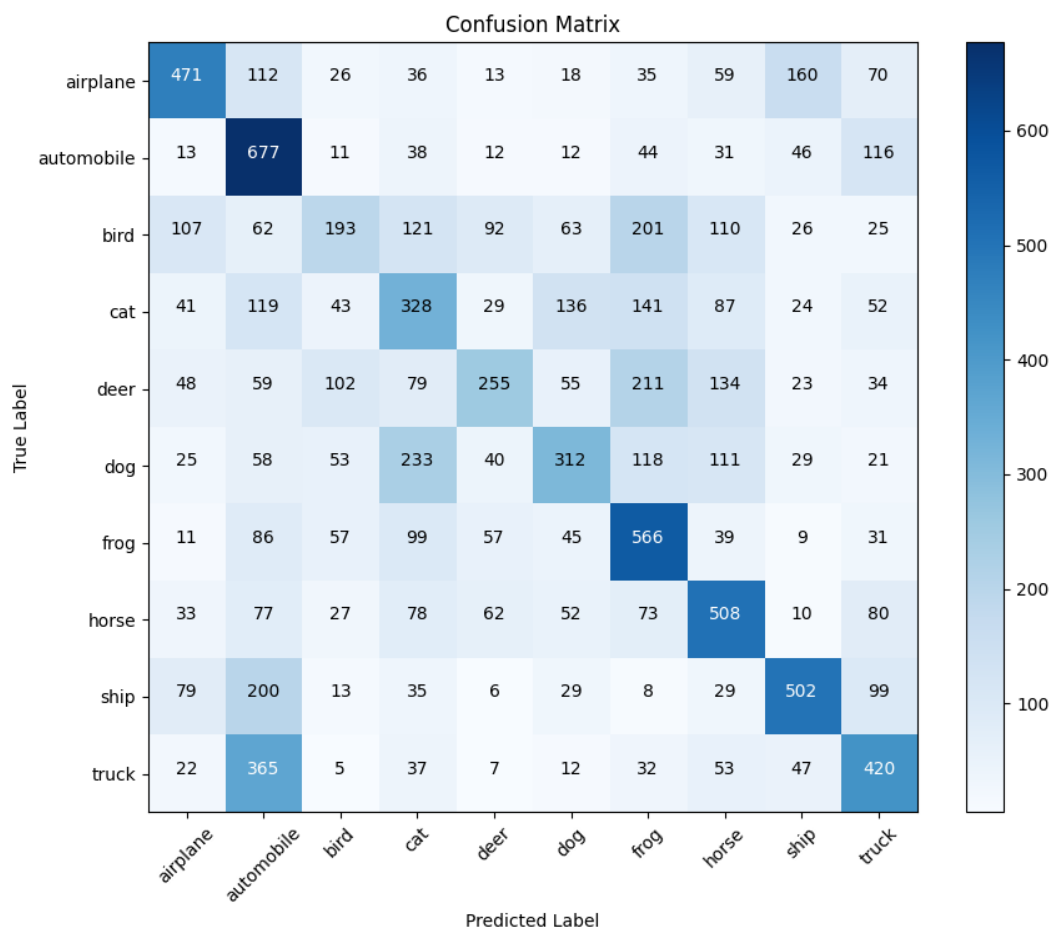


图 13: 基于 Mindspore 实现的神经网络在 CIFAR-10 数据集上的混淆矩阵

2. 优化算法：基于 Mindspore 实现的神经网络使用了 SGD、Adam 等高级的优化算法，收敛速度更快，性能更好。
3. 超参数调优：基于 Mindspore 实现的神经网络中的超参数经过充分调优，模型性能更佳。
4. 数据预处理：基于 Mindspore 实现的神经网络使用了更复杂的数据预处理技术，如归一化、标准化、变形等，能够更好地准备数据。
5. 自动微分：基于 Mindspore 实现的神经网络使用了自动微分技术，能够更方便地计算梯度，简化了模型训练过程。

