

Task1

- **什么是多层感知机（MLP）？其结构是怎样的？**

答：多层感知机（Multi-layer Perceptron，MLP）是一种深度学习模型，它由多个感知器组成，能够实现更复杂的分类任务。多层感知机由输入层、隐藏层和输出层组成。输入层负责接收外部输入的数据，隐藏层通过非线性激活函数（如ReLU、sigmoid等）将输入转换为更抽象的表示，输出层则根据隐藏层的输出做出最终的分类决策。通过调整隐藏层的参数，多层感知机能够学习到更复杂的分类规则。

- **数据在神经网络中扮演哪些角色？（数据集的 split 和处理）**

答：在随机选择神经网络训练数据的过程中，首要步骤是**数据采集**。数据采集需要确保数据的代表性和多样性。这意味着我们需要从不同的来源收集数据，并且应避免选择过于相似或重复的数据。

此外，在数据采集过程中，我们还需要关注**数据的预处理**。预处理包括数据清洗、标准化、归一化等步骤，旨在确保数据的质量和可靠性。

接下来是**数据选择**。在庞大的数据集中，我们不可能使用全部数据来训练神经网络。因此，我们需要选择一部分数据来进行训练。数据选择应基于训练任务的特定需求。

在**数据评估**方面，我们需要确定所选择的训练数据的质量和数量是否满足神经网络训练的需求。这涉及到对数据的全面评估，包括数据的准确性、多样性和完整性等方面。同时，我们还需要关注**数据的分布情况**，以确保所选择的训练数据能够反映总体的数据分布。

- **噪声是什么？特征是什么？标签是什么？**

答：在深度学习中，**噪声**是指数据或模型中的随机扰动或不准确性，这些扰动可能会影响模型的训练和预测性能。噪声通常来源于数据采集过程中的随机误差、标签错误或模型训练中的随机性。主要有三种可能的噪声，**数据噪声，模型噪声，训练噪声**

特征指的是数据中蕴含的、为解决特定任务（如分类、回归、生成等）有价值的信息。这些信息可以是数据的内在规律、模式或属性。它是模型从原始数据中“提炼”出的抽象表示，不同层级的特征对应数据不同层次的语义或结构。

数据标签，也称为数据注释，是指在数据集上进行详细的标记或标注的过程。这些标注在机器学习模型的训练阶段起到指导作用，帮助模型学习如何准确地进行预测和分类。标签的准确性决定了模型在处理新的、未标记数据时的表现。

- **Batch size 是什么？为什么堆叠成 Batch 可以提高运算速度？**

答：Batch Size，即**批量大小**，是指在训练神经网络时，每次迭代所使用的样本数量。当Batch Size为1时，称为在线学习（Online Learning）；当Batch Size等于整个数据集的大小时，称为全批量学习（Full-Batch Learning）；而当Batch Size介于1和整个数据集大小之间时，称为小批量学习（Mini-Batch Learning）。

Batch Size越大，每次迭代所使用的样本数量越多，训练速度越快。然而，过大的Batch Size可能导致内存不足，从而影响训练速度。

设置Batch

1. 数据集规模：对于小数据集，可以尝试使用较大的Batch Size，以充分利用计算资源并提高训练速度。对于大数据集，Batch Size可以相对较大，但仍需考虑内存限制。
2. 模型架构：不同的模型架构对Batch Size的敏感度不同。例如，一些模型在较小的Batch Size下表现更好，而另一些模型则可能需要较大的Batch Size来达到最佳性能。

3. 硬件资源：计算资源是设置Batch Size的重要考虑因素。如果GPU显存有限，过大的Batch Size可能导致内存溢出。因此，在设置Batch Size时，需要充分考虑硬件资源的限制。
4. 训练目标：不同的训练目标可能需要不同的Batch Size。例如，在追求快速收敛的情况下，可以尝试使用较大的Batch Size；而在追求模型性能的情况下，可能需要使用较小的Batch Size。

- **神经元是什么？**

答：在深度学习中，**神经元**（或称为人工神经元）是构成神经网络的基本单元。它们的设计灵感来源于生物神经元，旨在模拟人脑的信号传递机制。每个神经元接收来自其他神经元或外部输入的信号，并通过加权和激活函数处理这些信号，从而生成输出。

神经元的结构

一个典型的神经元包括以下几个部分：

- **输入**：神经元接收的信号，通常是来自其他神经元的输出或外部数据。
- **权重**：每个输入都有一个与之相关的权重，表示该输入对神经元输出的重要性。
- **偏置**：一个额外的参数，用于调整神经元的输出，使其能够独立于输入进行调整。
- **激活函数**：对加权输入的总和进行非线性变换，决定神经元的输出。

- **什么是激活函数？常见的激活函数有哪些？什么叫“非线性表达能力”？**

答：激活函数是人工神经网络中用于将神经元输入映射到输出的函数。它通过引入非线性特性，使神经网络能够学习和表达复杂的非线性关系。如果没有激活函数，神经网络的每一层仅是线性变换的叠加，无法处理复杂的非线性问题。

