

This can be run [run on Google Colab using this link](#)

Homework 2.1: Training a Deep Network

```
In [2]: # %shell wget 'https://raw.githubusercontent.com/CS7150/CS7150-Homework-2/refs/heads/main/hw2utils.py'
```

Overview

We will start by exploring the optimization aspects of deep network training. Throughout this journey, you will gain insights into:

Part 1:

- The fundamentals of simple gradient descent.
- The concept of weight decay.
- A deep understanding of PyTorch autograd and PyTorch optimizers.

Part 2:

- Analyzing raw gradients, means, and RMS (Root Mean Square).
- Delving into exponential moving averages.
- Exploring the workings of the ADAM optimization algorithm.

Part 3:

- Strategies for optimizing neural network parameters.
- Selecting appropriate nonlinearities, architectures, and layers to tackle the vanishing gradient problem.
- Leveraging techniques like regularization, parameterization, and specific layer choices to enhance generalization.
- Unpacking the roles and impacts of ADAM, ReLU activation, weight decay, network depth, network width, residual architectures, and batch normalization in deep learning.

Note

- **You do not need to tune hyper parameters in the regular tasks.**

- **You do not install any additional packages inside the Colab environment.**

- If you collaborate or get assistance from classmates, online resources, AI, or other sources, then you must then you must explicitly write down the sources that you used to credit them.
- Attend office hours and make post on Piazza if you have any questions.
- You have sufficient time to work on this assignment. Please refrain from asking for extensions.

Setup

In [3]: # Import Libraries

```
import torch
import numpy as np
import torch.nn as nn
import torch.optim as optim
import matplotlib
from torch.nn import Sequential, Module
from matplotlib import pyplot as plt
from scipy.stats import norm
from hw2utils import LossFunctionWithPlot, ConstantVectorNetwork
```

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\903424143.py:10: UserWarning: A NumPy version >=1.23.5 and <2.5.0 is required for this version of SciPy (detected version 1.23.1)
from scipy.stats import norm

Part 1: Simple Gradient Descent

Gradient descent is an iterative optimization algorithm for finding the minimum of a function. It works by starting at a point and then moving in the direction of the steepest descent until it reaches a minimum. The steepest descent is the direction in which the function is decreasing most rapidly.

To train a model using gradient descent, we start with a set of initial parameters. These parameters are the values of the variables in the model. We then repeatedly apply gradient descent to update the parameters. Each update moves the parameters in the direction of the steepest descent. This continues until the parameters converge to a minimum of the function.

The choice of the learning rate is important for gradient descent. The learning rate is the size of the steps that are taken in the direction of the steepest descent. If the learning rate is too small, the algorithm will converge slowly. If the learning rate is too large, the algorithm may diverge and never converge.

Let's say we have a function $\mathcal{L}(x)$ that we want to minimize. The gradient of $\mathcal{L}(x)$ is a vector that points in the direction of the steepest descent of $f(x)$. The gradient can be calculated using the following equation:

$$\nabla \mathcal{L}(x) = \begin{bmatrix} \frac{\partial \mathcal{L}(x)}{\partial x_1} \\ \frac{\partial \mathcal{L}(x)}{\partial x_2} \\ \vdots \\ \frac{\partial \mathcal{L}(x)}{\partial x_n} \end{bmatrix}$$

The gradient descent algorithm can be used to minimize $f(x)$ by repeatedly taking steps in the direction of the gradient. The update rule for gradient descent is given by the following equation:

$$x_{\text{new}} = x_{\text{old}} - \alpha \cdot \nabla \mathcal{L}(x_{\text{old}}) \quad (1)$$

where x_{old} is the current value of x , x_{new} is the new value of x , α is the learning rate, and $\nabla \mathcal{L}(x_{\text{old}})$ is the gradient of $\mathcal{L}(x)$ evaluated at x_{old} .

Simple Implementation of Gradient Descent on a quadratic loss

Below we demonstrate gradient descent optimization in action.

We iteratively update the `x` to try to minimize a quadratic function `L` defined by the `LossFunctionWithPlot` class. The trajectory of updates and corresponding losses are stored and plotted to visualize the optimization progress.

Provide an implementation of simple gradient descent below. In 21 steps you can make the loss go down to about 3 or better, and drive `x` somewhat towards the center of the minimum of the loss function.

- Use `loss.backward()` (read <https://pytorch.org/docs/stable/generated/torch.Tensor.backward.html>)
- Use `x.grad` to get the gradient.
- Update `x` in-place using `x -= something`, and know why `torch.no_grad()` is needed.
- Understand why gradients need to be zeroed: <https://stackoverflow.com/questions/48001598/why-do-we-need-to-call-zero-grad-in-pytorch>

Task 1.1 - Implement simple gradient descent (1 point)

```
In [4]: # Do not change the starting x.
x = torch.tensor([-1.5, 1.2])
x.requires_grad = True
L = LossFunctionWithPlot()

# Start by using a Learning rate of 0.1.
learning_rate = 0.05

for iter in range(50):
    loss = L(x)
    if iter % 7 == 0:
        print(f"Loss at step {iter} is {loss.item():.3f}")
    loss.backward()
    with torch.no_grad():
        x -= learning_rate * x.grad
```

```

#####
# Implement Simple Gradient Descent and update variable 'x'
# Read Documentation to compute a gradient of a parameter -
# https://pytorch.org/docs/stable/autograd.html
#####
x -= learning_rate * x.grad

# Equivalent to the above code without using torch.no_grad()
# x = (x - learning_rate * x.grad).detach()
# x.requires_grad = True

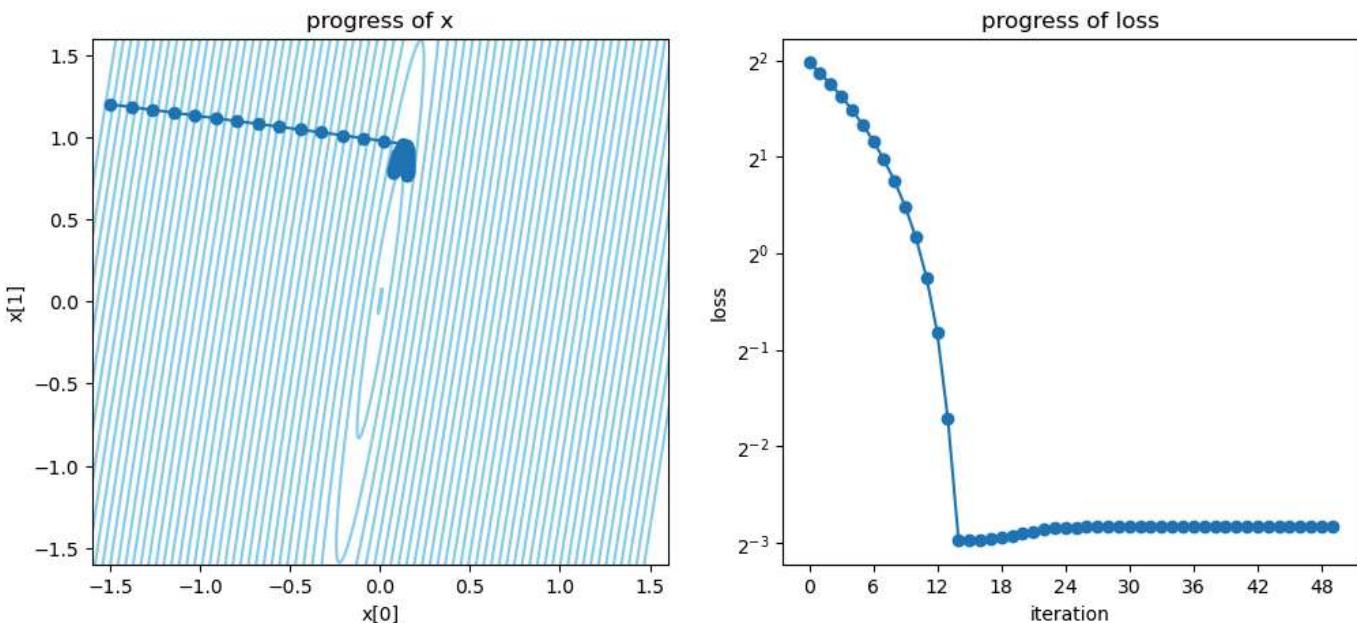
#####
# END OF YOUR CODE
#####
x.grad = None
L.plot_history()

```

Loss at step 0 is 3.937
 Loss at step 7 is 1.962
 Loss at step 14 is 0.127
 Loss at step 21 is 0.136
 Loss at step 28 is 0.140
 Loss at step 35 is 0.141
 Loss at step 42 is 0.141
 Loss at step 49 is 0.141

d:\uni\courses\S2025\CS_7150_Deep_Learning\HW\CS7150-Homework-2\hw2utils.py:72: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



One reason it is hard for gradient descent to move towards the bottom of the bowl is that there is no single ideal learning rate for all dimensions. In the problem in 1.1, notice that:

- the loss changes direction very quickly (the curvature is high) in the horizontal direction of `x[0]`
- the loss is very slow-changing (the curvature is lower) in the vertical direction of `x[1]`.

The "sawtooth" loss curves and the "zig-zag" paths are symptoms of an optimizer that is taking steps that are too large: the path could be repeatedly jumping over a valley in the loss surface and ending up at another point of high, or even higher loss. A lower learning rate can help, but it can lead to another problem. (What problem? Try it.)

Another fancy idea is to mix learning rates, with different learning rates for each parameter. Although you can still make a learning rate too high or too low.

Task 1.2 - Explore Simple Gradient Descent using various learning rates (_ points)

Now copy your code from 1.1 below here, but **experiment with learning rates**, including **unequal learning rates** for `x[0]` and `x[1]` by setting to a tensor with two values. Try to find a pair of learning rates that move x to the bottom of the bowl and stays there with near-zero loss.

In [5]:

```
#####
# Copy your solution from Task 1.1 here and attempt to discover a pair of
# Learning rates that guide the optimization process to place 'x' at the lowest
# point of the bowl and maintain it there with a nearly zero loss.
#####
# Do not change the starting x.
x = torch.tensor([-1.5, 1.2])
x.requires_grad = True
L = LossFunctionWithPlot()

# Start by using a learning rate of 0.1.
learning_rate = torch.tensor([0.1, 0.55])

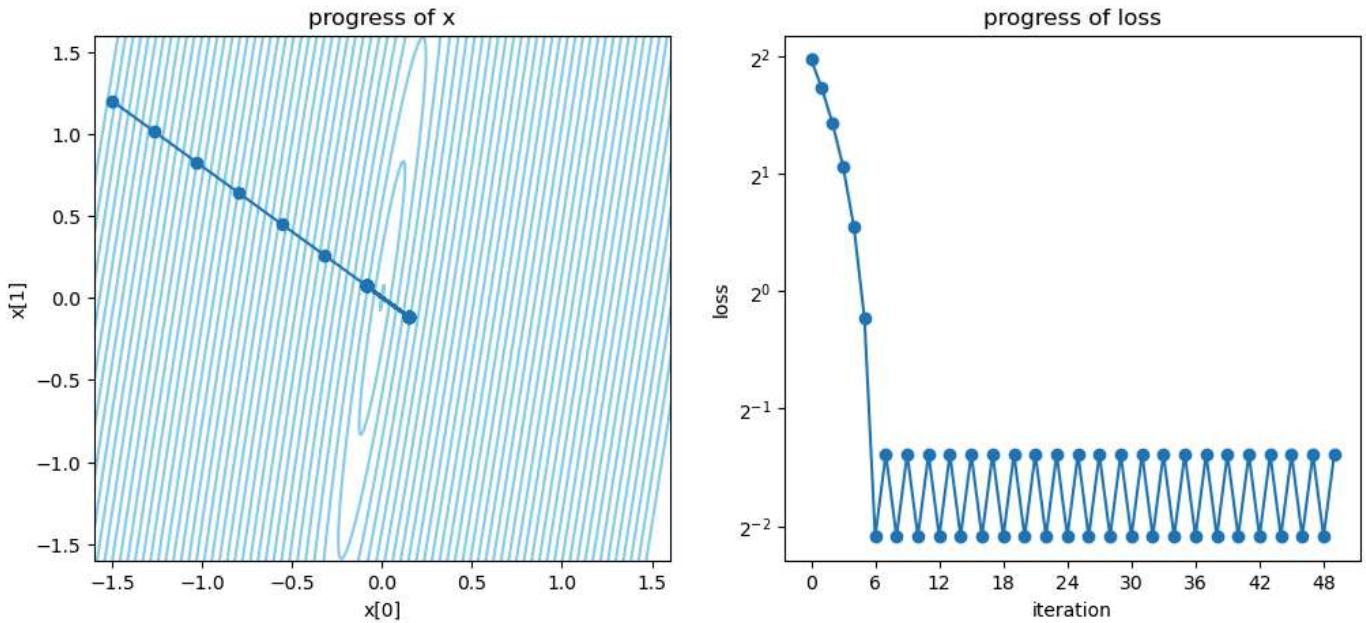
for iter in range(50):
    loss = L(x)
    if iter % 7 == 0:
        print(f"Loss at step {iter} is {loss.item():.3f}")
    loss.backward()
    with torch.no_grad():
        x -= learning_rate * x.grad
    x.grad = None
L.plot_history()

#####
# END OF YOUR CODE
#####

```

```
Loss at step 0 is 3.937
Loss at step 7 is 0.382
Loss at step 14 is 0.235
Loss at step 21 is 0.382
Loss at step 28 is 0.235
Loss at step 35 is 0.382
Loss at step 42 is 0.235
Loss at step 49 is 0.382
```

```
d:\uni\courses\S2025\CS_7150_Deep_Learning\HW\CS7150-Homework-2\hw2utils.py:72: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()
```



Inline Question 1.A (1 point): What bad things tend to happen when the learning rate is too high?

Answer:

When the learning rate is too high [It tends to over-adjust the parameters thus cause divergence of]

Inline Question 1.B (1 point): What bad things tend to happen when the learning rate is too low?

Answer:

When the learning rate is too low [It tends to be stuck into a local minima.]

Inline Question 1.C (1 point): Should a higher or lower learning rate be used on a dimension with higher curvature (with sharper changes)?

Answer:

When the dimension has higher curvature, the learning rate should generally be [lower to reduce overfitting]

II) Weight Decay, aka L2 regularization

Weight decay, also known as L2 regularization, adds the following term to the total loss:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{objective}}(\theta) + \frac{\lambda}{2\nu} |\theta|^2 \quad (2)$$

where:

- $\mathcal{L}(\theta)$ - is the regularized objective function.
- $\mathcal{L}_{\text{objective}}(\theta)$ - is the original objective function (without regularization).
- θ is the vector of model parameters.
- λ is the regularization parameter.

1.D) Inline Question (1 point): What is the gradient $\nabla \mathcal{L}(\theta)$ in terms of $\nabla \mathcal{L}_{\text{objective}}(\theta)$ and θ ?

$$\text{Answer: } \nabla \mathcal{L}(\theta) = \nabla \mathcal{L}_{\text{objective}}(\theta) + \frac{\lambda}{\nu} |\theta|$$

1.E) Inline Question (1 point): Suppose 'v' is the learning rate. Then what should be the update rule for Θ ?

$$\text{Answer: } \theta_t = \theta_{t-1} - v \nabla \mathcal{L}(\theta) = (1 - \lambda) \theta_{t-1} - v \nabla \mathcal{L}_{\text{objective}}(\theta_{t-1})$$

Task 1.3 - Implement a `SimpleGradientDescent` optimizer class in pytorch (2 point)

Pytorch encapsulates optimization algorithms into optimizer classes that hold on to a list of parameters being optimized, and that update the parameters in-place based on gradients using a `step()` method.

Here we see how that pattern works.

Implement the `SimpleGradientDescent` class as a pytorch-style optimizer by completing the code below.

Incorporate weight decay by adding the term you computed above, where λ is the `weight_decay` hyperparameter.

```
In [6]: class SimpleGradientDescent:
    def __init__(self, parameters, lr=0.1, weight_decay=0.0):
        self.lr = lr
        self.weight_decay = weight_decay
        self.parameters = []
        for x in parameters:
            self.parameters.append(x)

    def step(self):
        with torch.no_grad():
            for x in self.parameters:
                #####
                # Implement Simple Gradient Descent with the inclusion of
                # weight_decay, and then update the results in variable x.
                #####
                x -= self.lr * x.grad + self.weight_decay * x

                #####
```

```

#                                     END OF YOUR CODE                               #
#####
#####
```

```

def zero_grad(self):
    for x in self.parameters():
        x.grad = None

x = torch.tensor([-1.5, 1.2])
x.requires_grad = (
    True # This tells PyTorch that the x variable will be used to calculate gradients.
)

L = LossFunctionWithPlot()

learning_rate = 0.1
optimizer = SimpleGradientDescent([x], lr=learning_rate, weight_decay=1e-3)

for iter in range(50):
    loss = L(x)
    if iter % 7 == 0:
        print(f"Loss at step {iter} is {loss.item():.3f}")
    loss.backward()
    with torch.no_grad():
        optimizer.step()
    optimizer.zero_grad()

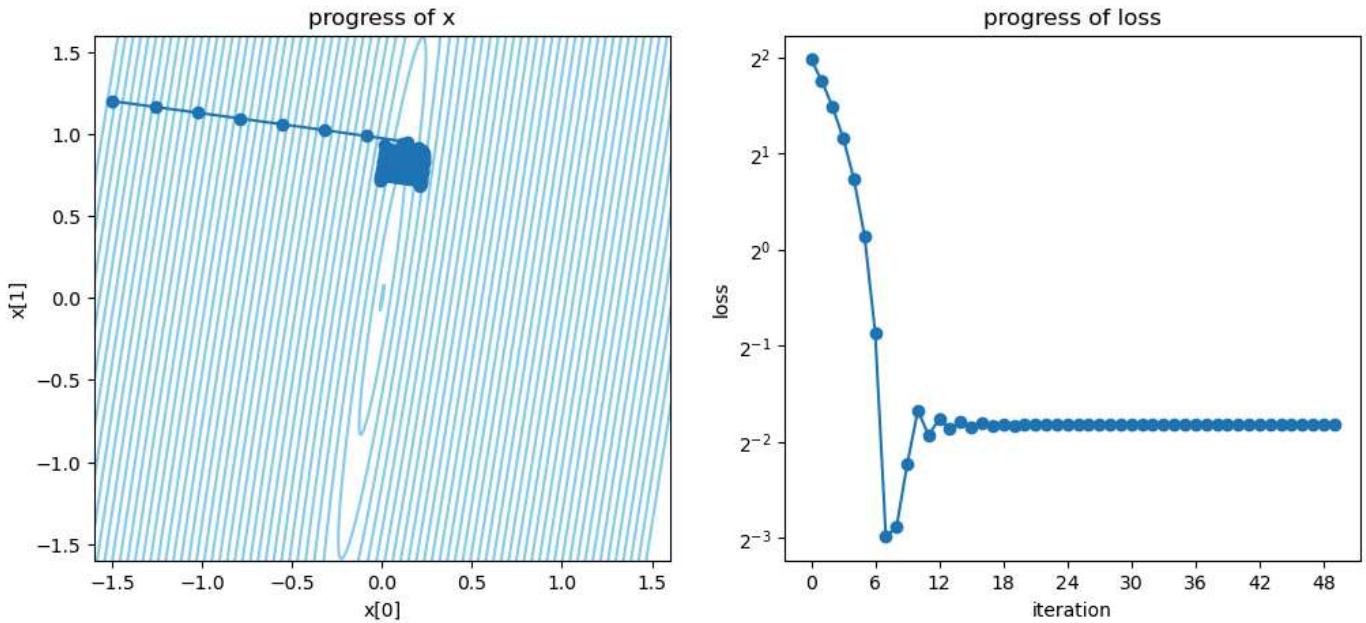
L.plot_history()
```

```

Loss at step 0 is 3.937
Loss at step 7 is 0.126
Loss at step 14 is 0.288
Loss at step 21 is 0.282
Loss at step 28 is 0.283
Loss at step 35 is 0.283
Loss at step 42 is 0.283
Loss at step 49 is 0.283
```

```
d:\uni\courses\S2025\CS_7150_Deep_Learning\HW\CS7150-Homework-2\hw2utils.py:72: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
```

```
fig.show()
```



Part 2: The ADAM optimizer

I) Analyzing Raw Gradients, Means, and RMS (Root Mean Square).