### 3.6 MindSpore 学习使用心得体会

使用华为 MindSpore 平台进行深度学习开发与常规的手动实现（如直接使用 NumPy 等基础库编写神经网络及其训练流程）相比，带来了若干明显的优势和区别。以下是基于 MindSpore 平台进行深度学习开发的一些心得体会，特别是在实现与训练 BP 神经网络对 CIFAR-10 数据集进行分类的任务上的体会：

1. 数据加载优势：MindSpore 的 dataset 模块提供了丰富的数据操作接口，可以方便地从各种来源加载数据。使用 GeneratorDataset 创建自定义数据生成器更为简单，能够快速接入数据流。与手动实现数据加载相比，该方法更高效、更灵活。
2. 数据预处理和增强：通过 map 操作和 MindSpore 内置的 vision 模块，能够在数据流管道中轻松加入图像的预处理和增强步骤，如归一化、标准化、变形等。这些操作内部进行了优化，执行速度快，并且调用接口简单明了，易于理解和操作。
3. 训练和 BP 计算过程区别：在 MindSpore 中使用 nn.Cell 来定义网络模型更为简洁，在 construct 方法中定义数据通过网络的流动。这使得模型的构造直观而易于管理，与手动实现（如使用 NumPy 自行编写前向传播和反向传播）相比，大大降低了开发难度。



4. 自动微分：MindSpore 的自动微分机制（通过 `value_and_grad` 函数）使得计算损失函数对于参数的梯度变得异常简单。开发者无需手动推导和编程实现参数的梯度计算，框架会自动完成，这一点是传统手动实现所不能比拟的。
5. 优化器和参数更新：在 MindSpore 中，优化器已被封装成了内置的类，如 `nn.SGD`，并且可以直接与模型参数关联起来，调用优化器的语法简单，使得在训练循环中参数更新的实现非常高效并且易于理解。

总结：使用基于 MindSpore 的深度学习开发相比于传统的手动实现方式，在数据加载、预处理、模型构建和训练等多个方面都具有明显的优势。一方面，MindSpore 的丰富数据接口和功能强大的数据预处理功能简化了数据准备工作。另一方面，通过自动微分及优化器提供的高层封装，极大地减少了梯度计算和参数更新的编程工作量，提升了开发效率。更不用说，MindSpore 平台针对华为 Ascend 处理器的性能优化，对于训练效率可能带来额外的提升。

## 3.7 代码附录

### 3.7.1 cifar-10-scratch.ipynb

```
1  import pickle
2  import numpy as np
3  import os
4  import matplotlib.pyplot as plt
5
6  # 加载 CIFAR-10 批次数据
7  def load_cifar10_batch(filename):
8      with open(filename, 'rb') as f:
9          datadict = pickle.load(f, encoding='bytes')
10         X = datadict[b'data']
11         Y = datadict[b'labels']
12         X = X.reshape(10000, 3, 32, 32).astype("float")
13         Y = np.array(Y)
14         return X, Y
15
```

```

16 # 加载所有 CIFAR-10 数据
17 def load_cifar10(ROOT):
18     xs = []
19     ys = []
20     for b in range(1, 6):
21         f = os.path.join(ROOT, 'data_batch_%d' % (b,))
22         X, Y = load_cifar10_batch(f)
23         xs.append(X)
24         ys.append(Y)
25     Xtr = np.concatenate(xs)
26     Ytr = np.concatenate(ys)
27     del X, Y
28     Xte, Yte = load_cifar10_batch(os.path.join(ROOT, 'test_batch'))
29     return Xtr, Ytr, Xte, Yte
30
31 ROOT = './cifar-10-batches-py'
32 X_train, y_train, X_test, y_test = load_cifar10(ROOT)
33
34 # 全连接神经网络类
35 class FullyConnectedNN:
36     def __init__(self, input_size, hidden_size, output_size):
37         self.W1 = np.random.randn(input_size, hidden_size) * 0.01
38         self.b1 = np.zeros((1, hidden_size))
39         self.W2 = np.random.randn(hidden_size, output_size) * 0.01
40         self.b2 = np.zeros((1, output_size))
41
42     def relu(self, Z):
43         return np.maximum(0, Z)
44
45     def softmax(self, Z):
46         expZ = np.exp(Z - np.max(Z))
47         return expZ / expZ.sum(axis=1, keepdims=True)
48
49     def forward(self, X):
50         self.Z1 = np.dot(X, self.W1) + self.b1
51         self.A1 = self.relu(self.Z1)
52         self.Z2 = np.dot(self.A1, self.W2) + self.b2
53         self.A2 = self.softmax(self.Z2)
54         return self.A2
55
56     def compute_loss(self, Y, Y_hat):