常见激活函数及其特点

Sigmoid 函数

Sigmoid 是一种 S 型函数，输出范围为 0 到 1，常用于二分类问题。其优点是梯度平滑，便于求导，但容易导致梯度消失问题，且计算复杂度较高。

Tanh 函数

Tanh 是双曲正切函数，输出范围为 -1 到 1。它以 0 为中心，相较于 Sigmoid 更适合实际应用，但仍存在梯度消失问题。

ReLU 函数

ReLU (Rectified Linear Unit) 是目前最常用的激活函数之一。它在输入为正时输出线性，计算简单且能有效解决梯度消失问题。但可能出现“神经元死亡”现象，即某些神经元在训练中永远不被激活。

Leaky ReLU 函数

Leaky ReLU 是 ReLU 的改进版本，允许负输入有一个小的非零斜率，从而解决 ReLU 的“神经元死亡”问题。

ELU 函数

ELU (Exponential Linear Unit) 在负输入时采用指数函数，能更好地控制均值和方差，但计算复杂度较高。

SELU 函数

SELU (Scaled Exponential Linear Unit) 通过自归一化特性实现内部归一化，适用于自归一化神经网络 (SNN)，能加速收敛。

Swish 函数

Swish 是一种平滑且非单调的激活函数，能增强模型的表达能力，适用于复杂任务。

Mish 函数

Mish 是 Swish 的改进版本，具有更高的平滑性，但计算复杂度更高。

Softmax 函数

Softmax 常用于多分类问题的输出层，将输出值映射为概率分布，便于分类任务。

非线性表达能力是指模型通过非线性方法来拟合复杂数据分布和特征，从而提高其学习和预测能力。

非线性模型能够更好地拟合复杂的函数曲线，这是因为它们具有多种不规则的表达方式。例如，深度学习中的非线性激活函数（如ReLU、sigmoid等）使得神经网络能够学习到数据中的复杂模式和特征，从而提高模型的泛化能力和性能。

• 什么是计算图？它和数据结构/离散数学中学的图有什么区别？怎么构建计算图？

答：

◦ 第一部分：什么是计算图 (Computation Graph)?

核心答案：计算图是一种用于**描述和评估数学表达式**的有向图。在深度学习的语境下，它特指一个**有向无环图 (Directed Acyclic Graph, DAG)**，其中：

- **节点 (Nodes)** 代表**数据 (如张量/Tensor) *或*数学运算 (如加法、乘法、矩阵乘法、激活函数)**。
- **边 (Edges)** 代表**数据的流动方向**，即一个节点的输出，作为另一个节点的输入。

本质上，任何复杂的神经网络模型，无论其结构多深、多宽，都可以被抽象成一个清晰的计算图。这个图精确地定义了输入数据是如何一步步地经过各种运算，最终转变为输出结果的。

比喻：计算图就像一个**“自动化工厂的流水线蓝图”**。

- **原材料 (输入数据)** 从流水线的一端进入。
- **每一个工位 (节点)** 代表一道加工工序（一个数学运算），比如“切割”、“焊接”、“喷漆”。
- **传送带 (边)** 则规定了半成品从一个工位流向下一个工位的顺序。
- 最终，从流水线末端出来的就是**成品 (输出预测)**。

◦ 第二部分：它和数据结构/离散数学中学的图有什么区别？

核心答案：计算图是图论概念在特定领域（数值计算）的一种**应用和特化**，它们的根本结构是一致的，但**核心区别在于节点的“含义”和图的“用途”**。

特 性	数据结构/离散数学中的“通用图”	深度学习中的“计算图”
节 点 含 义	可以是任何抽象实体，如城市、人、网页、任务等。	特指： 数学运算 (Operations) 或数据/张量 (Tensors)。
边 的 含 义	代表实体间的某种关系，如道路、友谊、链接、依赖等。	特指： 数据（张量）的流动方向。

特性	数据结构/离散数学中的“通用图”	深度学习中的“计算图”
图的结构	可以是任意结构，有向、无向、有环、无环皆可。	特指： 通常是有向无环图 (DAG)，因为计算过程是单向、不可逆的。
核心用途	用于建模和分析实体间的关系，如寻找最短路径、社区发现。	特指： 1. 前向传播： 按照图的顺序执行运算，得到最终结果。2. 反向传播： 沿着图的反方向，高效地计算梯度。

总结来说：计算图是“图”这个泛用数据结构的一个“特例”，它被赋予了特定的数学含义，并为深度学习中的两大核心算法——前向传播和反向传播——提供了理论基础和工程实现的框架。

○ **第三部分：怎么构建计算图？**

在现代深度学习框架（如 PyTorch, TensorFlow）中，计算图的构建过程被高度自动化了，主要分为两种方式：

1. **静态图 (Static Graph) - “先画图，再施工”** (代表：TensorFlow 1.x)

- **构建方式：**你需要先完整地定义好整个计算图的结构。就像先用蓝图软件把整个工厂的流水线都设计好，检查无误后，再把这个固定的蓝图交给工厂去执行。
- **执行方式：**一旦图被构建和“编译”好，它就会被优化并固化下来。之后你只需要向这个固定的图里“喂”数据即可，图的结构本身不能再改变。
- **优点：**便于优化，部署效率高。
- **缺点：**不够灵活，调试困难。