Measure mean of the gradient on each dimension

First, let us plot and examine the mean of the gradient.

Using the code below as a starting point, calculate and fill in the following mean gradient over the 50 iterations in the specific optimization below:

$$\text{mean_grads}[i] = \frac{1}{N} \sum_{t=1}^N \frac{\partial \mathcal{L}}{\partial x_i}(x^{(t)})$$

2.A) Inline Questions (2 points):

1. Mean gradient component with respect to $x[0] = \boxed{-0.2982}$
2. Mean gradient component with respect to $x[1] = \boxed{0.0940}$

Measure root-mean-square of the gradient on each dimension

Second, let's plot the root mean square (RMS) of the gradient.

Using the code below as a starting point, calculate and fill in the following root-mean-square gradient over the 50 iterations in the specific optimization below:

$$\text{rms_grads}[i] = \sqrt{\frac{1}{N} \sum_{t=1}^N \left(\frac{\partial \mathcal{L}}{\partial x_i}(x^{(t)}) \right)^2}$$

2.B) Inline Questions (2 points):

1. RMS gradient component with respect to $x[0] = \boxed{2.1092}$

2. RMS gradient component with respect to $x[1] = \boxed{0.3078}$

Understanding the challenges faced by simple gradient descent.

The optimization problem in Part I causes gradient descent to run into a few different problems:

- after initially descending the loss jumps back up.
- after making initial quick progress, x gets stuck instead of heading towards the middle.

Task 2.1 Insights from Raw Gradient Plots and Statistical Metrics (2 points)

To understand the problems, run the code below to see plots of the raw gradient, and then compute the Mean and root-mean-square (RMS) of each component of the gradient over all the iterations.

```
In [7]: x = torch.tensor([-1.5, 1.2])
x.requires_grad = (
    True # This tells PyTorch that the x variable will be used to calculate gradients.
)

L = LossFunctionWithPlot()

learning_rate = 0.1
optimizer = torch.optim.SGD([x], lr=learning_rate)
grads = []

for iter in range(50):
    loss = L(x)
    loss.backward()
    with torch.no_grad():
        optimizer.step()
    grads.append(x.grad.clone())
    optimizer.zero_grad()

grads = torch.stack(grads)

#####
# compute mean derivatives and root mean square derivatives for x[0] and x[1]
#####

mean_grads = grads.mean(
    dim=0
) # should be computed as a pair of means, one for each dimension
rms_grads = (
    grads.pow(2).mean(dim=0).sqrt()
) # computes RMS for each dimension separately
# should be computed as a pair of root-mean-squares, one for each dimension

#####
```

```

#                                                 END OF YOUR CODE #
#####
print("Mean grads for two dimensions separately:", mean_grads)
print("Root mean square grads for two dimensions separately:", rms_grads)

L.plot_history(grads=grads, mean_grads=mean_grads, rms_grads=rms_grads)

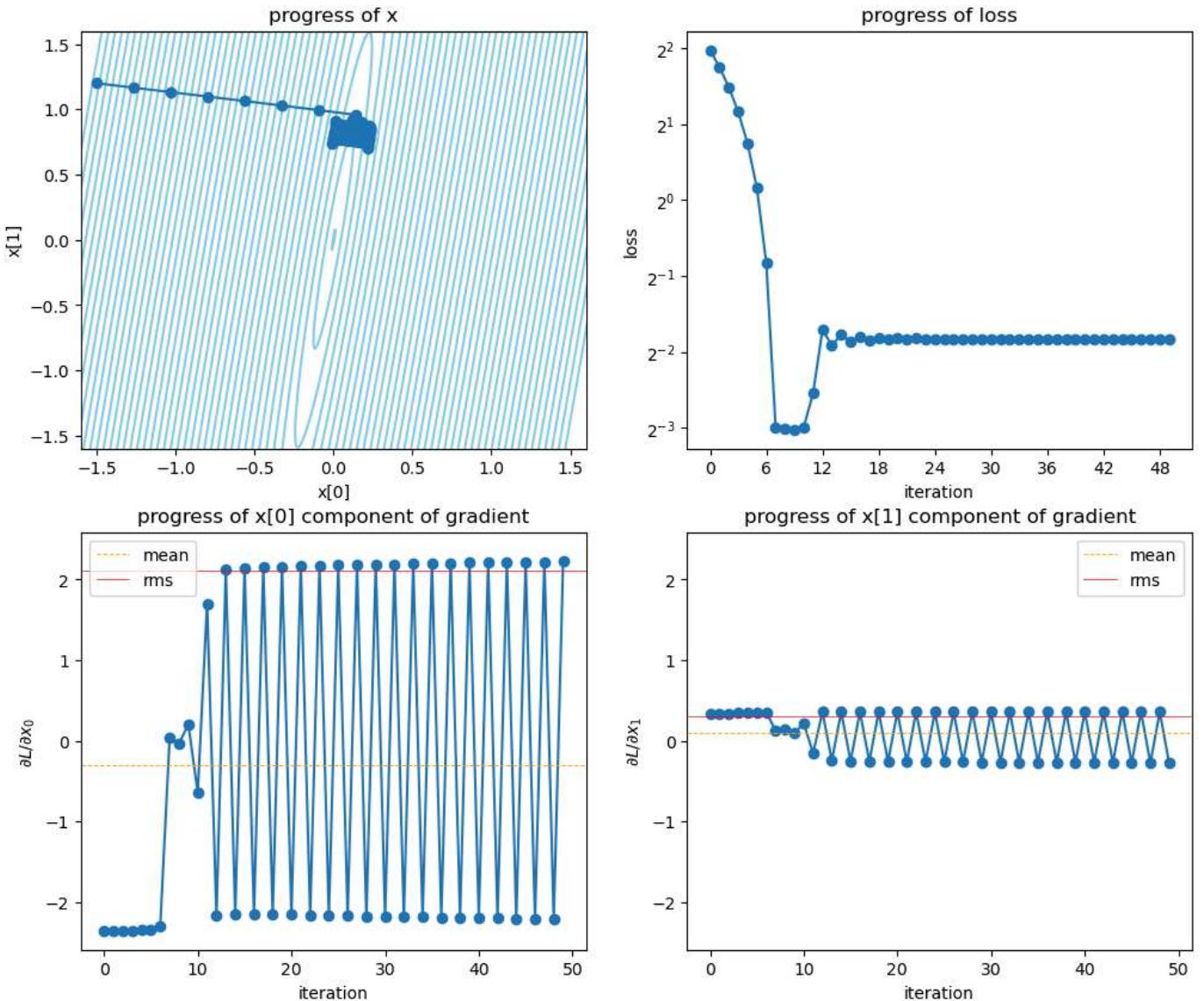
```

Mean grads for two dimensions separately: tensor([-0.2982, 0.0940])

Root mean square grads for two dimensions separately: tensor([2.1092, 0.3078])

d:\uni\courses\S2025\CS_7150_Deep_Learning\HW\CS7150-Homework-2\hw2utils.py:72: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Facts about mean and RMS

To understand the role of the mean and RMS, answer the following questions.

2.C) Inline Question (4 points):

- When they differ, which is guaranteed to be smaller: mean or RMS? =

2. Which is a better representation of the "average size", mean or RMS? = \text{RMS}
3. When the derivative is consistently positive, will the mean tend to be larger or smaller? = \text{}
4. When the derivative sign changes frequently, will the mean tend to be larger or smaller? = \text{}

Anticipating problems in optimization using Mean and RMS

Based on the problems we are seeing, we want the optimizer to slow down when the mean is much smaller than the RMS, and speed up when the mean and RMS are about the same size. That will:

1. Make sure the updates to be **small enough** when the gradient starts becoming bumpy, instead of oscillating.
2. Make sure the updates to be **large enough** when gradient is smooth but happens to be small.

The mean and RMS of the gradient can be used to deal with both these problems. In the next section, we will see how they are directly applied in the ADAM optimizer.

The idea of the ADAM optimizer

The ADAM optimizer automatically chooses a different learning rate for each parameter by using a heuristic that shrinks the update size in regions where the gradient is changing more quickly, while normalizing the update size so that it is a consistent size even in regions where the gradient is very small. Update magnitudes are calculated per-parameter, so ADAM can help deal with parameters that behave very differently from each other.

The idea behind ADAM is to choose an update that is proportional to a fraction between a weighted mean of the gradient and a weighted RMS of the gradient:

$$\Delta x = -\alpha \frac{\text{mean gradient}}{\text{rms gradient}} = -\alpha \frac{\sum_i w_i g_i}{\sqrt{\sum_i u_i g_i^2}}$$

In the definition above, the g_i are samples of the gradient from previous steps, and w_i and u_i are the weights to use for averaging.

Why one might divide the mean gradient by the RMS of the gradient?

To understand, answer the following question.

2.D) Inline Question (2 points):

Suppose there is a new problem which is scaled by some constant K so that all the new gradients are uniformly scaled larger. We want to understand whether ADAM speeds up or slows down when gradients are scaled up. Precisely: If gradients $\hat{g}_i = K g_i$ are scaled up with $K > 1$, how will the ADAM update $\Delta \hat{x}$ relate to the original problem's ADAM update Δx ? (e.g., which is larger?)

Answer.

$$\text{They should be equal. } \Delta\hat{x} = -\alpha \frac{\sum_i w_i K g_i}{\sqrt{\sum_i u_i (K g_i)^2}} = -\alpha \frac{K \cdot \sum_i w_i g_i}{K \cdot \sqrt{\sum_i u_i (g_i)^2}} = \Delta x$$

This property means that ADAM will not go too fast nor too slow just because the average size of the gradient is too large. The only thing that will cause ADAM to slow down is when the mean is much smaller than the RMS, which happens when the gradient frequently changes sign, for example, when the optimizer is oscillating around a minima.

Because the mean and RMS will change during the optimization, in practice the ADAM algorithm is based on using **exponential moving averages** which will adapt as optimization proceeds.

II) Exponential Moving Average (EMA) on Time Series Data

The exponential Moving Average (EMA) represents a moving average variant that assigns greater importance to the most recent data points within a time series. It achieves this by progressively diminishing the influence of older data points. Unlike simple moving averages, where all data points hold the same weight, EMA's differential weighting scheme enhances its sensitivity to recent data alterations, rendering it highly attuned to the latest fluctuations within the data.

EMA is a weighted average where, the weight of a sample of age $t - i$ is decayed exponentially by $\beta^{(t-i)}$, where $\beta < 1$ is the smoothing parameter. That is, using geometric series identities,

$$\text{EMA}_t = \frac{\beta^{t-1}x_1 + \beta^{t-2}x_2 + \dots + \beta x_{t-1} + x_t}{\beta^{t-1} + \beta^{t-2} + \dots + \beta + 1} = \frac{(1 - \beta) \cdot \sum_i \beta^{t-i} x_i}{1 - \beta^t}$$

As t gets large, the denominator becomes indistinguishable from one, and EMA can be estimated by computing just the numerator $\text{EMA}_t^* = (1 - \beta) \cdot \sum_i \beta^{t-i} x_i$. The numerator has the advantage that maintaining a running average only requires a single number be remembered: the most recent numerator EMA_{t-1}^* . The formula for calculating the numerator EMA_t^* is as follows:

$$\begin{aligned}\text{EMA}_0^* &= 0 \\ \text{EMA}_t^* &= \beta \cdot \text{EMA}_{t-1}^* + (1 - \beta) \cdot x_t\end{aligned}$$

When t is small, the numerator can be much smaller than one, so it must be included, so the full formula for the EMA is:

$$\text{EMA}_t = \frac{\text{EMA}_t^*}{1 - \beta^t}$$

Where,

- EMA_t - is the Exponential Moving Average at time t .
- EMA_t^* - is the Exponential Moving Average Numerator, which $\approx \text{EMA}_t$ when t is large.
- x_t - is the data point at time t that you want to include in the EMA calculation.
- EMA_{t-1}^* - is the EMA numerator calculated at the previous time step $t-1$.

- β is the smoothing factor

Task 2.2 Implement EMA (2 points)

Based on the definitions above, implement ema_update below, and produce the plot of the EMA of the synthetic time series data.

As you can see, unlike the ordinary mean, EMA adapts to changes in the data over time.

Also notice the difference between EMA* and EMA at the beginning of the series. Notice that EMA* is not unbiased: instead it has a clear bias towards zero.

```
In [8]: def ema_update(x_t, beta, t, ema_star_old):
    ema_star_t = 0.0
    ema = 0.0
    #####
    # 1) compute ema_star_t from ema_star_old, beta, and x_t.
    # 2) compute ema from ema_star_t, beta, and t
    #####
    ema_star_t = beta * ema_star_old + (1 - beta) * x_t
    ema = ema_star_t / (1 - beta**t)
    #####
    #
    # END OF YOUR CODE
    #
#####

return ema, ema_star_t

timestamps = 100
time_series = torch.cat([
    [
        torch.randn(timestamps) * 0.1 - 10.0,
        torch.randn(timestamps) + 0.5,
    ]
])
mean = time_series.mean()

beta = 0.9
ema_star = 0.0
history, history_star = [], []

for t, d in enumerate(time_series):

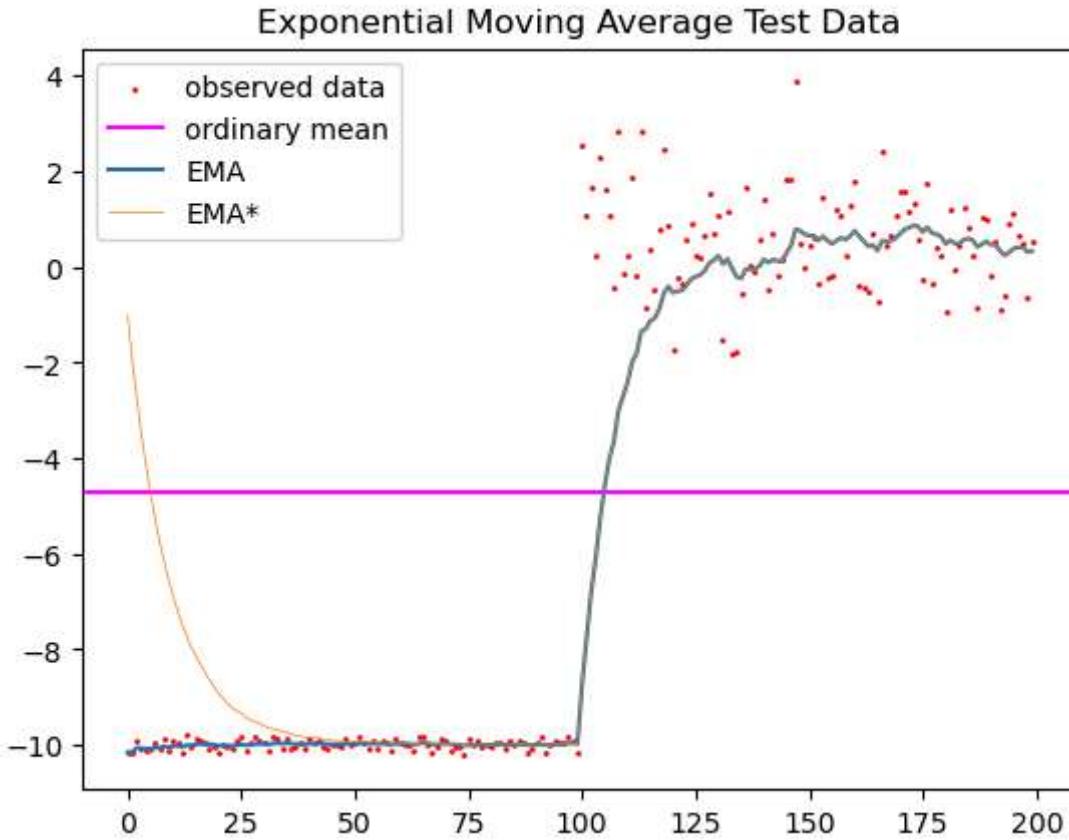
    # This is the ema update
    ema, ema_star = ema_update(d, beta, t + 1, ema_star)
    history.append(ema)
    history_star.append(ema_star)

plt.title("Exponential Moving Average Test Data")
plt.scatter(
    range(len(time_series)), time_series, s=1, color="red", label="observed data"
)
```

```

plt.axhline(mean, label="ordinary mean", color="magenta")
plt.plot(history, label="EMA")
plt.plot(history_star, label="EMA*", linewidth=0.5)
plt.legend()
plt.show()

```



III) ADAM Optimizer

ADAM (Adaptive Moment Estimation) is a popular optimization algorithm used for training machine learning and deep learning models. ADAM is known for its efficiency, robustness, and ability to handle a wide range of optimization problems.

Here are the key components and features of the ADAM optimizer:

- **Adaptive Learning Rates:** ADAM adapts the learning rates for each parameter during training. It maintains a separate learning rate for each parameter based on the historical gradients of that parameter. This adaptability helps the algorithm converge faster and handle sparse gradients effectively.
- **EMA Momentum:** ADAM incorporates the concept of momentum, similar to the SGD with momentum optimizer. It uses exponential moving averages of past gradients to help the optimization process. This momentum term smooths the optimization trajectory and accelerates convergence.