```

```

57     m = Y.shape[0]
58     log_likelihood = -np.log(Y_hat[range(m), Y])
59     loss = np.sum(log_likelihood) / m
60     return loss
61
62     def backward(self, X, Y, Y_hat):
63         m = X.shape[0]
64         dZ2 = Y_hat
65         dZ2[range(m), Y] -= 1
66         dZ2 /= m
67
68         dW2 = np.dot(self.A1.T, dZ2)
69         db2 = np.sum(dZ2, axis=0, keepdims=True)
70
71         dA1 = np.dot(dZ2, self.W2.T)
72         dZ1 = dA1 * (self.Z1 > 0)
73
74         dW1 = np.dot(X.T, dZ1)
75         db1 = np.sum(dZ1, axis=0, keepdims=True)
76
77         self.W1 -= self.learning_rate * dW1
78         self.b1 -= self.learning_rate * db1
79         self.W2 -= self.learning_rate * dW2
80         self.b2 -= self.learning_rate * db2
81
82     def compute_accuracy(self, X, Y):
83         Y_hat = self.forward(X)
84         predictions = np.argmax(Y_hat, axis=1)
85         accuracy = np.mean(predictions == Y)
86         return accuracy
87
88     def train(self, X_train, Y_train, X_val, Y_val, epochs=300,
89 learning_rate=0.01):
90         self.learning_rate = learning_rate
91         train_losses = []
92         val_losses = []
93         val_accuracies = []
94
95         for epoch in range(epochs):
96             # 打乱训练数据
97             indices = np.arange(X_train.shape[0])

```

```

97         np.random.shuffle(indices)
98         X_train = X_train[indices]
99         Y_train = Y_train[indices]
100
101     # 小批量梯度下降
102     for start_idx in range(0, X_train.shape[0], batch_size):
103         end_idx = min(start_idx + batch_size, X_train.shape[0])
104         X_batch = X_train[start_idx:end_idx]
105         Y_batch = Y_train[start_idx:end_idx]
106
107         # 前向传播
108         Y_hat_train = self.forward(X_batch)
109         train_loss = self.compute_loss(Y_batch, Y_hat_train)
110
111         # 反向传播
112         self.backward(X_batch, Y_batch, Y_hat_train)
113
114     # 每个 epoch 结束后计算验证集上的损失和准确率
115     Y_hat_val = self.forward(X_val)
116     val_loss = self.compute_loss(Y_val, Y_hat_val)
117     val_accuracy = self.compute_accuracy(X_val, Y_val)
118
119     # 存储损失和准确率
120     train_losses.append(train_loss)
121     val_losses.append(val_loss)
122     val_accuracies.append(val_accuracy)
123
124     # 打印每个 epoch 的损失和准确率
125     print(f'Epoch [{epoch + 1}], train_loss: {train_loss:.4f},
126           val_loss: {val_loss:.4f}, val_acc: {val_accuracy:.4f}')
127
128     return train_losses, val_losses, val_accuracies
129
130 # 预处理数据
131 X_train = X_train.reshape(X_train.shape[0], -1) / 255.0
132 X_test = X_test.reshape(X_test.shape[0], -1) / 255.0
133
134 # 将训练数据分成训练集和验证集
135 split_index = int(0.8 * X_train.shape[0])
136 X_val, y_val = X_train[split_index:], y_train[split_index:]
137 X_train, y_train = X_train[:split_index], y_train[:split_index]