2. **动态图 (Dynamic Graph) - “边施工，边画图”** (代表：PyTorch, TensorFlow 2.x)

- **构建方式：**你不需要预先定义图。计算图是在你的代码运行时被动态、即时地构建出来的。你每执行一行运算（比如 $c = a + b$ ），框架就在后台默默地为你画出相应的节点和边。
- **执行方式：**代码的执行过程，就是图的构建过程。这使得你可以使用 Python 的原生控制流（如 if 语句、for 循环）来动态地改变图的结构。
- **优点：**极其灵活，所见即所得，调试非常直观和方便。
- **缺点：**性能优化和部署相对静态图要更复杂一些。

在 PyTorch 中，你写的每一行模型代码，比如 `x = self.layer1(x)`，实际上就是在**动态地构建**计算图的一部分。这个过程对于开发者来说是**透明的、无感的**，这也是 PyTorch 如此受欢迎的重要原因。

● **怎么计算MLP的参数数量？什么是超参数？MLP有哪些超参数？**

○ **如何计算MLP的参数数量？**

核心答案：MLP的参数数量主要是其**所有线性层 (nn.Linear) 的权重 (weights) 和偏置 (biases) 数量的总和**。

详细解释:

一个从 in_features 个输入神经元连接到 out_features 个输出神经元的线性层，其参数量计算如下:

- **权重 (Weights):** 每个输入神经元都与每个输出神经元有一条连接，每条连接都有一个权重。所以权重数量是 $\text{in_features} * \text{out_features}$ 。
- **偏置 (Biases):** 每个输出神经元都有一个自己的偏置项。所以偏置数量是 out_features 。
- **总参数量 = $(\text{in_features} * \text{out_features}) + \text{out_features}$**

举例:

一个MLP有两层: 一个输入层 (784个节点), 一个隐藏层 (128个节点), 一个输出层 (10个节点)。

- 第一层 (784 -> 128) 的参数量 = $(784 * 128) + 128 = 100480$
- 第二层 (128 -> 10) 的参数量 = $(128 * 10) + 10 = 1290$
- **模型总参数量 = $100480 + 1290 = 101770$**

◦ 什么是超参数 (Hyperparameters)?

核心答案: 超参数是在**开始学习过程之前设置的、用于控制学习过程本身的外部参数**。它们不是模型通过训练学到的，而是由我们 (人类) 设定的。

比喻: 如果说模型的参数 (权重和偏置) 是学生通过学习记在脑子里的“**知识点**”，那么超参数就是我们为这个学生制定的“**学习计划**”和“**备考策略**” (比如每天学几小时、用什么参考书、做几套模拟题)。

◦ MLP有哪些超参数?

- **网络结构相关:**
 - **隐藏层的数量:** 决定了网络的深度。
 - **每个隐藏层的神经元数量:** 决定了每层的“宽度”。
- **训练过程相关:**
 - **学习率 (Learning Rate):** 控制模型每次更新参数的步长。
 - **优化器 (Optimizer):** 更新参数的具体算法，如 SGD, Adam。
 - **损失函数 (Loss Function):** 评估模型预测好坏的标准。
 - **激活函数 (Activation Function):** 如 ReLU, Sigmoid。
 - **训练的轮数 (Epochs):** 整个数据集被重复学习的次数。
 - **批次大小 (Batch Size):** 每次喂给模型进行学习的数据样本数量。

• 什么是隐藏层 (hidden layers) ? 它为什么叫这个名字?

◦ **核心答案:** 隐藏层是神经网络中**位于输入层和输出层之间的所有层**。

◦ 为什么叫“隐藏”层?

因为它对于使用者来说是**不可见的**。我们只能定义和观察到“输入” (比如一张图片的数据) 和最终的“输出” (比如识别结果是“猫”)，而数据在这些中间层所进行的复杂转换和提取的特征，就像发生在一个“黑箱”内部，我们无法直接看到，因此称之为“隐藏”层。

◦ **比喻:** 就像你的大脑进行思考。别人问你一个问题 (输入)，你给出一个答案 (输出)。但你脑海中神经元之间复杂的思考、联想、推理过程 (隐藏层)，对提问者来说是完全“隐藏”的。

• 什么是损失函数? 什么任务用什么损失函数?