The update rule for ADAM is as follows:

$$\begin{aligned}
 m_t^* &= \beta_1 \cdot m_{t-1}^* + (1 - \beta_1) \cdot g_t && \text{(First Moment)} \\
 m_t &= \frac{m_t^*}{1 - \beta_1^t} && \text{(Bias correction)} \\
 v_t^* &= \beta_2 \cdot v_{t-1}^* + (1 - \beta_2) \cdot g_t^2 && \text{(Second Moment)} \\
 v_t &= \frac{v_t^*}{1 - \beta_2^t} && \text{(Bias correction)} \\
 x_{t+1} &= x_t - \alpha \cdot \frac{m_t}{\sqrt{v_t} + \epsilon} && \text{(Parameter Update)}
 \end{aligned}$$

Where:

- m_t and v_t - are the EMA estimates of the first and second moments of the gradients at time step t, respectively.
- β_1 and β_2 - are exponential decay rates for the first and second moments, typically close to 1
- g_t - is the gradient of the parameter θ at time step t, which is used to calculate these moving averages.
- α is the learning rate.
- ϵ is a small constant, just to avoid division-by-zero.

This can also be written as

$$\begin{aligned}
 m_t, m_t^* &= \text{ema_update}(g_t, \beta_1, t, m_{t-1}^*) && \text{(EMA Update)} \\
 v_t, v_t^* &= \text{ema_update}(g_t^2, \beta_2, t, v_{t-1}^*) && \text{(EMA Update)} \\
 x_{t+1} &= x_t - \alpha \cdot \frac{m_t}{\sqrt{v_t} + \epsilon} && \text{(Parameter Update)}
 \end{aligned}$$

As a final detail, we can incorporate weight decay by defining g_t to include the weight decay term as worked out in a previous exercise.

Task 2.3 - Implement the ADAM optimizer (_ points)

Now implement your own ADAMOptimizer, using the following code as a starting point.

- Within the loop, you can use `ema_update` to make the code simpler.
- Remember that there are two smoothing factors, `beta_1` and `beta_2`.
- Remember to incorporate hyperparameters `learning_rate`, `epsilon`, and `weight_decay`

You can test your code by swapping between `torch.optim.Adam` and your own `ADAMOptimizer`. On this test, they should behave the same.

```
In [9]: class ADAMOptimizer:
    def __init__(self, parameters, lr=0.1, betas=[0.9, 0.999], eps=1e-8, weight_decay=0.0):
        self.lr = lr
        self.beta_1, self.beta_2 = betas
        self.eps = eps
```

```

    self.weight_decay = weight_decay

    # t is a running timestamp.
    self.t = 0

    self.parameters = []
    self.m_star = []
    self.v_star = []

    for x in parameters:

        self.parameters.append(x)

        self.m_star.append(torch.zeros_like(x))
        self.v_star.append(torch.zeros_like(x))

    def step(self):
        self.t += 1

        with torch.no_grad():
            for x, m_star, v_star in zip(self.parameters, self.m_star, self.v_star):
                g = x.grad + self.weight_decay * x

                # Decouple the weight decay from the gradient (Equivalent to AdamW)
                # g = x.grad
                #####
                # Implement ADAM optimization
                #

                # There are three steps:

                # (1) m and m_star are updated to incorporate g with smoothing
                # beta_1 using ema_update(...).
                m, m_star[...] = ema_update(g, self.beta_1, self.t, m_star)

                # (2) v and v_star are updated to incorporate g^2 with smoothing
                # beta_2 using ema_update(...).
                v, v_star[...] = ema_update(g.pow(2), self.beta_2, self.t, v_star)

                # (3) x is updated according to ratio between the mean and RMS
                # gradient estimates, using lr and epsilon.

                x -= self.lr * m / (v.sqrt() + self.eps)

                # Decouple the weight decay from the gradient (Equivalent to AdamW)
                # x -= self.lr * (m / (v.sqrt() + self.eps) + self.weight_decay * x)

                # Remember to advance the timestep t before each step.
                # Remember to update m_star and v_star in-place.

                # There is a difference between saying m_star = something and
                # m_star[...] = something.
                #####
                #####

```

```

#                                     END OF YOUR CODE #
#####
#####
```

def zero_grad(self):

for x **in** self.parameters:

x.grad = **None**

TEST CODE BELOW

'x' variable = is the current estimate which will be updated iteratively during the optimization.

'L' function = the quadratic loss function we will use. This one can also plot its inputs and outputs.

'learning_rate' = is the step size that is used to update the solution

x = torch.tensor([-1.5, 1.2])

x.requires_grad = (
 True # This tells PyTorch that the x variable will be used to calculate gradients.
)

L = LossFunctionWithPlot()

learning_rate = 0.1

weight_decay = 0.1

If you switch this for torch.optim.Adam, it should behave exactly the same.

optimizer = ADAMOptimizer([x], lr=learning_rate, weight_decay=weight_decay)
optimizer = torch.optim.Adam([x], lr=Learning_rate, weight_decay=weight_decay)

optimizer = torch.optim.AdamW([x], lr=Learning_rate, weight_decay=weight_decay)

for iter **in** range(50):

loss = L(x)

if iter % 7 == 0:
 print(f"Loss at step {iter} is {loss.item():.3f}")

loss.backward()

```

    with torch.no_grad():

        optimizer.step()

        optimizer.zero_grad()

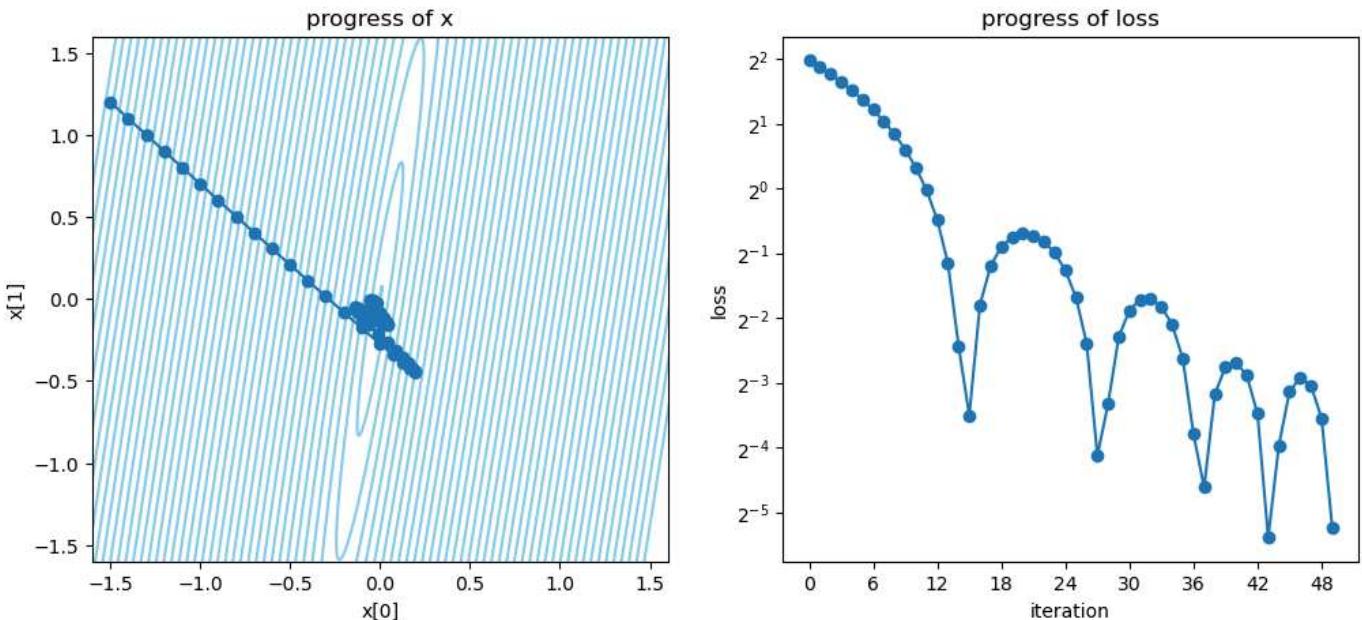
L.plot_history()

```

Loss at step 0 is 3.937
 Loss at step 7 is 2.054
 Loss at step 14 is 0.186
 Loss at step 21 is 0.605
 Loss at step 28 is 0.100
 Loss at step 35 is 0.162
 Loss at step 42 is 0.090
 Loss at step 49 is 0.027

d:\uni\courses\S2025\CS_7150_Deep_Learning\HW\CS7150-Homework-2\hw2utils.py:72: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Part 3: Training a Neural Network

Tiny Classification Problem

The following code loads raw data for a tiny classification problem.

The training data has 10000 samples, each a vector of 36 numbers along with a corresponding set of 10000 labels, assigning 0 or 1 to each sample. The test data has 2000 samples and labels that are disjoint from the training data.

```
In [10]: train_data, train_labels, test_data, test_labels = [
    torch.tensor(m[k]).float()
```

```

        for m in [np.load("tiny-classification.npz")]
            for k in "train_data train_labels val_data val_labels".split()
        ]

    print(
        f"The training data has {train_data.size(0)} samples, each a vector of
        {train_data.size(1)} numbers along with"
    )
    print(
        f"a corresponding set of {train_labels.size(0)} labels, assigning {train_labels.min()} or
        {train_labels.max()} to each sample."
    )

    print(
        f"The test data has {test_data.size(0)} samples and labels that are disjoint from the
        training data."
    )

```

The training data has 8000 samples, each a vector of 36 numbers along with
a corresponding set of 8000 labels, assigning 0.0 or 1.0 to each sample.
The test data has 1000 samples and labels that are disjoint from the training data.

The following code cell serves the purpose of training a neural network, tracking its performance throughout the training process, and generating visual representations of the training progress. It is crucial for students to grasp the functionality of the `run_test` function and `Supervise` class as it will be frequently utilized in the subsequent tasks.

Comments have been thoughtfully included to enhance code comprehension.

```
In [11]: def run_test(net, optmaker, test_every=10):
    # Set up the Loss Function and Optimizer
    optimizer = optmaker(
        net.parameters()
    ) # Initialize the optimizer with model parameters
    print(f"sum([p.numel() for p in net.parameters()]) parameters")
    train_losses, train_accs, test_accs = [], [], []

    for epoch in range(2000):
        loss = net(train_data.float(), train_labels.float())
        loss.backward()
        train_losses.append([epoch, loss.item()])
        optimizer.step() # Update model parameters using the optimizer's update rule
        if epoch % test_every == test_every - 1:
            grads = torch.stack([p.grad.abs().max() for p in net.parameters()])
            maxg, ming = grads.abs().max(), grads.abs().min()
            train_outputs = net.net(train_data.float())
            train_preds = (train_outputs.squeeze() > 0.5).float()
            train_accuracy = (train_preds == train_labels).float().mean()
            train_accs.append([epoch + 1, train_accuracy])
            net.eval()
            test_outputs = net.net(test_data.float())
            net.train()
            test_preds = (test_outputs.squeeze() > 0.5).float()
            test_accuracy = (test_preds == test_labels).float().mean()
            test_accs.append([epoch + 1, test_accuracy])
            print(f"Epoch {epoch+1}, Train Accuracy: {train_accuracy}, Test Accuracy: {test_accuracy}")

```

```

        f"Epoch {epoch+1}, Loss: {loss.item():.5f}, Grad range {maxg:.1e} to
{ming:.1e}, "
        f"Train Accuracy: {train_accuracy.item()}, Test Accuracy:
{test_accuracy.item()}" ,
        end="\r",
    )
    if test_accuracy.item() == 1.0:
        break
optimizer.zero_grad()

# Test the Model
with torch.no_grad():
    train_outputs = net.net(train_data.float())
    train_preds = (train_outputs.squeeze() > 0.5).float()
    train_accuracy = (train_preds == train_labels).float().mean()
    net.eval()
    test_outputs = net.net(test_data.float())
    net.train()
    test_preds = (test_outputs.squeeze() > 0.5).float()
    test_accuracy = (test_preds == test_labels).float().mean()
    print(
        f"\nTrain Accuracy: {train_accuracy.item():.5f}, Test Accuracy:
{test_accuracy.item():.5f}"
    )

# Visualiztion
fig, ax = plt.subplots()
ax2 = ax.twinx()
ax.plot(*zip(*train_losses), label="Training loss")
ax.set_yscale("log")
ax2.plot(*zip(*train_accs), color="orange", label="Training accuracy")
ax2.plot(*zip(*test_accs), color="red", label="Test accuracy")
ax2.set_ylim(0.0, 1.0)
for a in [ax, ax2]:
    for pos in "top right bottom left".split():
        a.spines[pos].set_visible(False)
ax.set_xlabel("Epochs")
ax.set_ylabel("Loss")
ax2.set_ylabel("Accuracy")
fig.legend(loc="lower left", bbox_to_anchor=(0, 0), bbox_transform=ax.transAxes)
fig.show()

print(
    f"Data width {train_data.size(1)}; Constant baseline accuracy {max(test_labels.sum(),
len(test_labels) - test_labels.sum()) / len(test_labels):.3f}"
)

```

Data width 36; Constant baseline accuracy 0.500

The `Supervise` class is a wrapper that combines a neural network model and a loss function to facilitate supervised learning tasks. It computes the loss by performing a forward pass through the neural network and comparing the predicted values to the true labels.

In [12]:

```

class Supervise(Module):
    def __init__(self, criterion, net):

```

```

super().__init__()
self.net = net
self.criterion = criterion

def forward(self, x, y):
    out = self.net(x).squeeze()
    return self.criterion(out, y)

```

Here is an example on how to train a model using Supervise and run_test function

In [13]:

```

#####
# Build a Neural Network which has the architecture as follows:-
# Hidden Dimension - 256
# Loss function - Mean Squared Error (MSE)
# Optimizer - Simple GradientDescent (lr=0.1) [Using the optimizer built in Task 1.3]
# Network Architecture - (Linear + Sigmoid) -> (Linear + Sigmoid) -> (Linear +
# Sigmoid)
#####

input_size = train_data.size(1)
hidden_dims = 256
output_dims = 1

run_test(
    Supervise(
        nn.MSELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.Sigmoid(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.Sigmoid(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: SimpleGradientDescent(p, lr=0.1),
)
#####
# END OF YOUR CODE
#####

```

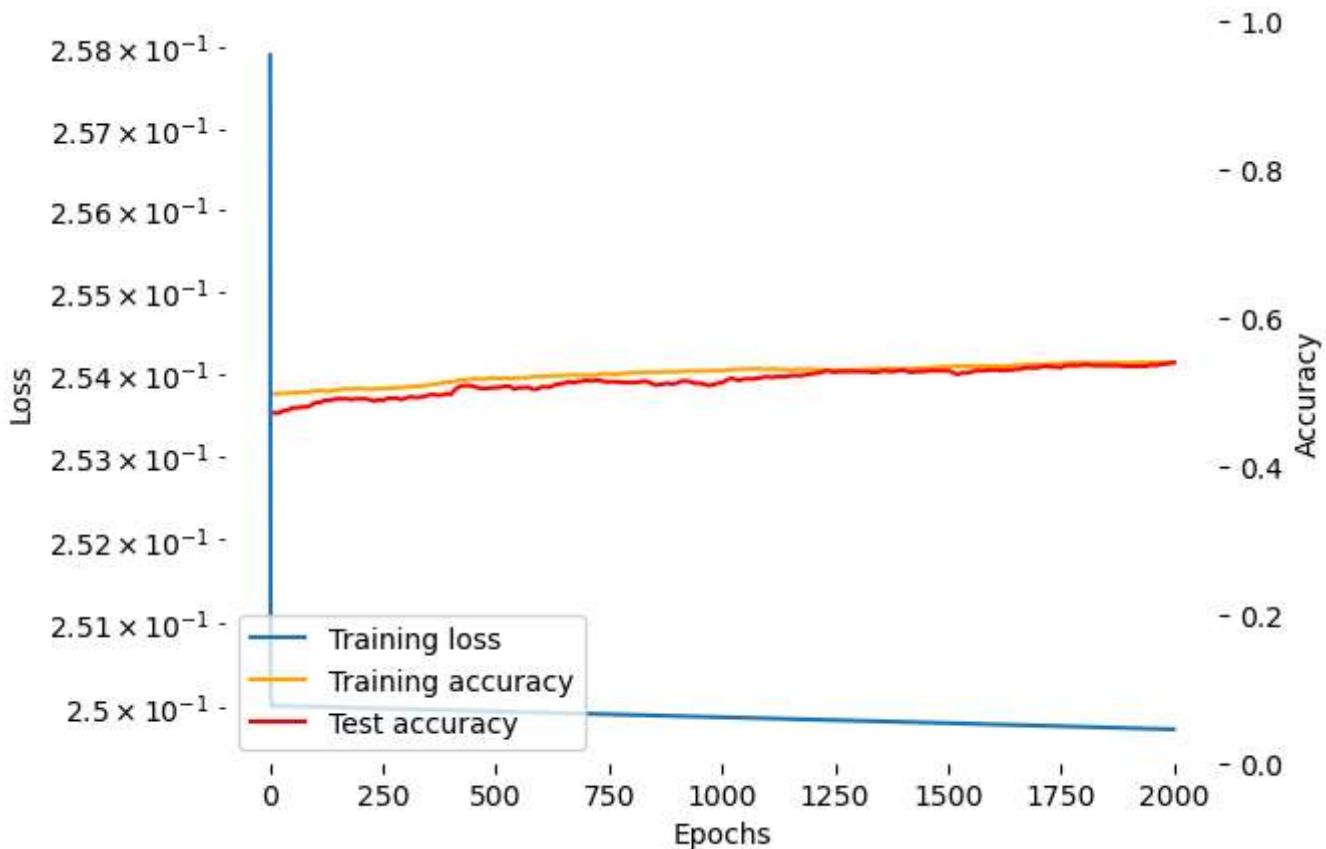
75521 parameters

Epoch 2000, Loss: 0.24973, Grad range 1.3e-04 to 2.3e-06, Train Accuracy: 0.5397499799728394, Test Accuracy: 0.5400000214576721

Train Accuracy: 0.53975, Test Accuracy: 0.54000

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Configuring and training various Neural Network Architectures

Let's construct a range of neural network architectures with varying configurations to explore their impact on training performance when considering factors such as:

- Activation functions
- Choice of optimizers
- Regularization
- Hidden layer dimensions
- Network depth
- Batch normalization
- Residual networks.

This experimentation will help us gain insights into how these factors influence the training process on our dataset.

Note - Utilize the previously defined `SimpleGradientDescent` and `ADAMOptimizer` optimizers exclusively, unless otherwise specified in the task to employ PyTorch's built-in optimizers.

1) Activation functions

Task 3.1 - Test sigmoid activations in a three-layer network with `run_Test` (1 point)

Let's apply our new `run_test` function together with your `SimpleGradientDescent` optimizer class to test (assuming that is the class name you gave it). You will need to put in your three-layer sigmoid network definition to make the example work.