```

```

137
138 # 超参数
139 input_size = 3072 # 32*32*3
140 hidden_size = 100
141 output_size = 10
142 learning_rate = 0.01
143 epochs = 100
144 batch_size = 64
145
146 # 初始化并训练模型
147 model = FullyConnectedNN(input_size, hidden_size, output_size)
148 train_losses, val_losses, val_accuracies = model.train(X_train, y_train,
149     X_val, y_val, epochs, learning_rate)
150
151 # 在测试集上进行预测
152 y_pred = np.argmax(model.forward(X_test), axis=1)
153 accuracy = np.mean(y_pred == y_test)
154 print(f'测试集准确率: {accuracy:.4f}')
155
156 # 绘制损失曲线
157 epochs_range = range(epochs)
158 plt.figure(figsize=(8, 4))
159 plt.plot(epochs_range, train_losses, label='训练损失')
160 plt.plot(epochs_range, val_losses, label='验证损失')
161 plt.xlabel('Epochs')
162 plt.ylabel('Loss')
163 plt.legend(loc='upper right')
164 plt.title('损失曲线')
165 plt.show()
166
167 # 绘制准确率曲线
168 plt.figure(figsize=(8, 4))
169 plt.plot(epochs_range, val_accuracies, label='验证准确率')
170 plt.xlabel('Epochs')
171 plt.ylabel('Accuracy')
172 plt.legend(loc='lower right')
173 plt.title('准确率曲线')
174 plt.show()
175
176 # 计算混淆矩阵

```

```

177 def compute_confusion_matrix(y_true, y_pred, num_classes):
178     cm = np.zeros((num_classes, num_classes), dtype=int)
179     for i in range(len(y_true)):
180         cm[y_true[i], y_pred[i]] += 1
181     return cm
182
183 # 绘制混淆矩阵
184 def plot_confusion_matrix(cm, classes, title='混淆矩阵', cmap=plt.cm.Blues):
185     plt.figure(figsize=(16, 12))
186     plt.imshow(cm, interpolation='nearest', cmap=cmap)
187     plt.title(title)
188     plt.colorbar()
189     tick_marks = np.arange(len(classes))
190     plt.xticks(tick_marks, classes, rotation=45)
191     plt.yticks(tick_marks, classes)
192
193     fmt = 'd'
194     thresh = cm.max() / 2.
195     for i, j in np.ndindex(cm.shape):
196         plt.text(j, i, format(cm[i, j], fmt),
197                 horizontalalignment="center",
198                 color="white" if cm[i, j] > thresh else "black")
199
200     plt.ylabel('真实标签')
201     plt.xlabel('预测标签')
202     plt.tight_layout()
203     plt.show()
204
205 # CIFAR-10 标签
206 cifar10_labels = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', '
    frog', 'horse', 'ship', 'truck']
207
208 # 计算并绘制混淆矩阵
209 cm = compute_confusion_matrix(y_test, y_pred, len(cifar10_labels))
210 plot_confusion_matrix(cm, classes=cifar10_labels)