- **什么是损失函数 (Loss Function)?**

核心答案：损失函数是一个用来衡量模型预测值与真实值之间“差距”或“误差”的函数。它的输出值（损失/Loss）越大，说明模型预测得越差。

比喻：损失函数就像一位“阅卷老师”。你（模型）交上你的答案（预测值），老师会根据标准答案（真实值）给你打一个“扣分”（损失值）。你的目标就是通过学习，让这个“扣分”越低越好。

- **什么任务用什么损失函数？**

- **回归任务 (Regression Task)** - 预测一个连续的数值（比如房价、温度）：

- **均方误差损失 (Mean Squared Error, MSE)**：最常用。它计算的是预测值与真实值之差的平方的平均值。对较大的误差给予更重的惩罚。

- **分类任务 (Classification Task)** - 预测一个类别（比如是猫还是狗）：

- **交叉熵损失 (Cross-Entropy Loss)**：最常用。

- **二元交叉熵 (Binary Cross-Entropy)**：用于二分类问题（是/否，垃圾邮件/非垃圾邮件）。

- **分类交叉熵 (Categorical Cross-Entropy)**：用于多分类问题（图片是猫、狗、还是马）。

- **前向传播、梯度、学习率、反向传播、优化器是什么？**

- **前向传播 (Forward Propagation)**：数据从输入层开始，逐层穿过网络（包括所有隐藏层），经过计算和激活函数处理，最终到达输出层得到一个**预测结果**的过程。这是**模型进行预测**的步骤。

- **梯度 (Gradient)**：计算出损失后，梯度告诉我们“**为了让损失变小，每个参数（权重/偏置）应该朝着哪个方向、调整多少**”。它本质上是损失函数对每个参数的偏导数，指向了损失上升最快的方向，我们沿着它的反方向更新参数就能减小损失。

- **学习率 (Learning Rate)**：是一个超参数，它控制着我们根据梯度更新参数时的“**步子迈多大**”。学习率太高，可能会“一步迈过头”，导致不稳定；太低，则学习速度太慢。

- **反向传播 (Backward Propagation)**：是一个高效计算网络中所有参数梯度的**算法**。它从最终的损失值开始，像链式反应一样，一层一层地把“误差信号”反向传播回网络中的每一个参数，并计算出每个参数对应的梯度。这是**计算如何修正**的步骤。

- **优化器 (Optimizer)**：是**执行参数更新**的具体算法。它拿到反向传播计算出的梯度和我们设定的学习率，然后用某种策略（比如最简单的SGD，或者更智能的Adam）去更新模型的每一个参数。常见的有：

- **SGD (Stochastic Gradient Descent)**：最基础的优化器。

- **Adam (Adaptive Moment Estimation)**：目前最流行、最通用的优化器，它能为不同参数自动适应学习率。

比喻总结：

你做了一道题（**前向传播**），老师告诉你扣了20分（**损失**），并告诉你错在哪里、该怎么改（**梯度**）。你决定这次要认真听讲、改正错误的80%（**学习率**）。老师把这个扣分点从最终答案一步步反推到你最初的草稿，让你明白每一步错在哪里（**反向传播**）。最后，你拿出你的“错题本”（**优化器**），用你自己的学习方法（比如 Adam 策略）把知识点记下来（更新参数）。

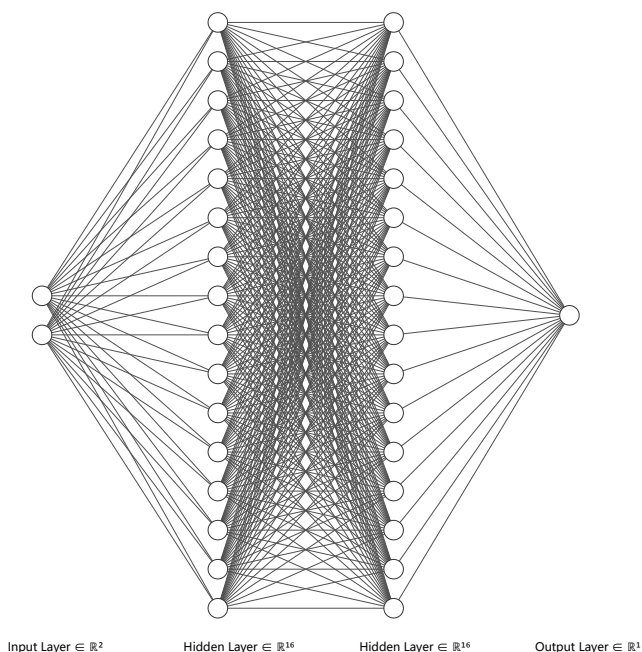
- **归一化是什么？正则化是什么？**

- **归一化 (Normalization)**：是一种**数据预处理**技术。它将输入数据的不同特征缩放到一个相同的、标准的范围（如0到1，或均值为0、方差为1）。