In [14]:

```
torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Create an architecture similar to the example but with reduced hidden
# dimension
# Hidden Dimension - 128
# Loss - MSELoss()
# Optimizer - Simple GradientDescent - (lr = 1.0) [Use the Optimizer built above]
# Network Architecture - (Linear + Sigmoid) -> (Linear + Sigmoid) -> (Linear + Sigmoid)
#####
hidden_dims = 128

run_test(
    Supervise(
        nn.MSELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.Sigmoid(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.Sigmoid(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: SimpleGradientDescent(p, lr=1.0),
)

#####
# END OF YOUR CODE
#####
```

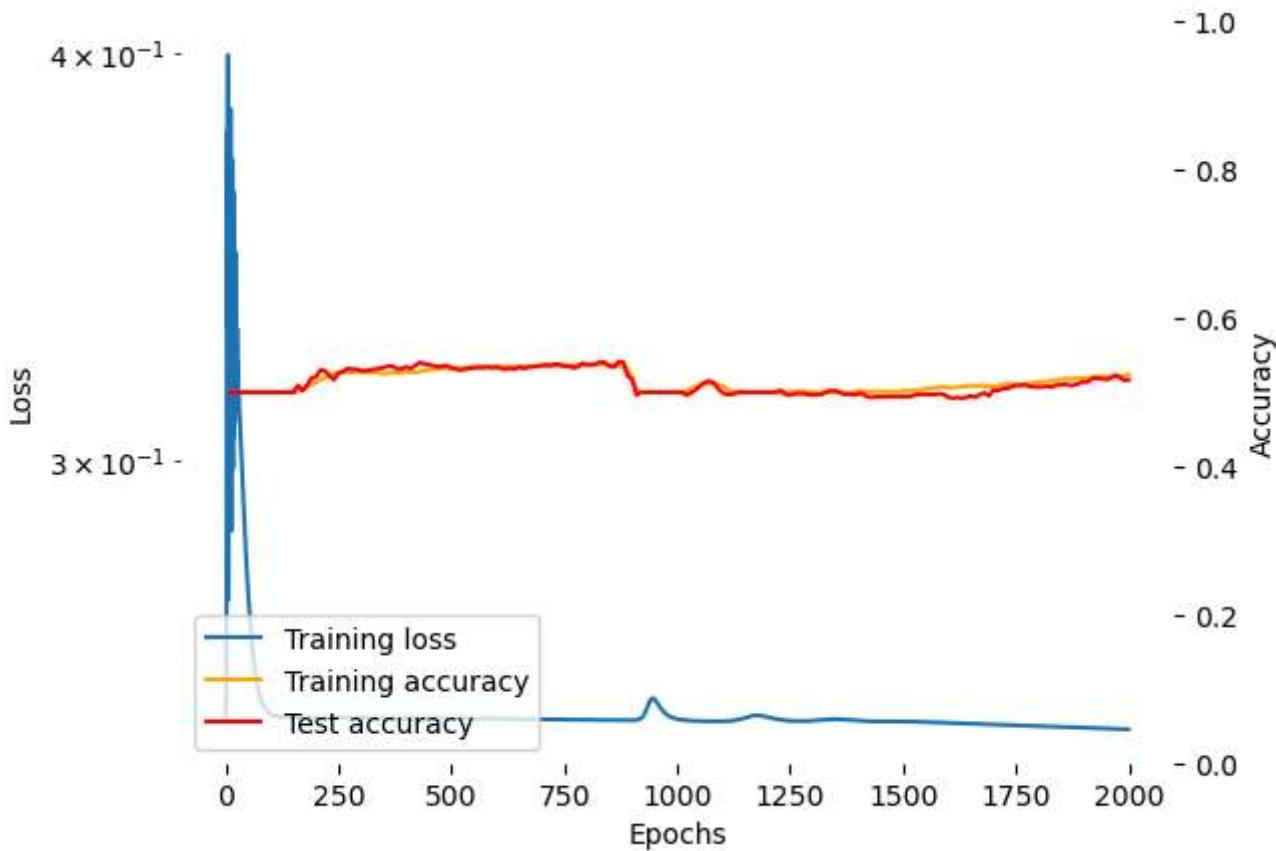
21377 parameters

Epoch 2000, Loss: 0.24789, Grad range 1.5e-02 to 3.6e-04, Train Accuracy: 0.5241249799728394, Test Accuracy: 0.5170000195503235

Train Accuracy: 0.52412, Test Accuracy: 0.51700

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```



Task 3.2 - Test Tanh activations in a three-layer network (1 point)

Implement the same neural network architecture

Read documentation - <https://pytorch.org/docs/stable/generated/torch.nn.Tanh.html>

Now let us compare Sigmoid to other activations. Switch every `nn.Sigmoid()` to `nn.Tanh()`. Tanh balances positive and negative outputs and sometimes work better. Implement and decide whether it is helping in this case.

```
In [15]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Create an architecture similar to the previous one, but replace every
# nn.Sigmoid() activation function with nn.Tanh() (except for the last one).
# Hidden Dimension - 128
# Loss - MSELoss()
# Optimizer - Simple GradientDescent - (lr = 1.0) [Use the Optimizer built above]
# Network Architecture - (Linear + Tanh) -> (Linear + Tanh) -> (Linear + Sigmoid)
#####

hidden_dims = 128

run_test(
    Supervise(
        nn.MSELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.Tanh(),
```

```

        nn.Linear(hidden_dims, hidden_dims),
        nn.Tanh(),
        nn.Linear(hidden_dims, output_dims),
        nn.Sigmoid(),
    ),
),
lambda p: SimpleGradientDescent(p, lr=1.0),
)

#####
# END OF YOUR CODE #
#####

```

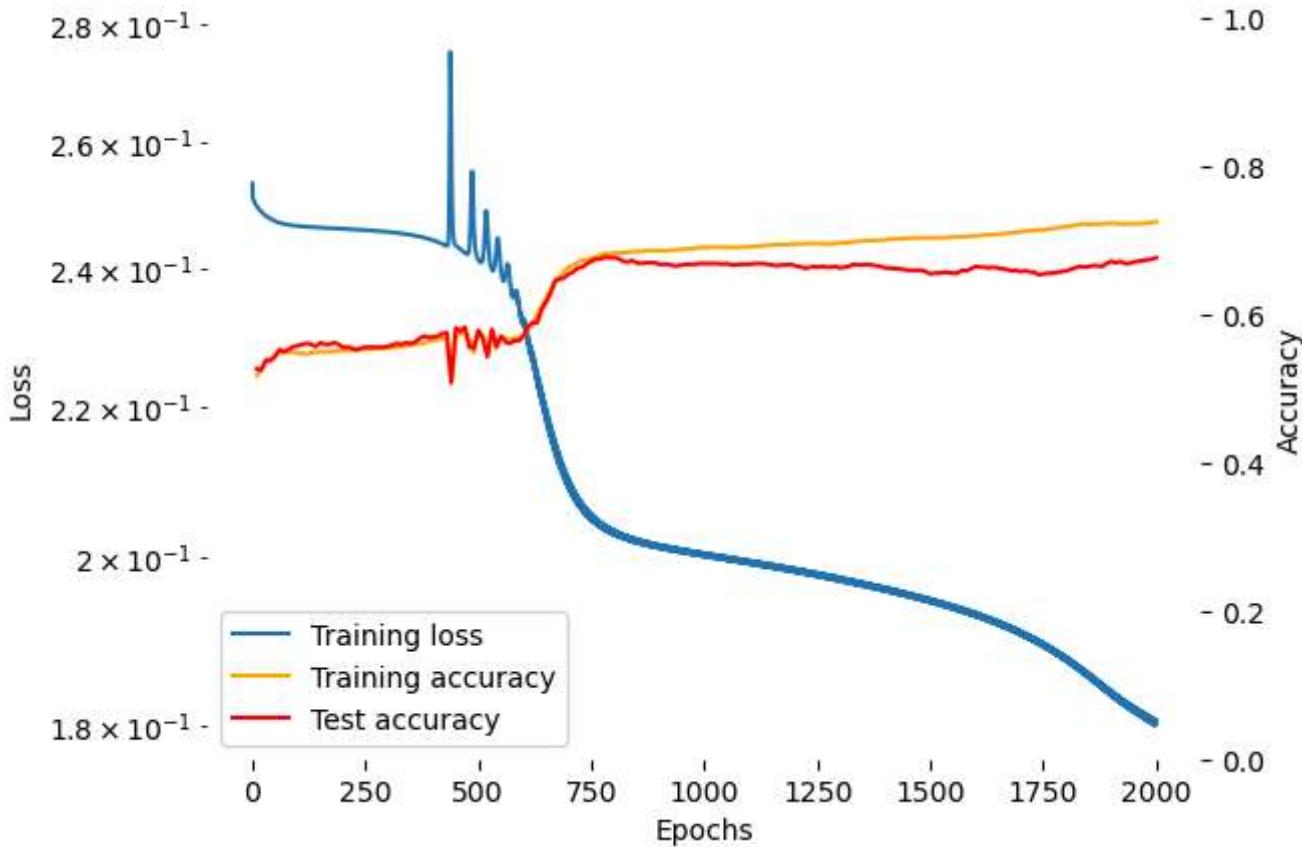
21377 parameters

Epoch 2000, Loss: 0.18052, Grad range 2.6e-02 to 7.4e-03, Train Accuracy: 0.7256249785423279, Test Accuracy: 0.6769999861717224

Train Accuracy: 0.72562, Test Accuracy: 0.67700

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Task 3.3 - Test ReLU activations in a three-layer network (1 point)

The ReLU activation was studied closely by Glorot, which we have discussed in class. It tends to be very effective at avoiding vanishing gradients, because on the positive side it never saturates. Replace all your nonlinearities with ReLU while keeping the architecture otherwise the same. Does ReLU help in this case?

In [16]: `torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)`

```

# Create an architecture similar to the previous one, but replace every
# nn.Tanh() activation function with nn.ReLU() (except for the last one).
# Hidden Dimension - 128
# Loss - MSELoss()
# Optimizer - SimpleGradientDescent - (Lr = 1.0)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####
hidden_dims = 128

run_test(
    Supervise(
        nn.MSELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: SimpleGradientDescent(p, lr=1.0),
)

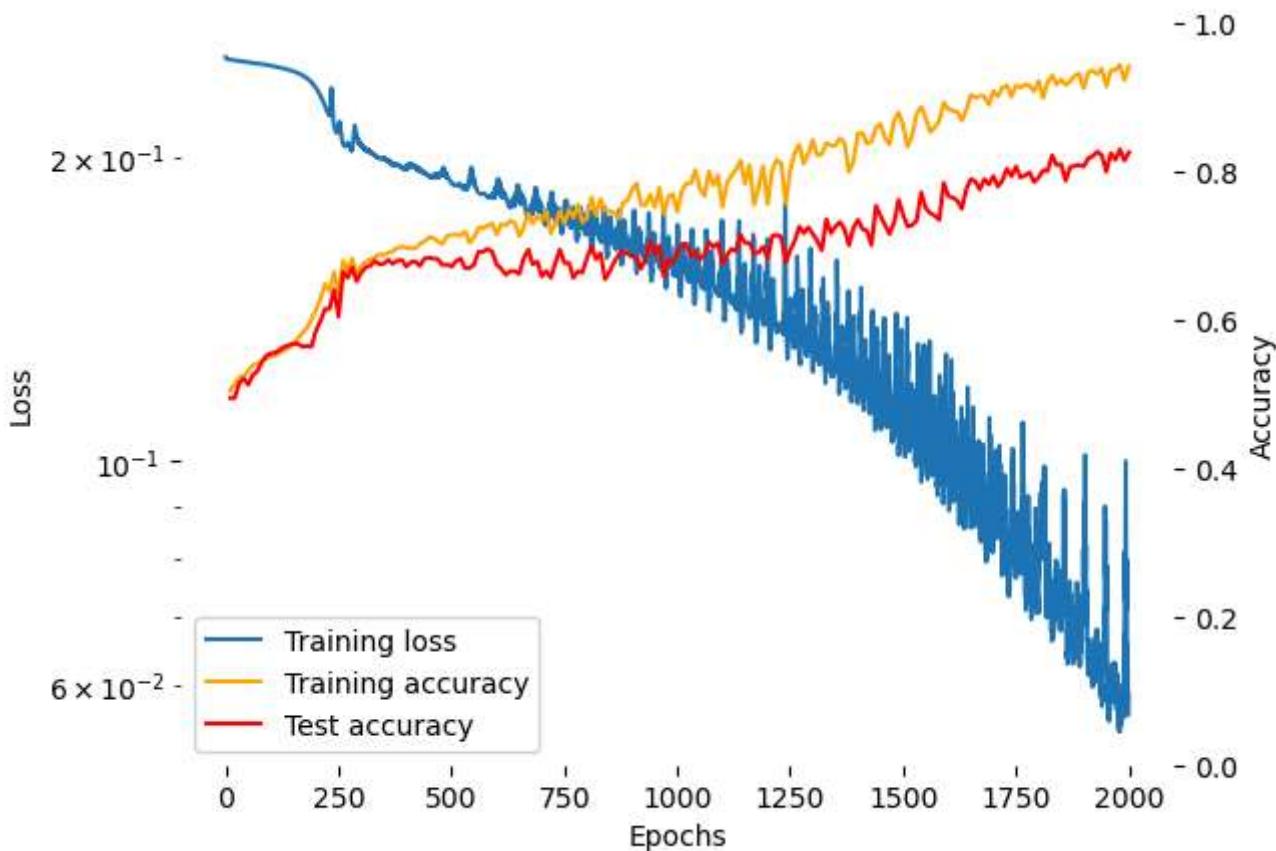
#####
# END OF YOUR CODE
#
#####

```

21377 parameters
Epoch 2000, Loss: 0.05873, Grad range 2.5e-02 to 7.4e-03, Train Accuracy: 0.9421250224113464, Test Accuracy: 0.8259999752044678
Train Accuracy: 0.94213, Test Accuracy: 0.82600

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



3.A) Inline Question (1 point):

In the above Tasks, we employed three different activation functions—sigmoid, tanh, and relu—in our neural network architecture. Describe the difference in the behavior of the optimization process that you observed between Sigmoid, TanH, and ReLU. Do your results confirm or contradict the results that Glorot reported in his 2010 study?

Sigmoid struggled during the process and end up with the least progress.

Tanh did some progress until around epoch 750, and then the test acc started dropping.

ReLU made the most progress consistently over the whole training.

The results conform to the study by Glorot.

Task 3.4 - Implement Binary Cross Entropy Loss (1 point)

In a classification setting, we often prefer to interpret the outputs as probabilities and drive the probability distribution towards the true distribution. The standard way to achieve that is to use the cross-entropy loss. Cross-entropy (as seen in HW1) also test to avoid saturation when compared to MSE, when used in combination with softmax.

Replace the supervision with `BCELoss` rather than mean square error, and observe any differences.

Read documentation for BCE Loss - <https://pytorch.org/docs/stable/generated/torch.nn.BCELoss.html>

```
In [17]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Build a neural network with following architecture:
# Hidden Dimension - 128
# Loss - Binary Cross Entropy
# Optimizer - SimpleGradientDescent - (lr = 1.0)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####
hidden_dims = 128

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: SimpleGradientDescent(p, lr=1.0),
)

#####
#           END OF YOUR CODE
#####
```

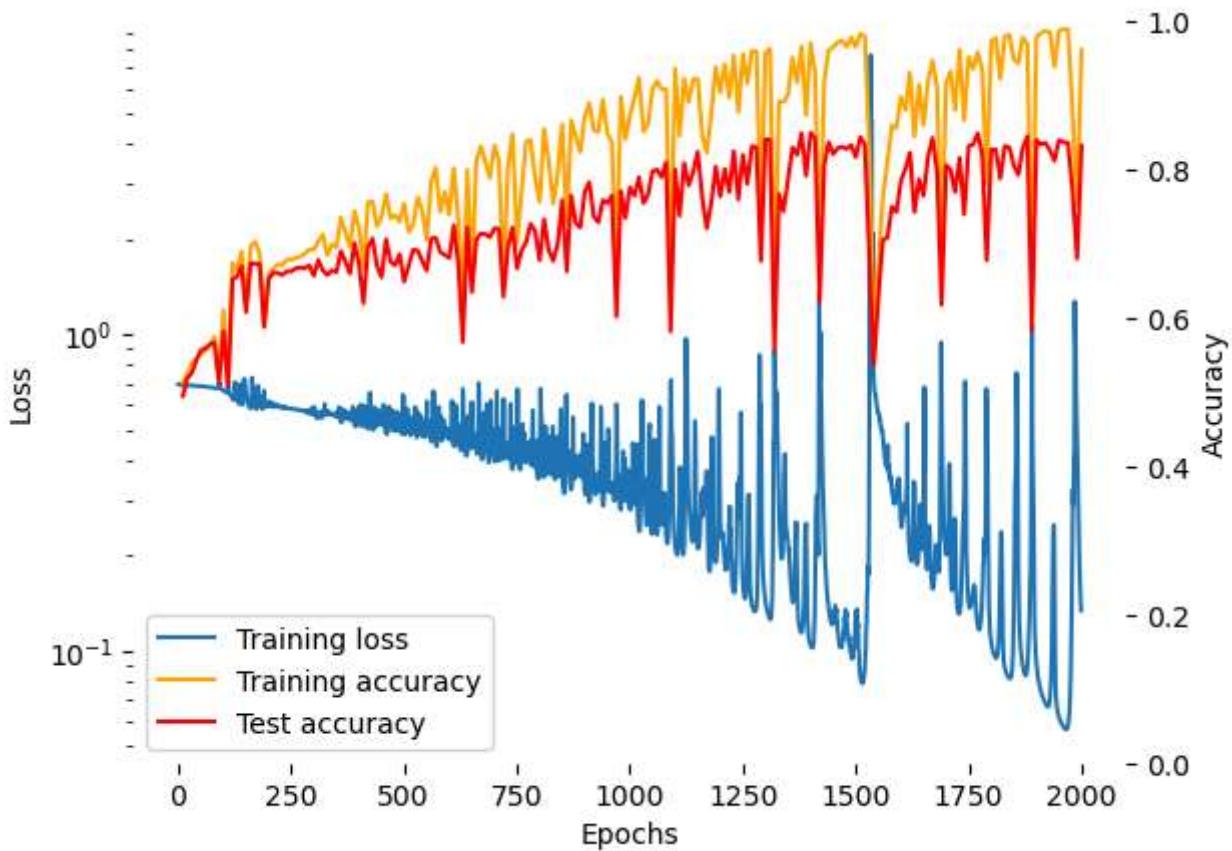
21377 parameters

Epoch 2000, Loss: 0.13423, Grad range 9.6e-03 to 1.1e-03, Train Accuracy: 0.9613749980926514, Test Accuracy: 0.8320000171661377

Train Accuracy: 0.96137, Test Accuracy: 0.83200

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```



2) Optimizers

Task 3.5 Implement using SGD Optimizer (1 point)

So far we have been using our own `SimpleGradientDescent`. Now try comparing results with pytorch's built-in `torch.optim.SGD` class. How does your implementation compare? Is it the same?