```

### 3.7.2 bp\_mindspore.ipynb

```

1 import os
2 import pickle
3 import numpy as np

```

```

4  import mindspore
5  from mindspore import nn
6  from mindspore import dataset as ds
7  from mindspore.dataset import vision, transforms
8  import matplotlib.pyplot as plt
9
10 path = './cifar-10-batches-py'
11
12 def load_cifar10_batch(file):
13     with open(file, 'rb') as fo:
14         dict = pickle.load(fo, encoding='bytes')
15         data = dict[b'data']
16         labels = dict[b'labels']
17         data = data.reshape(len(data), 3, 32, 32).transpose(0, 2, 3, 1)
18         return data, labels
19
20 def load_cifar10(path):
21     x_train = []
22     y_train = []
23     for i in range(1, 6):
24         data, labels = load_cifar10_batch(os.path.join(path, f'data_batch_{i}'))
25         x_train.append(data)
26         y_train.append(labels)
27     x_train = np.concatenate(x_train)
28     y_train = np.concatenate(y_train)
29
30     x_test, y_test = load_cifar10_batch(os.path.join(path, 'test_batch'))
31
32     return (x_train, y_train), (x_test, y_test)
33
34 def create_dataset(data, labels, batch_size, shuffle=True):
35     def generator():
36         for i in range(len(data)):
37             yield data[i], labels[i]
38
39     dataset = ds.GeneratorDataset(generator, ["image", "label"], shuffle=
shuffle)
40     image_transforms = [
41         vision.Rescale(1.0 / 255.0, 0),

```

```

42         vision.Normalize(mean=(0.4914, 0.4822, 0.4465), std=(0.2023, 0.1994,
43         0.2010)),
44         vision.HWC2CHW()
45     ]
46     label_transform = transforms.TypeCast(mindspore.int32)
47
48     dataset = dataset.map(operations=image_transforms, input_columns='image'
49     )
50     dataset = dataset.map(operations=label_transform, input_columns='label')
51     dataset = dataset.batch(batch_size)
52
53     return dataset
54
55 (train_images, train_labels), (test_images, test_labels) = load_cifar10(path
56 )
57 train_dataset = create_dataset(train_images, train_labels, batch_size=64)
58 test_dataset = create_dataset(test_images, test_labels, batch_size=64)
59
60 class Network(nn.Cell):
61     def __init__(self):
62         super().__init__()
63         self.flatten = nn.Flatten()
64         self.dense_relu_sequential = nn.SequentialCell(
65             nn.Dense(32*32*3, 512),
66             nn.ReLU(),
67             nn.Dense(512, 512),
68             nn.ReLU(),
69             nn.Dense(512, 10)
70         )
71
72     def construct(self, x):
73         x = self.flatten(x)
74         logits = self.dense_relu_sequential(x)
75         return logits
76
77 model = Network()
78 epochs = 10
79 batch_size = 64
80 learning_rate = 1e-2
81 loss_fn = nn.CrossEntropyLoss()
82 optimizer = nn.SGD(model.trainable_params(), learning_rate=learning_rate)

```



```

80
81 # Define forward function
82 def forward_fn(data, label):
83     logits = model(data)
84     loss = loss_fn(logits, label)
85     return loss, logits
86
87 # Get gradient function
88 grad_fn = mindspore.value_and_grad(forward_fn, None, optimizer.parameters,
89                                     has_aux=True)
89
90 # Define function of one-step training
91 def train_step(data, label):
92     (loss, _), grads = grad_fn(data, label)
93     optimizer(grads)
94     return loss
95
96 def train_loop(model, dataset):
97     size = dataset.get_dataset_size()
98     model.set_train()
99     total_loss = 0
100     for batch, (data, label) in enumerate(dataset.create_tuple_iterator()):
101         loss = train_step(data, label)
102         total_loss += loss.asnumpy()
103
104     return total_loss / size
105
106 def test_loop(model, dataset, loss_fn):
107     num_batches = dataset.get_dataset_size()
108     model.set_train(False)
109     total, test_loss, correct = 0, 0, 0
110     all_preds = []
111     all_labels = []
112     for data, label in dataset.create_tuple_iterator():
113         pred = model(data)
114         total += len(data)
115         test_loss += loss_fn(pred, label).asnumpy()
116         correct += (pred.argmax(1) == label).asnumpy().sum()
117         all_preds.extend(pred.argmax(1).asnumpy())
118         all_labels.extend(label.asnumpy())
119     test_loss /= num_batches