- **目的**：消除不同特征之间量纲（单位和尺度）的影响，让模型能更公平、更快速地学习。
- **比喻**：在计算一个人的健康指数时，如果身高用“米”（约1.7），体重用“克”（约60000），那么体重这个特征的数值会主导整个计算。归一化就像是把身高和体重都转换成一个0到100的评分，让它们在同一起跑线上被模型学习。
- **正则化 (Regularization)**：是一种在模型训练过程中，用于**防止过拟合**的技术。它通过在损失函数中添加一个“惩罚项”来限制模型的复杂度。
 - **目的**：让模型学习到更通用、更简单的规律，而不是死记硬背训练数据的细节和噪声。
 - **比喻**：老师在教学生时，不仅要求他们“答对题”（最小化损失），还要求他们“解题思路要简洁、普适”（最小化正则化惩罚项），不许用那些只能解这道题的“歪门邪道”。常见的有 L1 和 L2 正则化。
- **什么是欠拟合？什么是过拟合？**
 - **欠拟合 (Underfitting)**：指模型**过于简单**，以至于无法捕捉到数据中的基本规律。它在训练集和测试集上**表现都很差**。
 - **比喻**：一个没怎么学习的学生，给他做练习题（训练集）和期末考试题（测试集），他都考不及格。
 - **过拟合 (Overfitting)**：指模型**过于复杂**，以至于把训练数据中的噪声和偶然特征都当作规律学了进去。它在训练集上**表现极好**，但在新的、未见过的数据（测试集）上**表现很差**。
 - **比喻**：一个只会死记硬背的学生，把练习册上的所有题（包括答案的错别字）都背下来了，所以练习题能考100分。但期末考试题型稍微一变，他就完全不会了，考了不及格。他“记住”了数据，却没有“理解”规律。

代码实践部分，相关代码已经上传到Github

- 我实现的神经网络是一个全连接神经网络即MLP，具体结构见图，使用NN.SVG生成含有一个输入层，含有两个神经元结点，两个隐藏层每个隐藏层含有16个神经元结点，一个输出层含有一个节点



- **尝试调整noise大小观察训练结果**
 - **第一轮**：设置noise=0.2，训练结果如下

- | | | | | | |
|-------|-----------|-------|--------|-----------|--------|
| Epoch | [10/100] | Loss: | 0.5803 | Accuracy: | 0.8213 |
| Epoch | [20/100] | Loss: | 0.3803 | Accuracy: | 0.8313 |
| Epoch | [30/100] | Loss: | 0.2778 | Accuracy: | 0.8750 |
| Epoch | [40/100] | Loss: | 0.2356 | Accuracy: | 0.8962 |
| Epoch | [50/100] | Loss: | 0.1977 | Accuracy: | 0.9175 |
| Epoch | [60/100] | Loss: | 0.1580 | Accuracy: | 0.9463 |
| Epoch | [70/100] | Loss: | 0.1255 | Accuracy: | 0.9613 |
| Epoch | [80/100] | Loss: | 0.1022 | Accuracy: | 0.9637 |
| Epoch | [90/100] | Loss: | 0.0886 | Accuracy: | 0.9675 |
| Epoch | [100/100] | Loss: | 0.0803 | Accuracy: | 0.9700 |

- 第二轮：设置noise=0.3，训练结果如下

- | | | | | | |
|-------|-----------|-------|--------|-----------|--------|
| Epoch | [10/100] | Loss: | 0.5256 | Accuracy: | 0.7712 |
| Epoch | [20/100] | Loss: | 0.3707 | Accuracy: | 0.8175 |
| Epoch | [30/100] | Loss: | 0.3299 | Accuracy: | 0.8625 |
| Epoch | [40/100] | Loss: | 0.3229 | Accuracy: | 0.8588 |
| Epoch | [50/100] | Loss: | 0.3076 | Accuracy: | 0.8662 |
| Epoch | [60/100] | Loss: | 0.2933 | Accuracy: | 0.8725 |
| Epoch | [70/100] | Loss: | 0.2745 | Accuracy: | 0.8812 |
| Epoch | [80/100] | Loss: | 0.2531 | Accuracy: | 0.8938 |
| Epoch | [90/100] | Loss: | 0.2340 | Accuracy: | 0.9038 |
| Epoch | [100/100] | Loss: | 0.2201 | Accuracy: | 0.9150 |

- 第三轮：设置noise=0.4，训练结果如下

- | | | | | | |
|-------|-----------|-------|--------|-----------|--------|
| Epoch | [10/100] | Loss: | 0.5369 | Accuracy: | 0.8200 |
| Epoch | [20/100] | Loss: | 0.3847 | Accuracy: | 0.8363 |
| Epoch | [30/100] | Loss: | 0.3582 | Accuracy: | 0.8488 |
| Epoch | [40/100] | Loss: | 0.3383 | Accuracy: | 0.8562 |
| Epoch | [50/100] | Loss: | 0.3252 | Accuracy: | 0.8625 |
| Epoch | [60/100] | Loss: | 0.3159 | Accuracy: | 0.8588 |
| Epoch | [70/100] | Loss: | 0.3119 | Accuracy: | 0.8638 |
| Epoch | [80/100] | Loss: | 0.3099 | Accuracy: | 0.8650 |
| Epoch | [90/100] | Loss: | 0.3083 | Accuracy: | 0.8625 |
| Epoch | [100/100] | Loss: | 0.3077 | Accuracy: | 0.8625 |