Read Documentation for SGD Optimizer using Pytorch -

<https://pytorch.org/docs/stable/generated/torch.optim.SGD.html>

```
In [18]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Build a neural network with following architecture, Careful, Hidden
# Dimension has changed.
# Hidden Dimension - 128
# Loss - Binary Cross Entropy
# Optimizer - SGD - (lr = 0.5)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####

hidden_dims = 128

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
```

```

        nn.Linear(hidden_dims, hidden_dims),
        nn.ReLU(),
        nn.Linear(hidden_dims, output_dims),
        nn.Sigmoid(),
    ),
),
lambda p: torch.optim.SGD(p, lr=0.5),
)

#####
#           END OF YOUR CODE
#####

```

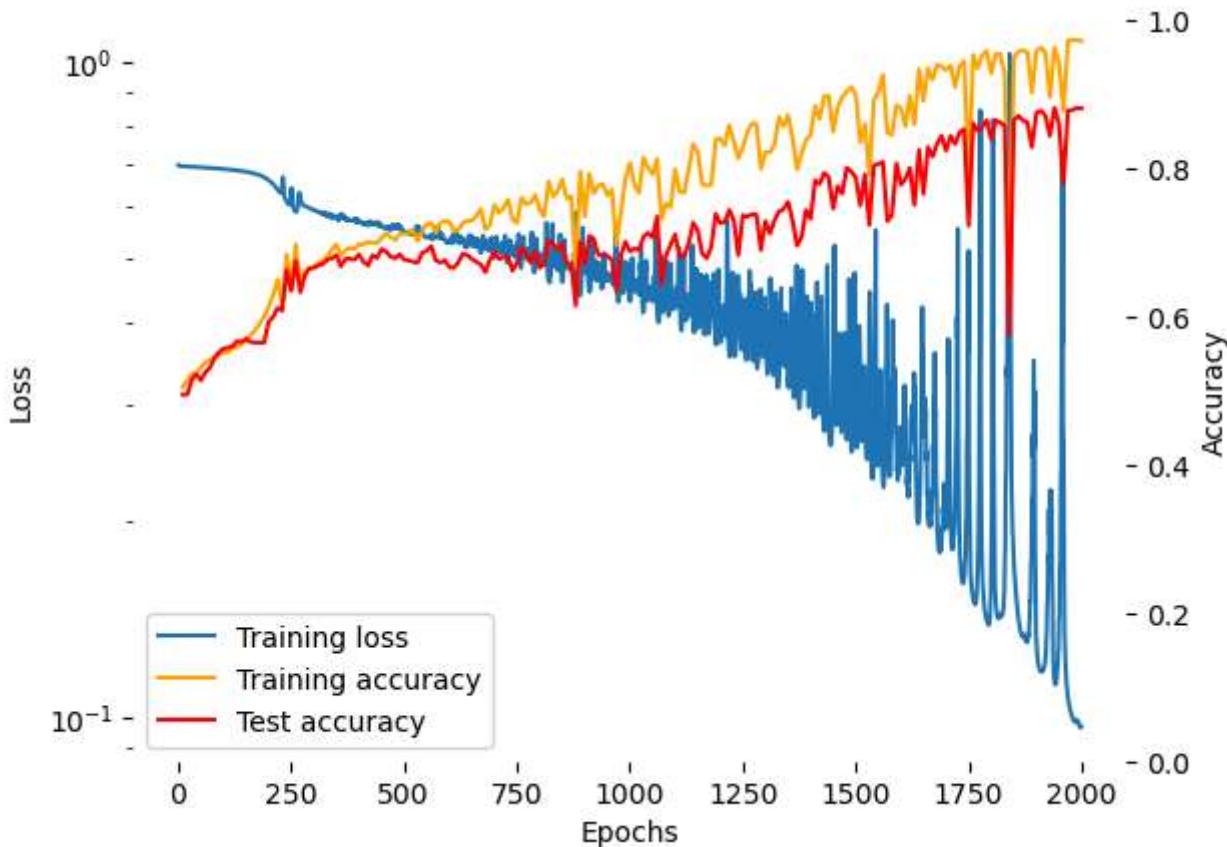
21377 parameters

Epoch 2000, Loss: 0.09705, Grad range 1.4e-02 to 6.8e-03, Train Accuracy: 0.971875011920929, Test Accuracy: 0.8809999823570251

Train Accuracy: 0.97188, Test Accuracy: 0.88100

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Task 3.6 - ADAM Optimizer (1 point)

ADAM is a very powerful optimizer and should improve results.

Use your own `ADAMOptimizer` class here to see how it behaves. If you wish to debug against the standard ADAM optimizer, then you can try it out as well, but when you hand in your results, show what your `ADAMOptimizer` does. Ideally, they should behave the same.

Read Documentation for ADAM Optimizer using Pytorch -
<https://pytorch.org/docs/stable/generated/torch.optim.Adam.html>

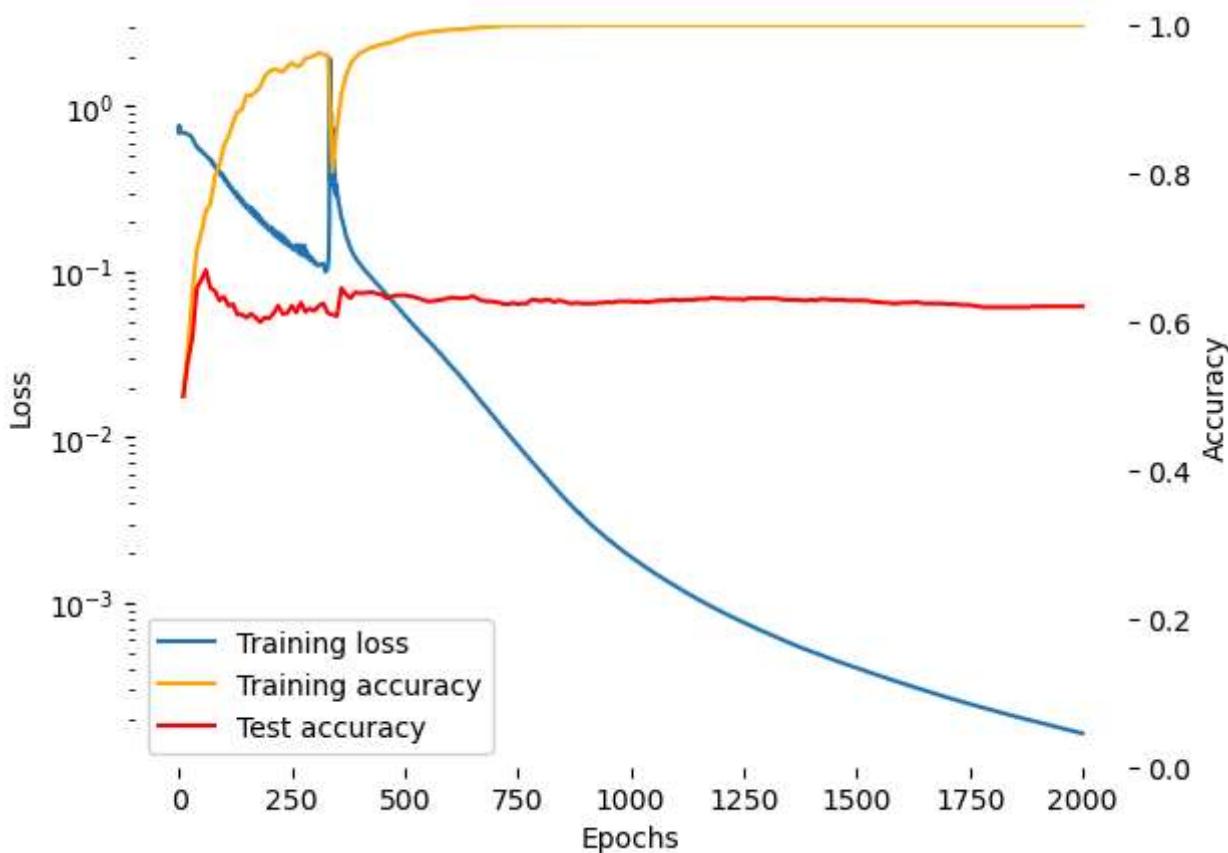
```
In [19]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Build a neural network with following architecture:
# Hidden Dimension - 128
# Loss - Binary Cross Entropy
# Optimizer - ADAM (weight decay = 0 ,lr = 0.01)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####

hidden_dims = 128

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=0.0),
)

#####
# END OF YOUR CODE
#
```

```
21377 parameters
Epoch 2000, Loss: 0.00016, Grad range 1.1e-04 to 3.8e-06, Train Accuracy: 1.0, Test Accuracy: 0.6209999918937683 9871253967
Train Accuracy: 1.00000, Test Accuracy: 0.62100
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()
```



3.B) Inline Question (2 points): Does ADAMOptimizer do better? Explain how the optimization behaves differently than you observed with Simple Gradient Descent.

No, it doesn't.

Adam achieved equally good performance on training data, but failed to generalize to test data.

The Battle Against Overfitting

Let's apply some techniques to enhance and refine our model.

A) The Art of Regularization

Task 3.7 - Add a penalty term to the loss function that discourages large weights. This helps prevent the model from fitting noise in the data. (1 point)

```
In [20]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Design a neural network with the specified architecture and introduce a
# regularization term in the loss function to discourage the growth of large
# weights.
# Hidden Dimension - 128
# Loss - Binary Cross Entropy
# Optimizer - Adam with weight decay - (lr = 0.01, weight_decay=1e-3)
```

```

# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####
hidden_dims = 128

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)
#####

# END OF YOUR CODE #
#####

```

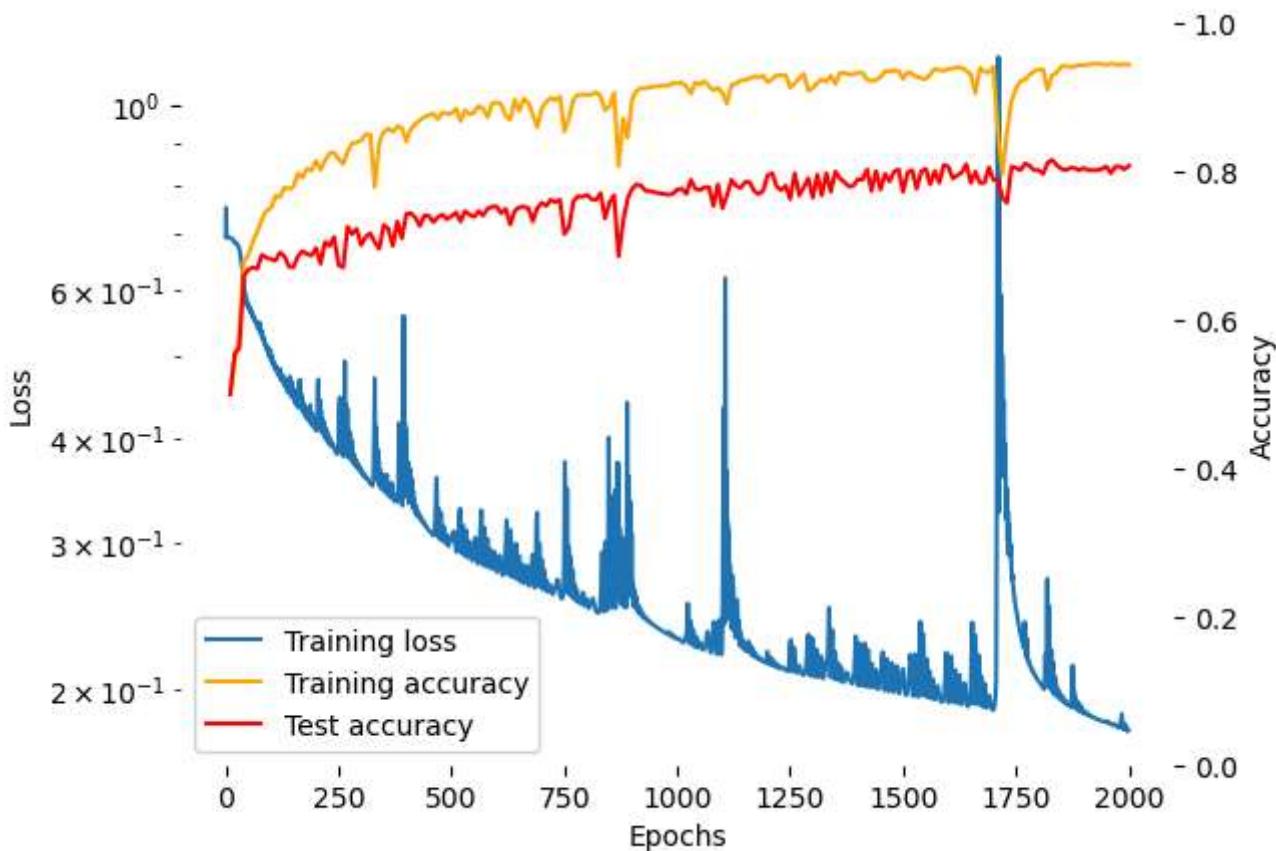
21377 parameters

Epoch 2000, Loss: 0.17888, Grad range 5.9e-02 to 4.9e-03, Train Accuracy: 0.9441249966621399, Test Accuracy: 0.8080000281333923

Train Accuracy: 0.94412, Test Accuracy: 0.80800

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



B) Slimming Down for Success

Reducing Hidden Dimension - Fewer parameters means that the model has less flexibility to fit the training data and it is forced to learn simpler features that are more likely to generalize to new data.

Task 3.8 - Going Deeper with Shrinking Hidden Dimensions (1 point)

Let's explore the effectiveness of dimension reduction as a technique by reducing the number of hidden dimensions from 128 to 64.

```
In [21]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 64
# Loss - Binary Cross Entropy
# Optimizer - Adam with weight decay - (Lr = 0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####

hidden_dims = 64

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
```

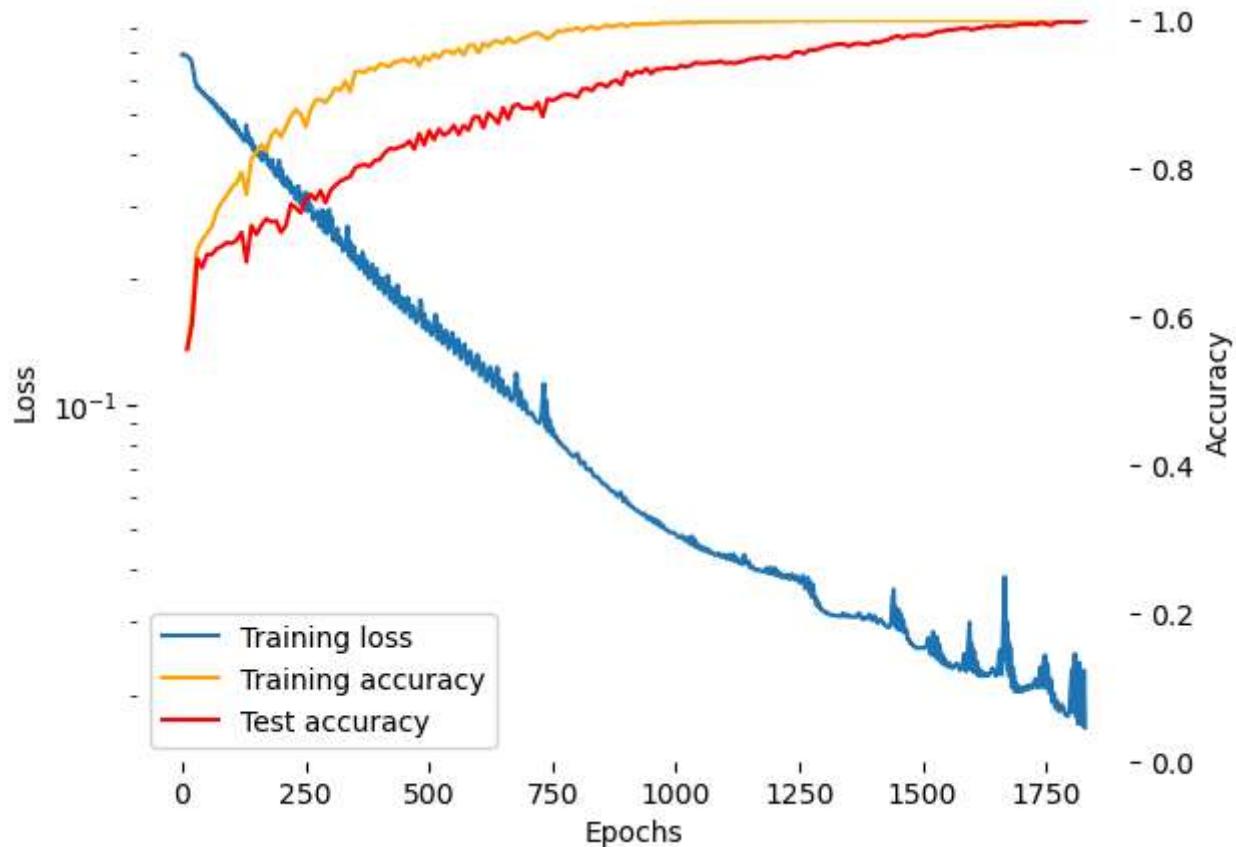
```

        nn.ReLU(),
        nn.Linear(hidden_dims, output_dims),
        nn.Sigmoid(),
    ),
),
lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

#####
#           END OF YOUR CODE
#
#####

6593 parameters
Epoch 1830, Loss: 0.01667, Grad range 1.4e-02 to 1.6e-03, Train Accuracy: 1.0, Test Accuracy:
1.0  000257492065  00257492065
Train Accuracy: 1.00000, Test Accuracy: 1.00000
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()

```



Task 3.10 - Going more Deeper with Shrinking Hidden Dimensions (1 point)

Dimension reduction improved our outcome. Let's continue by further reducing the hidden dimension from 64 to 32.

```
In [22]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 32
```

```

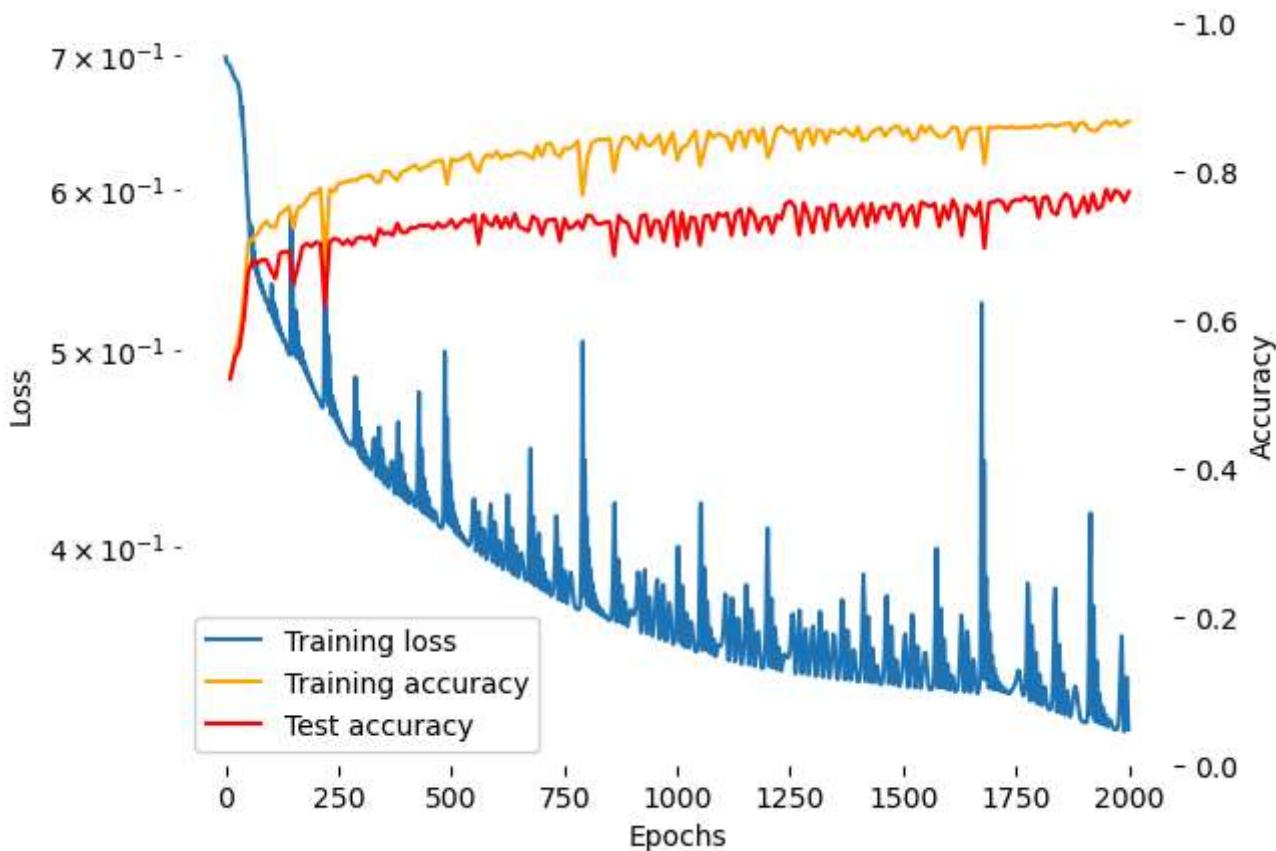
# Loss - Binary Cross Entropy
# Optimizer - Adam with weight decay - (lr = 0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####
hidden_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

#####
# END OF YOUR CODE
#####

```

2273 parameters
Epoch 2000, Loss: 0.32464, Grad range 1.1e-01 to 1.5e-02, Train Accuracy: 0.8679999709129333, Test Accuracy: 0.7730000019073486
Train Accuracy: 0.86800, Test Accuracy: 0.77300
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()



Task 3.11 - The Last Shrink (1 point)

Did it help to reduce to 32 dimensions?