```

```

120     correct /= total
121     return test_loss, correct, all_preds, all_labels
122
123 def compute_confusion_matrix(y_true, y_pred, num_classes):
124     conf_matrix = np.zeros((num_classes, num_classes), dtype=np.int32)
125     for true, pred in zip(y_true, y_pred):
126         conf_matrix[true, pred] += 1
127     return conf_matrix
128
129 train_losses = []
130 test_losses = []
131 test_accuracies = []
132 all_preds = []
133 all_labels = []
134
135 for t in range(epochs):
136     train_loss = train_loop(model, train_dataset)
137     test_loss, test_accuracy, epoch_preds, epoch_labels = test_loop(model,
138 test_dataset, loss_fn)
139     train_losses.append(train_loss)
140     test_losses.append(test_loss)
141     test_accuracies.append(test_accuracy)
142     all_preds.extend(epoch_preds)
143     all_labels.extend(epoch_labels)
144     print(f"Epoch [{t}], train_loss: {train_loss:.4f}, val_loss: {test_loss
145 :.4f}, val_acc: {test_accuracy:.4f}")
146 print("Done!")
147
148 # Plotting
149 epochs_range = range(1, epochs + 1)
150 plt.figure(figsize=(12, 6))
151
152 plt.subplot(1, 2, 1)
153 plt.plot(epochs_range, train_losses, label='Train Loss')
154 plt.plot(epochs_range, test_losses, label='Test Loss')
155 plt.xlabel('Epoch')
156 plt.ylabel('Loss')
157 plt.title('Training and Testing Loss')
158 plt.legend()
159
160 plt.subplot(1, 2, 2)

```

```

159 plt.plot(epochs_range, [acc * 100 for acc in test_accuracies], label='Test
    Accuracy')
160 plt.xlabel('Epoch')
161 plt.ylabel('Accuracy (%)')
162 plt.title('Test Accuracy')
163 plt.legend()
164
165 plt.show()
166
167 # Confusion Matrix
168 cifar10_labels = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', '
    frog', 'horse', 'ship', 'truck']
169 conf_matrix = compute_confusion_matrix(all_labels[:10000], all_preds
    [:10000], 10)
170 plt.figure(figsize=(10, 8))
171 plt.imshow(conf_matrix, interpolation='nearest', cmap=plt.cm.Blues)
172 plt.title('Confusion Matrix')
173 plt.colorbar()
174 tick_marks = np.arange(10)
175 plt.xticks(tick_marks, cifar10_labels, rotation=45)
176 plt.yticks(tick_marks, cifar10_labels)
177 plt.xlabel('Predicted Label')
178 plt.ylabel('True Label')
179
180 # Add text annotations to the confusion matrix
181 thresh = conf_matrix.max() / 2
182 for i in range(conf_matrix.shape[0]):
183     for j in range(conf_matrix.shape[1]):
184         plt.text(j, i, format(conf_matrix[i, j], 'd'),
185                 horizontalalignment='center',
186                 color='white' if conf_matrix[i, j] > thresh else 'black')
187
188 plt.show()

```

## 4 心得体会

本学期的三个实验很好的锻炼了我的代码能力，我通过这些代码的手动实现，更加深刻的理解了课上的专业知识点，如朴素贝叶斯实现分类和预测，KNN 实现分类和预测，神经网络的实现等。同时，通过实验我也学会了如何使用 mindspore 平台。

在使用华为 mindspore 平台时，我觉得数据集的加载相比于手动实现的方式更加方便，同时 mindspore 平台提供了丰富的 API 接口，可以方便地实现神经网络的训练和评估。在本次实验中，我实现了一个基于全连接神经网络的分类器，并在 CIFAR-10 数据集上进行了训练和评估。通过本次实验，我对神经网络的训练过程有了更深入的理解，同时也学习了如何使用 mindspore 平台实现神经网络模型。希望在以后的实验中，我能够更加熟练地使用 mindspore 平台，实现更加复杂的神经网络模型。但是相比于 PyTorch, Tensorflow 等主流深度学习算法框架，mindspore 仍有部分数据集的加载仍然比较复杂，希望有关开发人员能够进一步优化 API 接口，提高 mindspore 的易用性。