- 第四轮：设置noise=0.1，训练结果如下

- | | | | | | |
|-------|-----------|-------|--------|-----------|--------|
| Epoch | [10/100] | Loss: | 0.5589 | Accuracy: | 0.7925 |
| Epoch | [20/100] | Loss: | 0.3796 | Accuracy: | 0.8100 |
| Epoch | [30/100] | Loss: | 0.2612 | Accuracy: | 0.8612 |
| Epoch | [40/100] | Loss: | 0.1871 | Accuracy: | 0.9200 |
| Epoch | [50/100] | Loss: | 0.1502 | Accuracy: | 0.9375 |
| Epoch | [60/100] | Loss: | 0.1090 | Accuracy: | 0.9537 |
| Epoch | [70/100] | Loss: | 0.0711 | Accuracy: | 0.9812 |
| Epoch | [80/100] | Loss: | 0.0440 | Accuracy: | 0.9938 |
| Epoch | [90/100] | Loss: | 0.0279 | Accuracy: | 0.9962 |
| Epoch | [100/100] | Loss: | 0.0191 | Accuracy: | 0.9962 |

- 第五轮：设置noise=0.05，训练结果如下

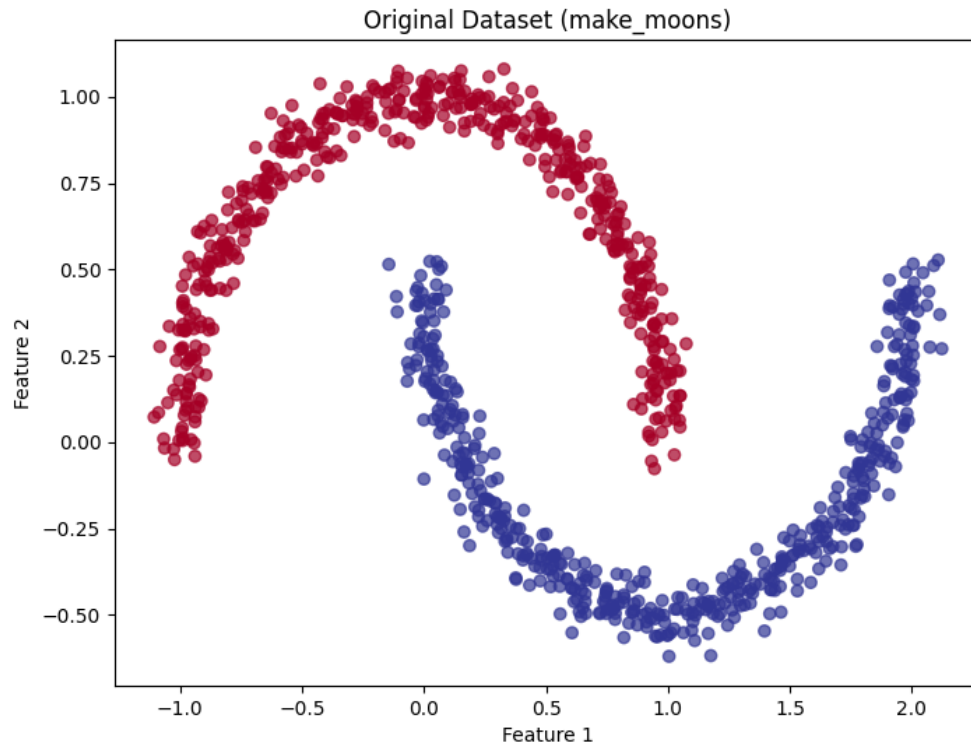
-


```
Epoch [10/100], Loss: 0.5707, Accuracy: 0.8213
Epoch [20/100], Loss: 0.3748, Accuracy: 0.8237
Epoch [30/100], Loss: 0.2742, Accuracy: 0.8500
Epoch [40/100], Loss: 0.2170, Accuracy: 0.8938
Epoch [50/100], Loss: 0.2070, Accuracy: 0.8938
Epoch [60/100], Loss: 0.2031, Accuracy: 0.8962
Epoch [70/100], Loss: 0.2016, Accuracy: 0.9025
Epoch [80/100], Loss: 0.2010, Accuracy: 0.9000
Epoch [90/100], Loss: 0.2008, Accuracy: 0.9012
Epoch [100/100], Loss: 0.2006, Accuracy: 0.9025
```