Let's try one more time, maybe it will work? (Below, please try reducing the hidden dimension from 32 to 16.)

```
In [23]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 16
# Loss - Binary Cross Entropy
# Optimizer - Adam with weight decay - (lr = 0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####

hidden_dims = 16

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
)
```

```
        ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)
```

```
#####
# END OF YOUR CODE
#####

```

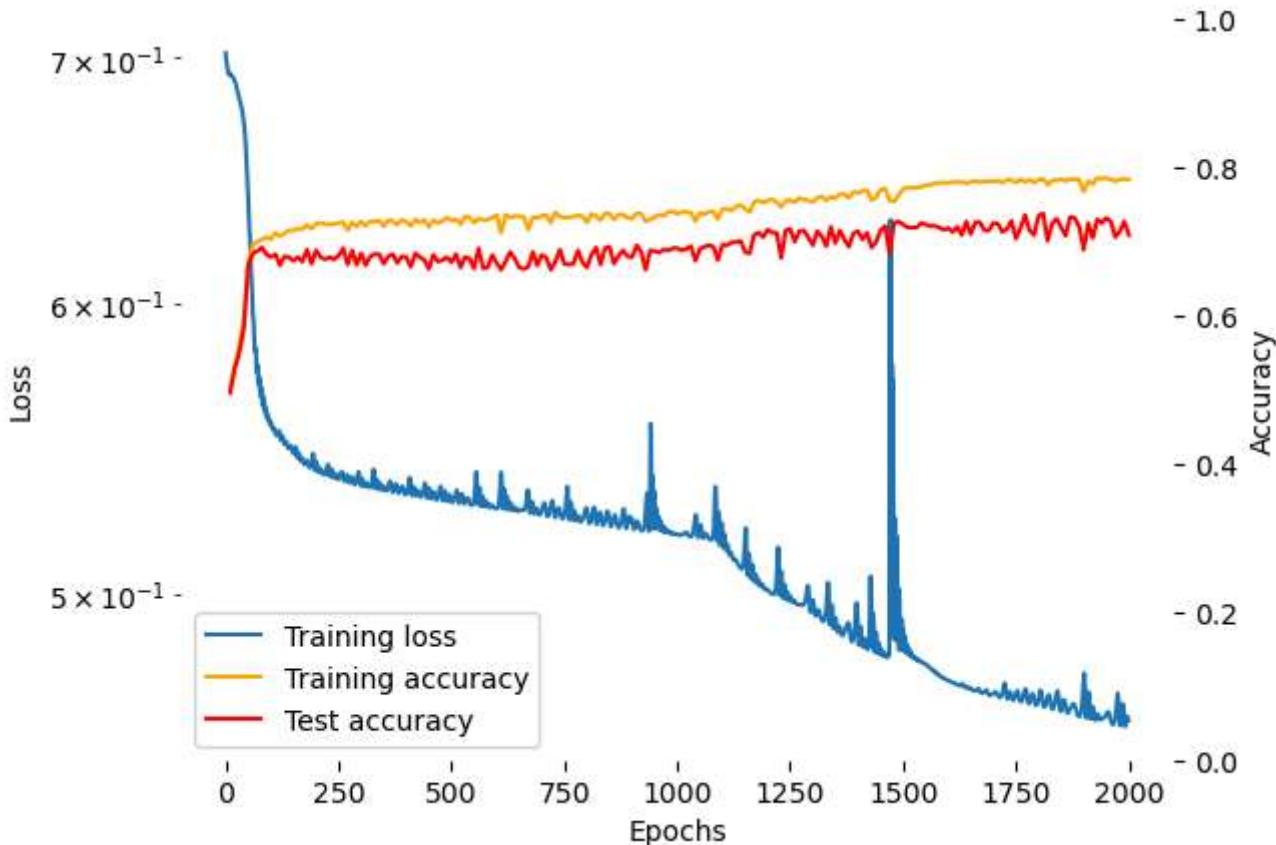
881 parameters

Epoch 2000, Loss: 0.46186, Grad range 8.4e-02 to 2.4e-02, Train Accuracy: 0.7836250066757202, Test Accuracy: 0.7089999914169312

Train Accuracy: 0.78363, Test Accuracy: 0.70900

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```



3.C) Inline Question - In the above Tasks [3.8 - 3.11] (Slimming Down for success) what were your key observations? (1 point)

The less the dimensions:

- The more epochs it takes the peak performance on train data.
- The less gap (overfitting) between train and test accuracy.

C) Layer by Layer

Increasing the number of layers.

Geoff Hinton likes to assert that deeper layers can capture increasingly abstract and high-level features in the data. This hierarchy allows the network to focus on relevant patterns and discard noise, making it less prone to fitting random variations in the training data.

Is it true? Let's try it.

Task 3.12 - Increasing the Depth of the network (1 point)

Let's enhance our model by introducing an additional layer, creating a network with four layers.

```
In [24]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 64
# Loss - Binary Cross Entropy
# Optimizer - Adam (lr = 0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + ReLU)
# -> (Linear + Sigmoid)
#####

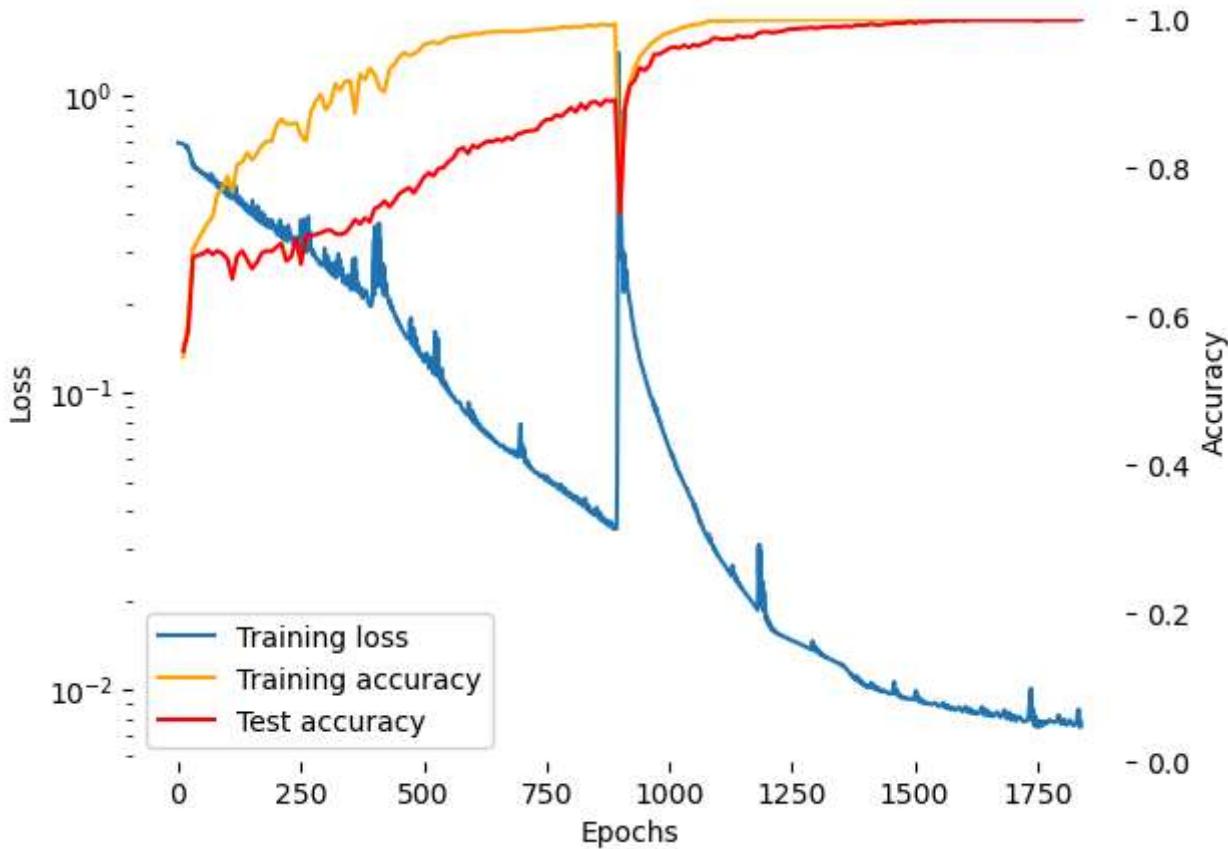
hidden_dims = 64

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)
#####

#                                     END OF YOUR CODE
#####
```

```
10753 parameters
Epoch 1840, Loss: 0.00777, Grad range 2.5e-02 to 2.1e-04, Train Accuracy: 1.0, Test Accuracy:
1.0  000128746033  00047683716
Train Accuracy: 1.00000, Test Accuracy: 1.00000
```

```
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()
```



Task 3.13 - Increasing the Depth and reducing the width of the network (1 point)

The recent modification yielded remarkable results. Now, let's take it a step further by enhancing our model's architecture: we'll increase the number of layers from 4 to 6 and reduce the hidden dimension from 64 to 32.

```
In [26]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 32
# Loss - Binary Cross Entropy
# Optimizer - Adam with weight decay - (Lr=0.005, weight_decay=1e-4)
# Network Architecture - (Linear + ReLU) -> (Linear + ReLU) -> (Linear + ReLU)
# -> (Linear + ReLU) -> (Linear + ReLU) -> (Linear + Sigmoid)
#####

hidden_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, 1),
            nn.Sigmoid()
        )
    )
)
```

```

        nn.Linear(hidden_dims, hidden_dims),
        nn.ReLU(),
        nn.Linear(hidden_dims, hidden_dims),
        nn.ReLU(),
        nn.Linear(hidden_dims, output_dims),
        nn.Sigmoid(),
    ),
),
lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

#####
# END OF YOUR CODE
#
#####

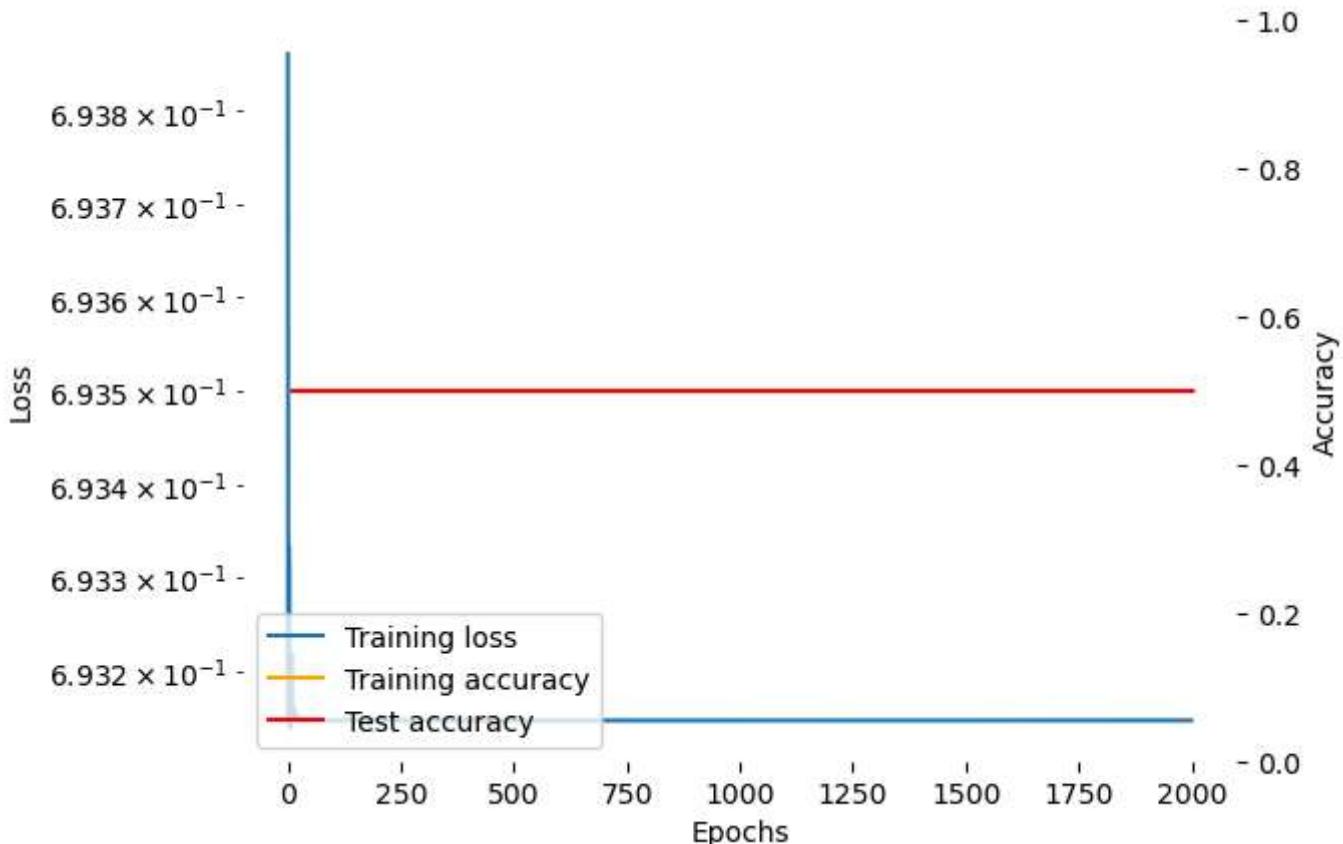
```

5441 parameters

Epoch 2000, Loss: 0.69315, Grad range 4.7e-10 to 0.0e+00, Train Accuracy: 0.5, Test Accuracy: 0.5
Train Accuracy: 0.50000, Test Accuracy: 0.50000

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



D) Batch Normalization

Batch Normalization is a technique used to improve the training of deep neural networks. It works by normalizing the activations of each layer, which helps to prevent the network from becoming too sensitive to the initialization of the weights and the order of the training data."

Pytorch documentation -

<https://pytorch.org/docs/stable/generated/torch.nn.BatchNorm1d.html#torch.nn.BatchNorm1d>

Task 3.14 Adding Batch Normalization (1 point)

In the previous tasks, our model with a hidden dimension of 32 didn't deliver the desired performance. To address this, let's incorporate Batch Normalization into that architecture and assess whether it can enhance its performance.

Read documentation - <https://pytorch.org/docs/stable/generated/torch.nn.BatchNorm1d.html>

In [27]:

```
torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# Let's create a neural network with following architecture:
# Hidden Dimension - 32
# Loss - Binary Cross Entropy
# Optimizer - Adam (lr = 0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU + batch_normalization) ->
# (Linear + ReLU) -> (Linear + ReLU) -> (Linear + ReLU) ->
# (Linear + ReLU + batch_normalization) -> (Linear + Sigmoid)
#####

hidden_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.BatchNorm1d(hidden_dims),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.BatchNorm1d(hidden_dims),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

#####
# END OF YOUR CODE
#####
```

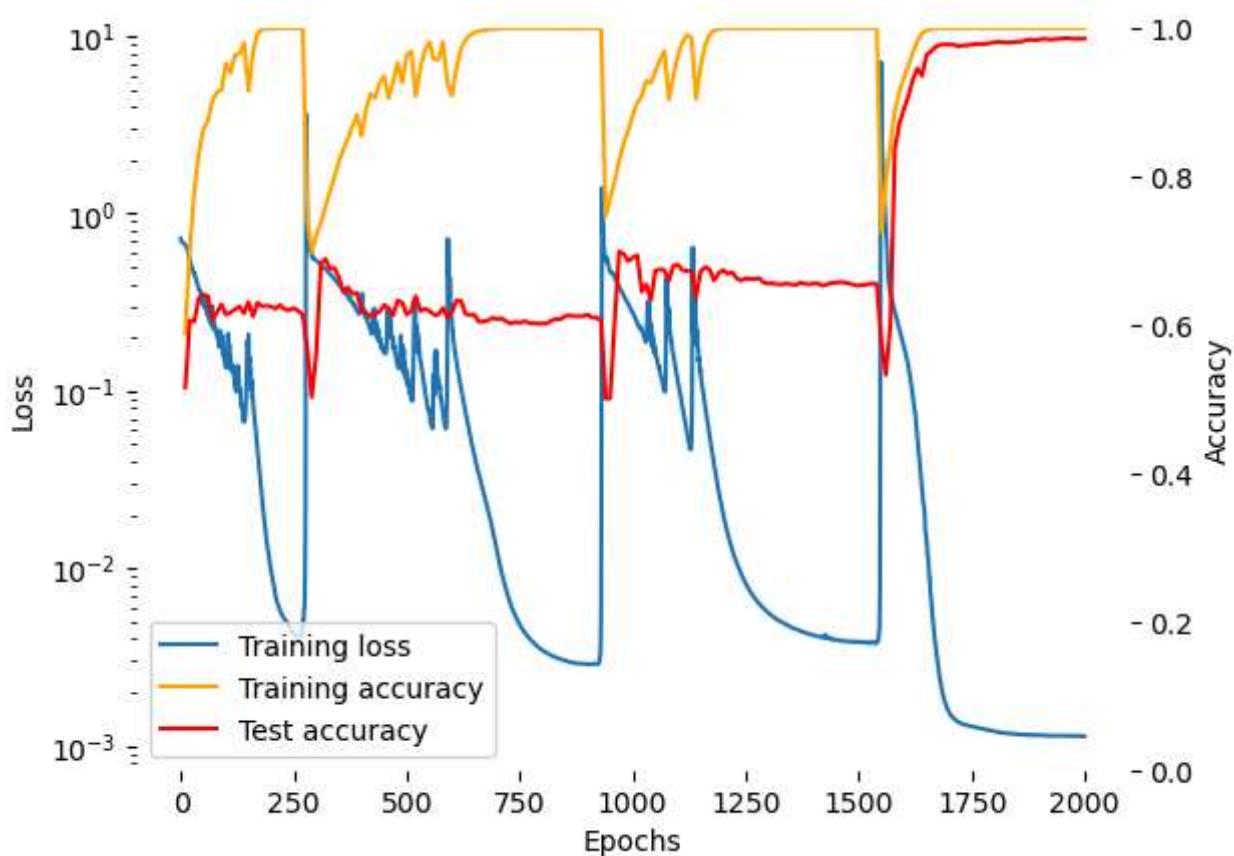
5569 parameters

Epoch 2000, Loss: 0.00113, Grad range 4.2e-03 to 2.1e-04, Train Accuracy: 1.0, Test Accuracy:

0.986000014305115 00238418579

Train Accuracy: 1.00000, Test Accuracy: 0.98600

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.



E) Residual Networks

Residual Networks (ResNets) are a type of deep neural network that are designed to address the problem of vanishing gradients.

Task 3.15 - Modify the code below so that that `ResidualSequence` does not just implement $y = f(x)$ but instead implements $y = f(x) + x$. The template code has a bug and only implement $y=x$. (1 point)

In [28]:

Example of Residual Block Architecture

```
torch.nn.Sequential(  
    ...  
    nn.Linear(train_data.size(1), hidden_dims),  
    ResidualSequence(  
        ...  
        nn.ReLU(),  
        ...  
        nn.Linear(hidden_dims, fan_out_dims),  
        nn.ReLU(),  
        nn.Linear(fan_out_dims, hidden_dims),  
    ),  
    nn.Linear(hidden_dims, 1),  
    nn.Sigmoid()  
)
```

Task 3.16 - Implement Residual blocks in a Neural Network architecture (2 point)

Design a neural network with four Residual blocks, each composed according to the specifications outlined below.

In [29]:

```
torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)  
#####  
# In this task, we will configure a neural network consisting of four  
# residual blocks  
# Hidden Dimension = 16  
# fan_out_dims = 32  
# Loss - Binary Cross Entropy  
# Optimizer - Adam- (lr=0.01, weight_decay=1e-3)  
# Residual Block - [ReLU → Linear → ReLU → Linear]  
# Network Architecture - Linear -> Residual Block -> Residual Block ->  
# Residual Block -> Residual Block -> (Linear + Sigmoid)  
#####  
  
hidden_dims = 16  
fan_out_dims = 32  
  
run_test(  
    Supervise(  
        nn.BCELoss(),  
        nn.Sequential(  
            nn.Linear(input_size, hidden_dims),  
            ResidualSequence(  
                nn.ReLU(),  
                nn.Linear(hidden_dims, fan_out_dims),  
                nn.ReLU(),  
                nn.Linear(fan_out_dims, hidden_dims),  
            ),  
            ResidualSequence(  
                nn.ReLU(),  
                nn.Linear(hidden_dims, fan_out_dims),  
            ),  
            nn.Linear(hidden_dims, 1),  
            nn.Sigmoid()  
        )  
    )  
)
```

```

        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    ResidualSequence(
        nn.ReLU(),
        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    ResidualSequence(
        nn.ReLU(),
        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    nn.Linear(hidden_dims, 1),
    nn.Sigmoid(),
),
),
lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

```

```

#####
#           END OF YOUR CODE
#####

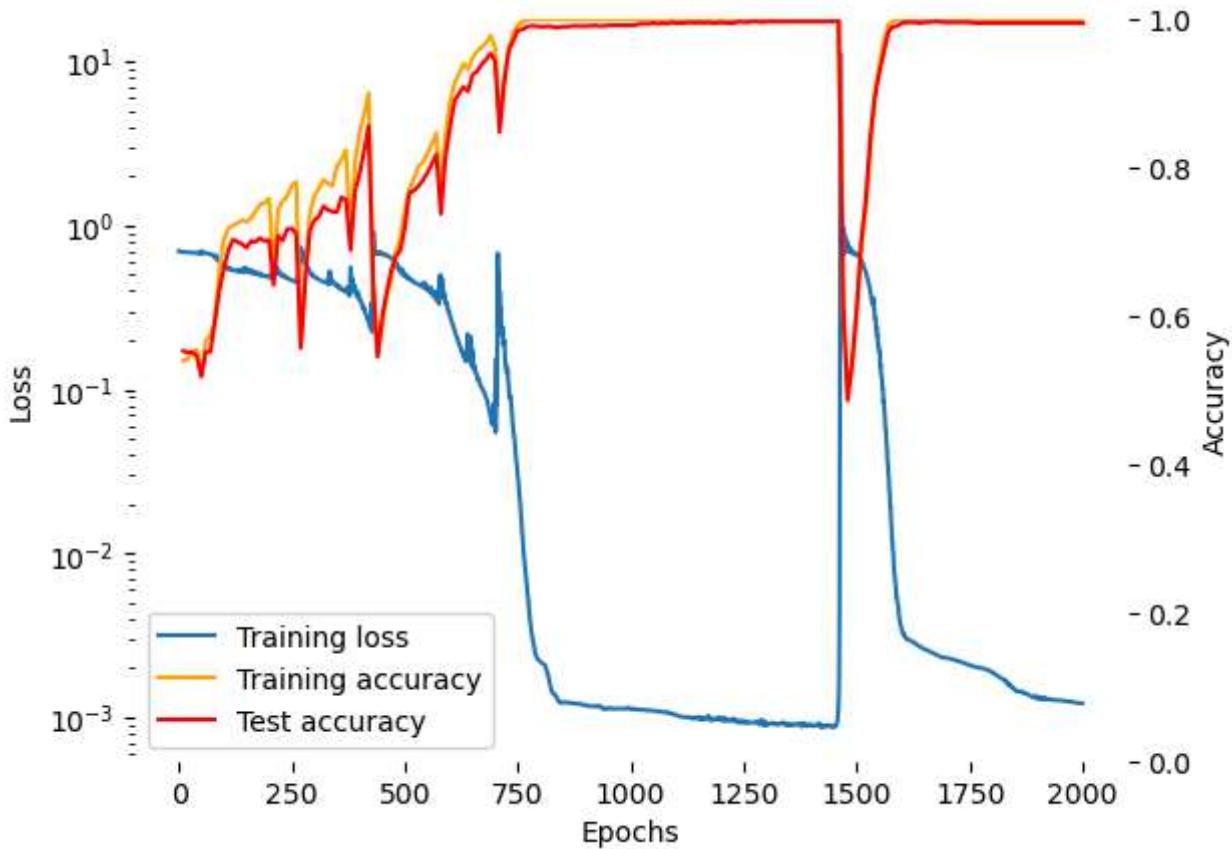
```

4897 parameters

Epoch 2000, Loss: 0.00122, Grad range 2.2e-03 to 1.5e-04, Train Accuracy: 1.0, Test Accuracy: 0.9950000047683716 00047683716
Train Accuracy: 1.00000, Test Accuracy: 0.99500

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Task 3.17 - Reducing Residual Blocks (2 points)

Let's decrease the number of Residual Blocks and observe whether it has any impact on our performance.

```
In [30]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
hidden_dims = None
fan_out_dims = None
#####
# In this task, we will configure a neural network consisting of two
# residual blocks
# Hidden Dimension = 16
# fan_out_dims = 32
# Loss - Binary Cross Entropy
# Optimizer - Adam - (Lr=0.01, weight_decay=1e-3)
# Residual Block - [ReLU > Linear > ReLU > Linear]
# Network Architecture - Linear -> Residual Block ->
# (Linear + Sigmoid)
#####

hidden_dims = 16
fan_out_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            ResidualSequence(
                nn.ReLU(),
                nn.Linear(hidden_dims, fan_out_dims)
            )
        )
    )
)
```

```

        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    ResidualSequence(
        nn.ReLU(),
        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    nn.Linear(hidden_dims, 1),
    nn.Sigmoid(),
),
),
lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

```

```

#####
#           END OF YOUR CODE
#####

```

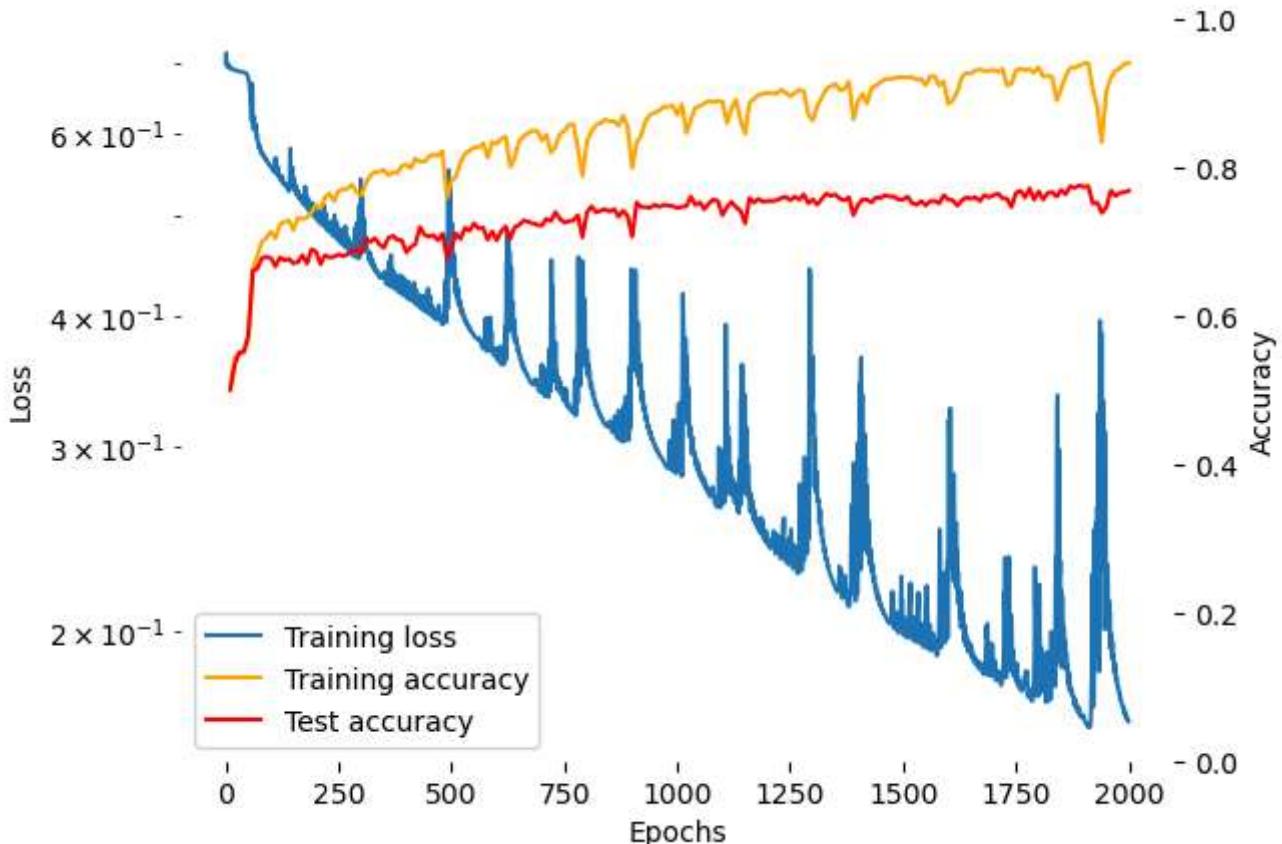
2753 parameters

Epoch 2000, Loss: 0.16382, Grad range 3.8e-02 to 2.4e-03, Train Accuracy: 0.9413750171661377, Test Accuracy: 0.7689999938011169

Train Accuracy: 0.94138, Test Accuracy: 0.76900

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Task 3.18 Develop a Neural Network with a combination of BatchNorm and Residual Blocks (3 points)

The synergy of these two components often results in improved model performance because BatchNorm stabilizes activations and enables the use of deeper networks, while Residual connections facilitate the training of deep networks and prevent degradation in performance. Together, they can enhance the model's ability to learn intricate patterns and improve its generalization to unseen data. However, it's important to strike a balance and avoid overly complex models, as they may lead to overfitting if not properly regularized.

We will create a neural network incorporating both Batch Normalization and Residual Blocks to evaluate if we can achieve favorable outcomes.

```
In [31]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)

#####
# In this task, we will configure a neural network consisting of two
# residual blocks,
# Hidden Dimension - 16
# fan_out_dims = 32
# Loss - Binary Cross Entropy
# Optimizer - Adam - (lr = 0.01, weight_decay=1e-3)
# Residual Block - [Batch_Norm -> Linear -> ReLU -> Linear]
# Network Architecture - Linear -> Residual Block -> Residual Block ->
# (Batch_Norm + Linear + Sigmoid)
#####

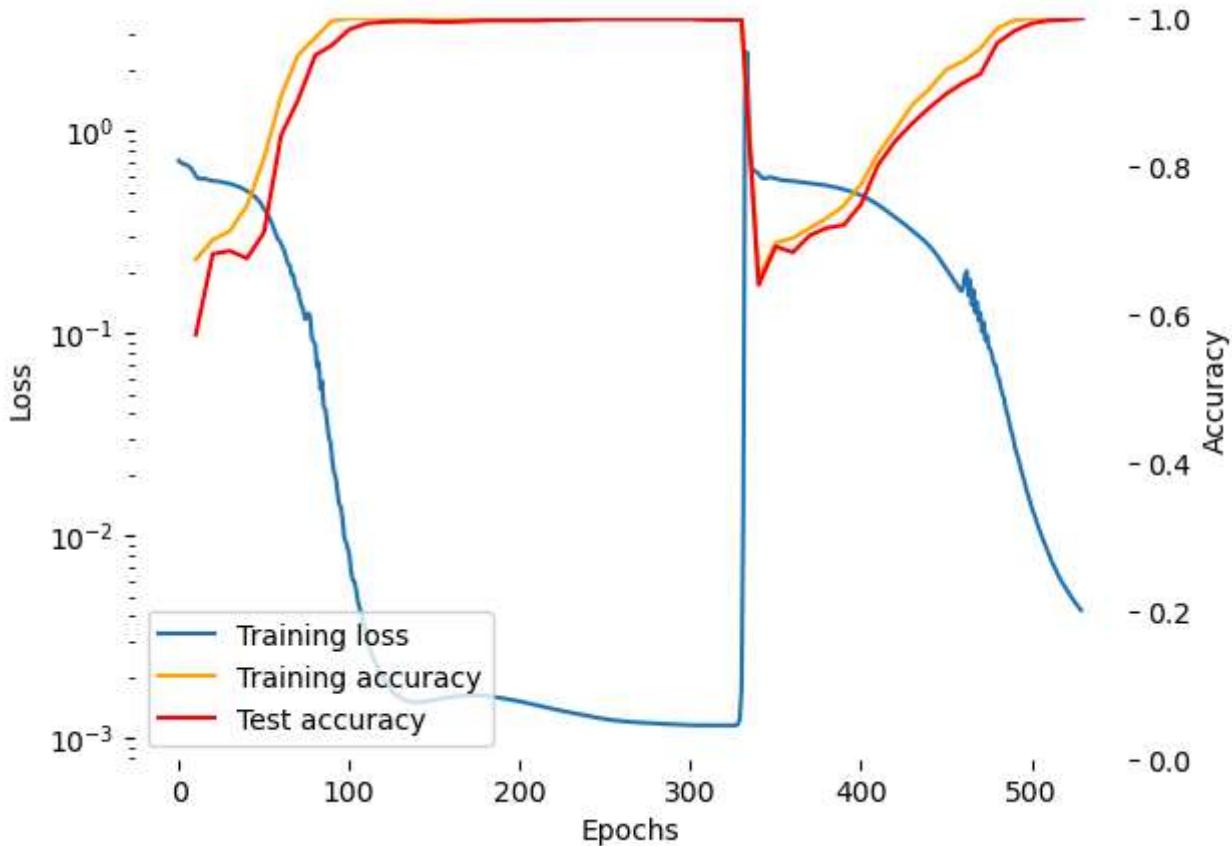
hidden_dims = 16
fan_out_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            ResidualSequence(
                nn.BatchNorm1d(hidden_dims),
                nn.Linear(hidden_dims, fan_out_dims),
                nn.ReLU(),
                nn.Linear(fan_out_dims, hidden_dims),
            ),
            ResidualSequence(
                nn.BatchNorm1d(hidden_dims),
                nn.Linear(hidden_dims, fan_out_dims),
                nn.ReLU(),
                nn.Linear(fan_out_dims, hidden_dims),
            ),
            nn.BatchNorm1d(hidden_dims),
            nn.Linear(hidden_dims, 1),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
```

```

)
#####
#           END OF YOUR CODE
#####
2849 parameters
Epoch 530, Loss: 0.00428, Grad range 3.7e-03 to 1.4e-09, Train Accuracy: 1.0, Test Accuracy: 1.0
0 000257492065 9979019165
Train Accuracy: 1.00000, Test Accuracy: 1.00000
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
fig.show()

```



Task 3.19 - Develop a neural network that incorporates Batch Normalization, Residual connections, and an increased number of layers. (3 points)