- 由上述实验可知，noise过大或者过小都会导致损失率变大，准确率降低，在当前模型下 noise=0.1时损失率最低，准确率最高

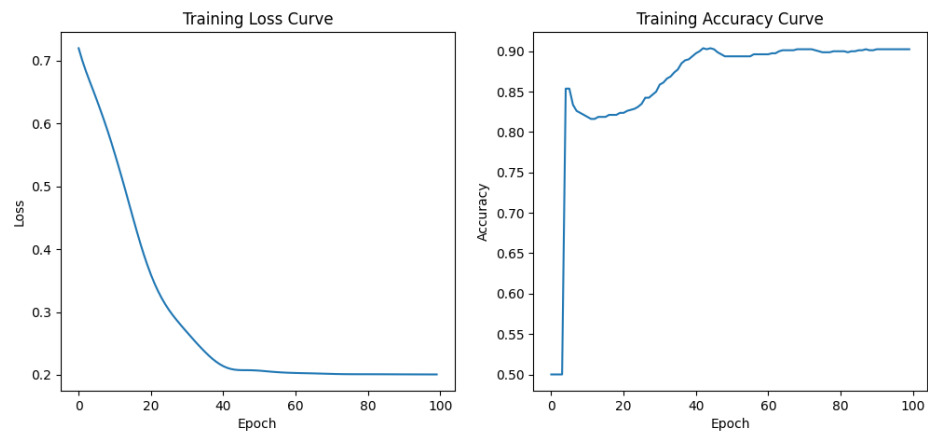
- 数据集可视化

-



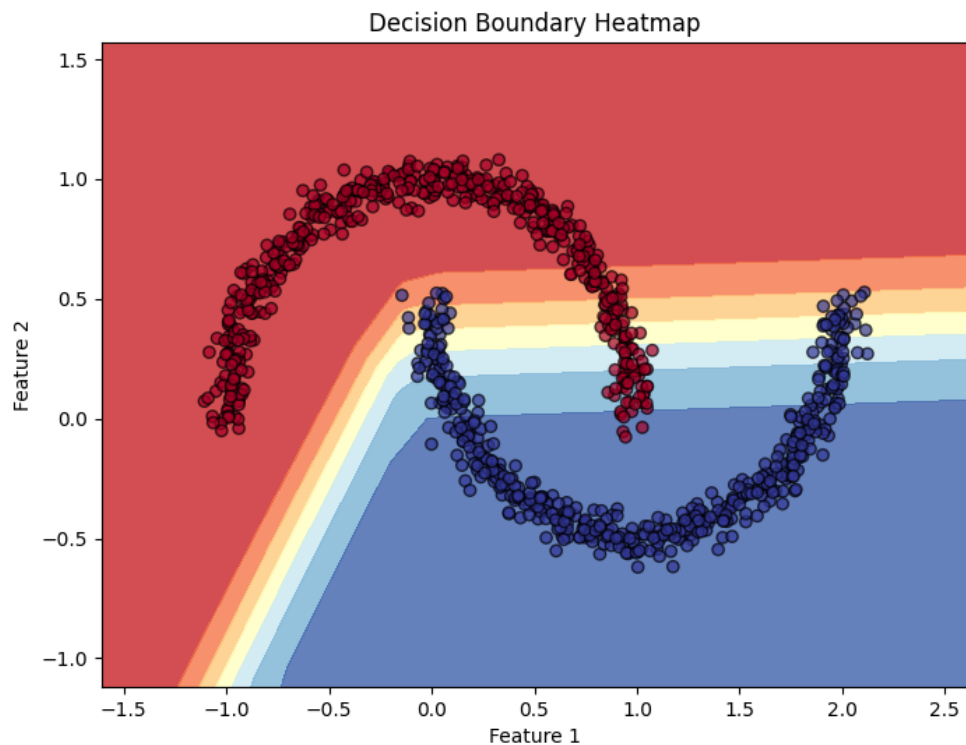
- **数据集的数学表示：**一个由两个类别组成的二维点集 $\{(x_1, x_2)_i\}$ 。其中每个类别的点，都采样自一个经过平移和旋转的半圆形曲线，并叠加了一个均值为零、方差由noise 参数控制的二维高斯噪声。这个数据集的核心特性是，两个类别的决策边界是一条非线性的曲线，无法用单一的直线进行有效分割。
- 在cpp中我是使用Eigen矩阵实现的数据集，首先把pytorch_impl.py使用的数据集导出为.csv文件，在cpp文件中读取该文件并且加载到Eigen矩阵中
- 损失下降曲线和分类准确率曲线

○



- 每个点推理结果可视化 (热力图)

○



- 在我的代码中，我是通过在神经网络的最后一层，使用一个特殊的激活函数来实现这一点的，这里由于是二分类问题所以使用的是sigmoid函数，如果是多分类问题则使用的是softmax函数进行归一化处理
- PyTorch vs. NumPy vs. C++(Eigen): 性能对比与深度剖析

在完成了三个版本的神经网络实现后，我们可以从**开发效率**和**运行效率**两个维度进行对比。虽然我们是在 CPU 上运行，但其性能差异的底层原理同样适用于 GPU。

一、性能排序 (从快到慢)

在同等硬件 (CPU) 和相同模型配置下，三个版本的**理论运行速度** (主要指训练/推理的核心计算部分) 排序如下：

1. C++ (Eigen): 最快
2. PyTorch: 很快 (与 C++ 版本差距不大，有时甚至可能因底层优化更好而反超)

3. NumPy: 最慢 (与前两者有数量级的差距)

而开发效率的排序则恰恰相反: `PyTorch > NumPy > C++`。

二、性能差异的深度原因分析

1. C++ (Eigen): 为什么最快? —— “原生战斗机”

◦ 底层原理:

- **直接编译为机器码:** C++ 是一种编译型语言。你的代码 (`main.cpp`) 通过编译器 (`g++`) 被直接翻译成了 CPU 可以**直接理解和执行的、最底层的机器指令**。这个过程中没有任何中间解释层, 几乎没有额外的性能开销。
- **高度优化的数学库 (Eigen):** 我们使用的 Eigen 库, 其本身就是用极其高效的 C++ 模板元编程技术写成的。它在编译时就能知道你的矩阵大小和类型, 从而生成高度优化的、专门针对你硬件的数学运算代码。它会充分利用 CPU 的**向量化指令集 (SIMD, 如 SSE/AVX)**, 让 CPU 在一个指令周期内, 同时对多个数据进行相同的运算 (比如同时完成4个浮点数的加法), 极大地提升了计算效率。

◦ 硬件利用情况:

- **CPU 亲和性最高:** 能够最大化地榨干 CPU 的计算能力, 包括多核心 (需要手动或用 OpenMP 等库开启多线程) 和 SIMD 指令集。
- **内存控制最精细:** C++ 允许你手动管理内存, 避免不必要的内存分配和拷贝, 对于性能极其敏感的应用来说至关重要。

- **结论:** C++ 版本就像一架为速度而生的原生战斗机, 每一个零件都为性能而设计, 没有任何多余的负重。

2. PyTorch: 为什么很快? —— “披着 Python 外衣的 C++ 战斗机”

◦ 底层原理:

- **混合编程模型:** PyTorch 是一个典型的“前店后厂”模型。
 - **前端 (Python):** 你用来写代码的 `torch.nn.Linear`, `loss.backward()` 等接口, 都是用 Python 写的。这使得它非常**灵活、易用、开发效率极高**。
 - **后端 (C++/CUDA):** 然而, 所有真正消耗计算资源的操作 (如大规模的矩阵乘法、卷积、梯度计算等), **都不是由 Python 执行的**。当你调用一个 PyTorch 函数时, Python 解释器会立刻把这个请求, 连同数据 (Tensor), 一起打包“扔”给它背后那个**用 C++ 和 CUDA 写成的高性能计算核心**。
- **优化的计算图:** PyTorch 会将你的模型操作构建成一个计算图。对于推理, 它甚至可以通过 TorchScript 等技术将整个图进行静态化和优化, 进一步减少 Python 解释器的开销。

◦ 硬件利用情况:

- **极致的硬件优化:** PyTorch 的后端核心, 是由全世界最顶尖的性能工程师编写的。它不仅会像 Eigen 一样利用 CPU 的 SIMD 指令, 还会链接到 Intel MKL、cuDNN 等厂商专门提供的、针对其硬件“量身定制”的、闭源的数学库。这些库的优化程度, 甚至超过了通用的开源库。
- **无缝的 GPU 加速:** PyTorch 的设计核心就是 GPU 计算, 能够极其高效地利用 CUDA 进行大规模并行计算。

- **结论：**PyTorch 的速度来自于它把 **Python 的灵活性和 C++/CUDA 的极致性能完美地结合了起来**。你只是在用 Python “发号施令”，而真正干重活的，是背后那个身经百战的 C++ 雇佣兵团。它和 C++ 版本的速度差距，主要来自于 Python 与 C++ 之间“发号施令”的这点通信开销。

3. NumPy：为什么最慢？——“能干的 Python 员工，但调度开销大”

- **底层原理：**

- **C 语言核心，但被 Python 逐行调用：**和 PyTorch 类似，NumPy 的核心数学运算（如 `np.dot`）也是用 C 语言写的，速度很快。**那为什么整个流程下来，它比 PyTorch 慢得多呢？**
- **问题的关键在于“反向传播”：**在我们的 NumPy 实现中，整个反向传播的过程是由**多行 Python 代码组成的一个序列**。

```
# 每一行都是一次独立的 Python->C 的调用
dz3 = a3 - y_train
dw3 = 1/m * np.dot(a2.T, dz3)
db3 = ...
dz2 = np.dot(dz3, w3.T) * relu_derivative(a1)
# ... 更多行
```

- **Python 全局解释器锁 (GIL) 和调用开销：**每一行 NumPy 操作，都涉及一次从 Python 解释器到 C 底层代码的**调用开销 (Overhead)**。Python 解释器需要解析这行代码，准备数据，调用 C 函数，然后再把结果拿回来。当你的反向传播链条很长时，这种“**Python 逐行发指令 -> C 执行 -> Python 再发下一条指令**”的模式，其累计的调度开销会变得非常巨大。
 - **对比 PyTorch：**PyTorch 的 `loss.backward()`，是你只**下达了一次命令**。然后，整个反向传播的复杂计算链条，**完全在 C++ 后端内部一气呵成**，完全没有 Python 的中间介入。
- **硬件利用情况：**
 - NumPy 同样能很好地利用底层的数学库（如 BLAS, LAPACK）和 CPU 的 SIMD 指令。
 - 但它天生不具备 GPU 加速能力，并且受制于 Python 的调度模式，无法像 PyTorch 或 C++ 那样形成一个高效的、端到端的计算流。
 - **结论：**NumPy 慢，不是慢在单次数学运算上，而是慢在了**由 Python 主导的、碎片化的、高调度开销的执行模式上**。它就像一个能力很强的员工，但他的经理（Python）必须一条一条地给他下指令，无法让他独立地、连续地完成一整套复杂的工作。