We will construct a neural network with 6 residual blocks to assess whether we can further enhance performance by increasing its depth.

```
In [32]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)
#####
# In this task, we will configure a neural network consisting of two
# residual blocks.
# Hidden Dimension - 16
# fan_out_dims = 32
# Loss - Binary Cross Entropy
# Optimizer - Adam - (lr = 0.01, weight_decay=1e-3)
```

```

# Residual Block - [Batch_Norm -> Linear -> ReLU -> Linear]
# Network Architecture - Linear -> Residual Block -> Residual Block ->
# Residual Block -> Residual Block -> Residual Block -> Residual Block
# -> (Batch_Norm + Linear + Sigmoid)
#####
hidden_dims = 16
fan_out_dims = 32

run_test(
    Supervise(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            ResidualSequence(
                nn.BatchNorm1d(hidden_dims),
                nn.Linear(hidden_dims, fan_out_dims),
                nn.ReLU(),
                nn.Linear(fan_out_dims, hidden_dims),
            ),
            nn.BatchNorm1d(hidden_dims),
            nn.Linear(hidden_dims, 1),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

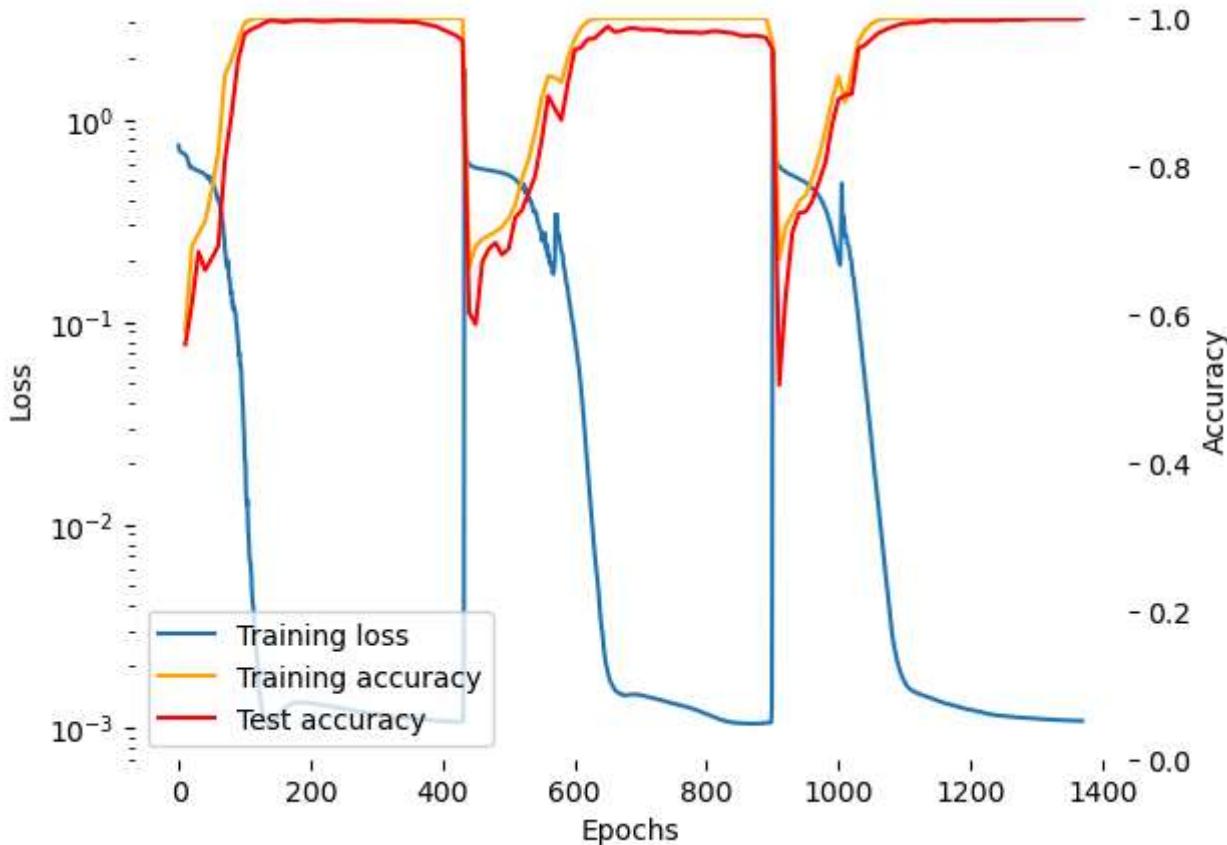
```

```
#####
#                               END OF YOUR CODE
#####
#####
```

```
7265 parameters
Epoch 1370, Loss: 0.00108, Grad range 1.1e-03 to 6.3e-09, Train Accuracy: 1.0, Test Accuracy:
1.0 000128746033 00027179718
Train Accuracy: 1.00000, Test Accuracy: 1.00000
```

```
C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.
```

```
fig.show()
```



Part 4: Weight Initialization

Weight initialization in deep learning refers to the procedure of assigning initial values to a neural network's weights, which are the tunable parameters learned during training. The manner in which these weights are initialized plays a critical role in shaping the training process and ultimately impacts the network's performance. Several weight initialization techniques have been developed to combat challenges such as vanishing or exploding gradients, with the goal of establishing a solid foundation for training deep neural networks. Some of the weight initialization methods are :-

1. Zero Weight Initialization
2. Random Weight Initialization
3. Xavier Initialization (Glorot Initialization)
4. He Initialization (often used in deep CNN's)

1) Zero Initialization

Zero weight initialization initializes a neural network's weights to zero, which can be effective for specific neural network types. This initialization can mitigate overfitting by making it harder for the model to fit training data perfectly.

However, it has drawbacks. It hinders learning complex input-output relationships from scratch and makes the model sensitive to hyperparameters like the learning rate. This approach is generally discouraged due to the potential emergence of symmetric neurons and slow convergence as a result of weight symmetry issues.

The code visualizes layer-wise activation distributions in a neural network for a Zero Weight Initialization. It offers insights into the network's learning and can reveal potential issues.

Task 4.1 - Initialize zero weights for each layer (_ points)

```
In [33]: num_layers = 7
layer_dims = [2048] * num_layers

input_data = np.random.randn(16, layer_dims[0])

activations = []

weights = []

for i in range(num_layers - 1):
    #####
    # Initialize zero weights for each layer (except the last one)
    # and store it in 'W' variable
    #####
    W = np.zeros((layer_dims[i], layer_dims[i + 1]))
    #####
    # END OF YOUR CODE
    #####
    weights.append(W)

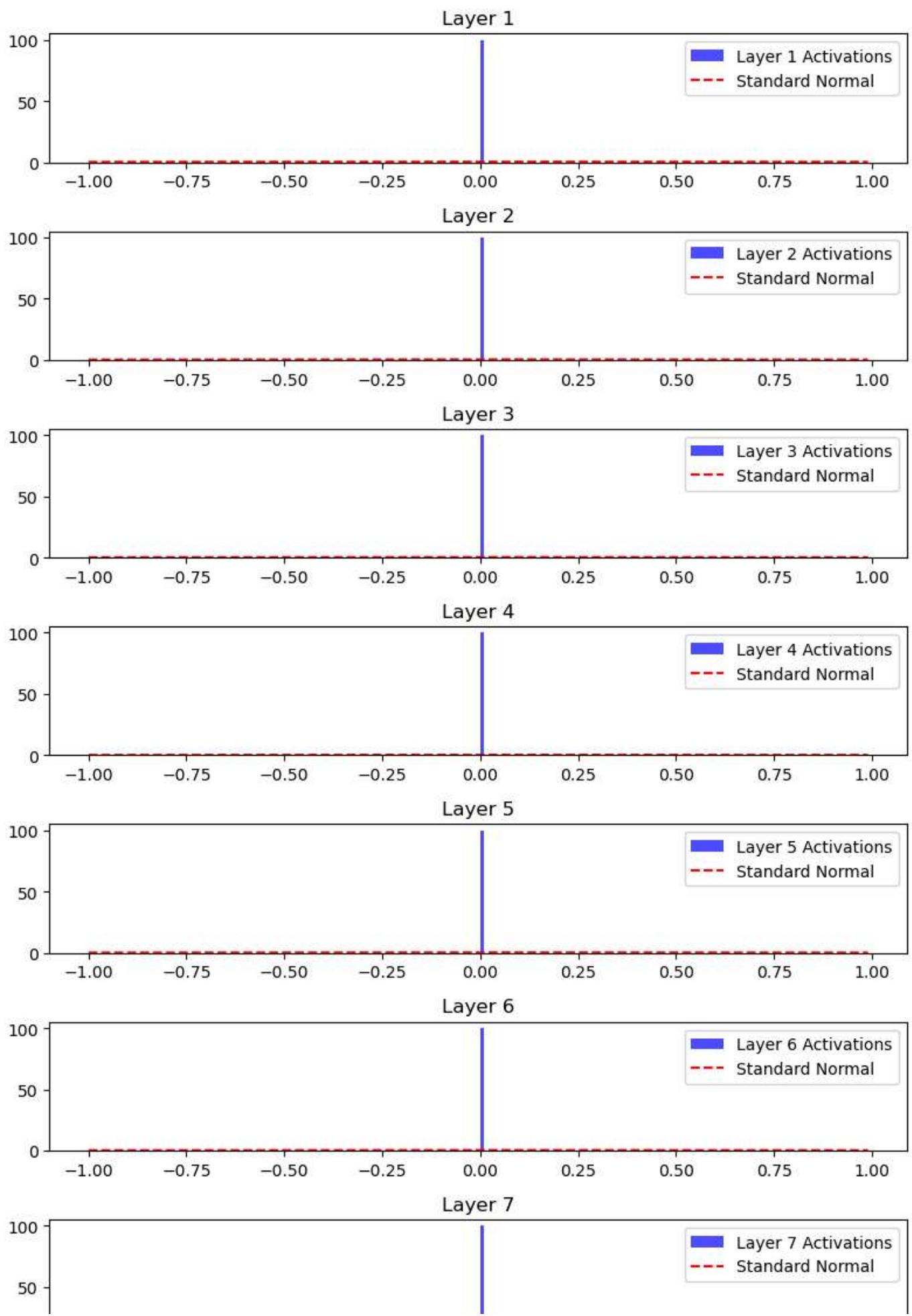
# Forward pass through the network
for i in range(num_layers):
    W = weights[i] if i < num_layers - 1 else None
    layer_input = input_data if i == 0 else activations[i - 1]
    if W is not None:
        layer_output = np.tanh(np.dot(layer_input, W))
    else:
        layer_output = layer_input
    activations.append(layer_output)

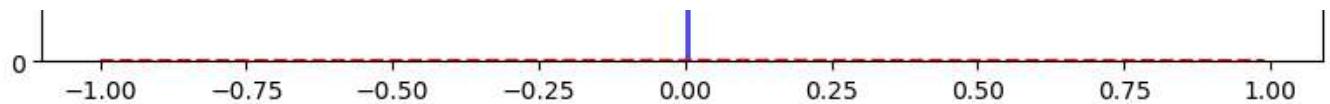
# Create a figure with subplots for each Layer's activation distribution
fig, axes = plt.subplots(num_layers, 1, figsize=(8, 12))

x_axis = np.arange(-1, 1, 0.01)
```

```
for i in range(num_layers):
    ax = axes[i]
    ax.hist(
        activations[i].flatten(),
        bins=100,
        density=True,
        alpha=0.7,
        color="blue",
        label=f"Layer {i+1} Activations",
    )
    ax.plot(
        x_axis,
        norm.pdf(x_axis, 0, 1),
        color="red",
        linestyle="--",
        label="Standard Normal",
    )
    ax.set_title(f"Layer {i+1}")
    ax.legend()

plt.tight_layout()
plt.show()
```





4.A) Inline Question (2 points): Have you noticed any of the drawbacks in the results mentioned earlier? If so, please highlight your observations.

Since all the activations are collapsed to 0, this makes them all update identically, the model might struggle to break this kind of symmetry to learn more complex correlations.

2. Random Weight Initialization

Random Weight Initialization in deep learning sets the initial weights of a neural network to random values, drawn from a distribution. It's simple and encourages diverse starting points for training, breaking symmetry among neurons. However, it can lead to vanishing/exploding gradients in deep networks, making it less effective for them. It's sensitive to initialization values and lacks control compared to specialized methods like Xavier or He initialization.

Task 4.2 Initialize Random weights for each layer (_ points)

The below code visualizes layer-wise activation distributions for a Random Weight Initialization in a neural network. It offers insights into the network's learning and can reveal potential issues.

```
In [34]: # Define the number of Layers and their dimensions
num_layers = 7
layer_dims = [2048] * num_layers

input_data = np.random.randn(16, layer_dims[0])

activations = []
weights = []

for i in range(num_layers - 1):
    #####
    # Initialize random weights for each Layer (except the last one)
    # and the store value in 'W' variable
    #####
    W = np.random.randn(layer_dims[i], layer_dims[i + 1])
    #####
    # END OF YOUR CODE
    #####
    weights.append(W * 0.01)

# Forward pass through the network
for i in range(num_layers):
    W = weights[i] if i < num_layers - 1 else None
    layer_input = input_data if i == 0 else activations[i - 1]
    if W is not None:
        layer_output = np.tanh(np.dot(layer_input, W))
```

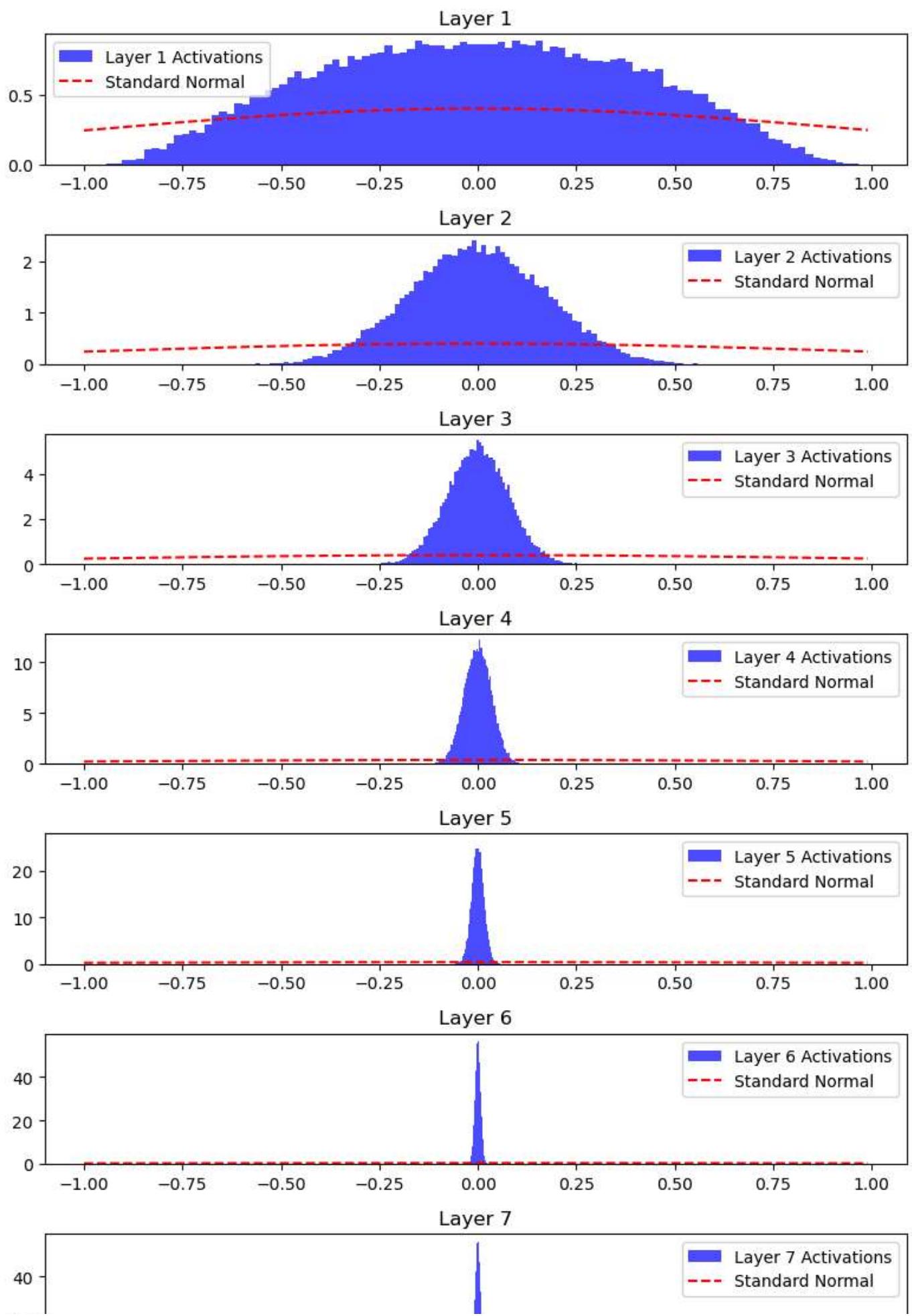
```
    else:
        layer_output = layer_input
        activations.append(layer_output)

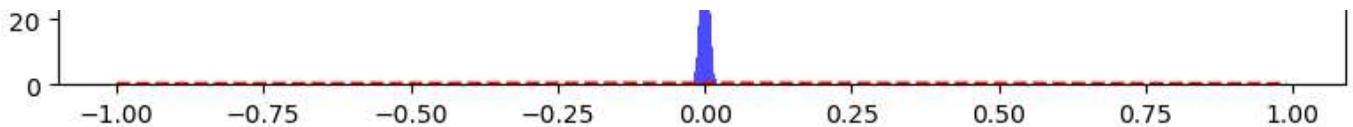
# Create a figure with subplots for each layer's activation distribution
fig, axes = plt.subplots(num_layers, 1, figsize=(8, 12))

x_axis = np.arange(-1, 1, 0.01)

for i in range(num_layers):
    ax = axes[i]
    ax.hist(
        activations[i].flatten(),
        bins=100,
        density=True,
        alpha=0.7,
        color="blue",
        label=f"Layer {i+1} Activations",
    )
    ax.plot(
        x_axis,
        norm.pdf(x_axis, 0, 1),
        color="red",
        linestyle="--",
        label="Standard Normal",
    )
    ax.set_title(f"Layer {i+1}")
    ax.legend()

plt.tight_layout()
plt.show()
```





4.B) Inline Question (2 points): We can see all the activations tend to zero for deeper network layers. what can be expected regarding the gradients dL/dW , and is there still potential for learning?

The gradients will become extremely small, leading to vanishing gradients, essentially decreasing the updates speed to be extremely slow.

I can hear you all say try increasing weights which might resolve the issue. So, let's proceed by multiplying our weights by a factor of 5 and see if it helps.

Task 4.3 - Initialize large Random weights for each layer (_ points)

In [35]:

```

num_layers = 7
layer_dims = [2048] * num_layers

input_data = np.random.randn(16, layer_dims[0])

activations = []

weights = []

for i in range(num_layers - 1):

    ##### Initialize random weights for each Layer similar to the Last
    # task but here you scale them up by a factor of 5
    #####
    W = np.random.randn(layer_dims[i], layer_dims[i + 1]) * 5
    #####
    # END OF YOUR CODE
    #####
    weights.append(W * 0.01)

# Forward pass through the network
for i in range(num_layers):
    W = weights[i] if i < num_layers - 1 else None
    layer_input = input_data if i == 0 else activations[i - 1]
    if W is not None:
        layer_output = np.tanh(np.dot(layer_input, W))
    else:
        layer_output = layer_input
    activations.append(layer_output)

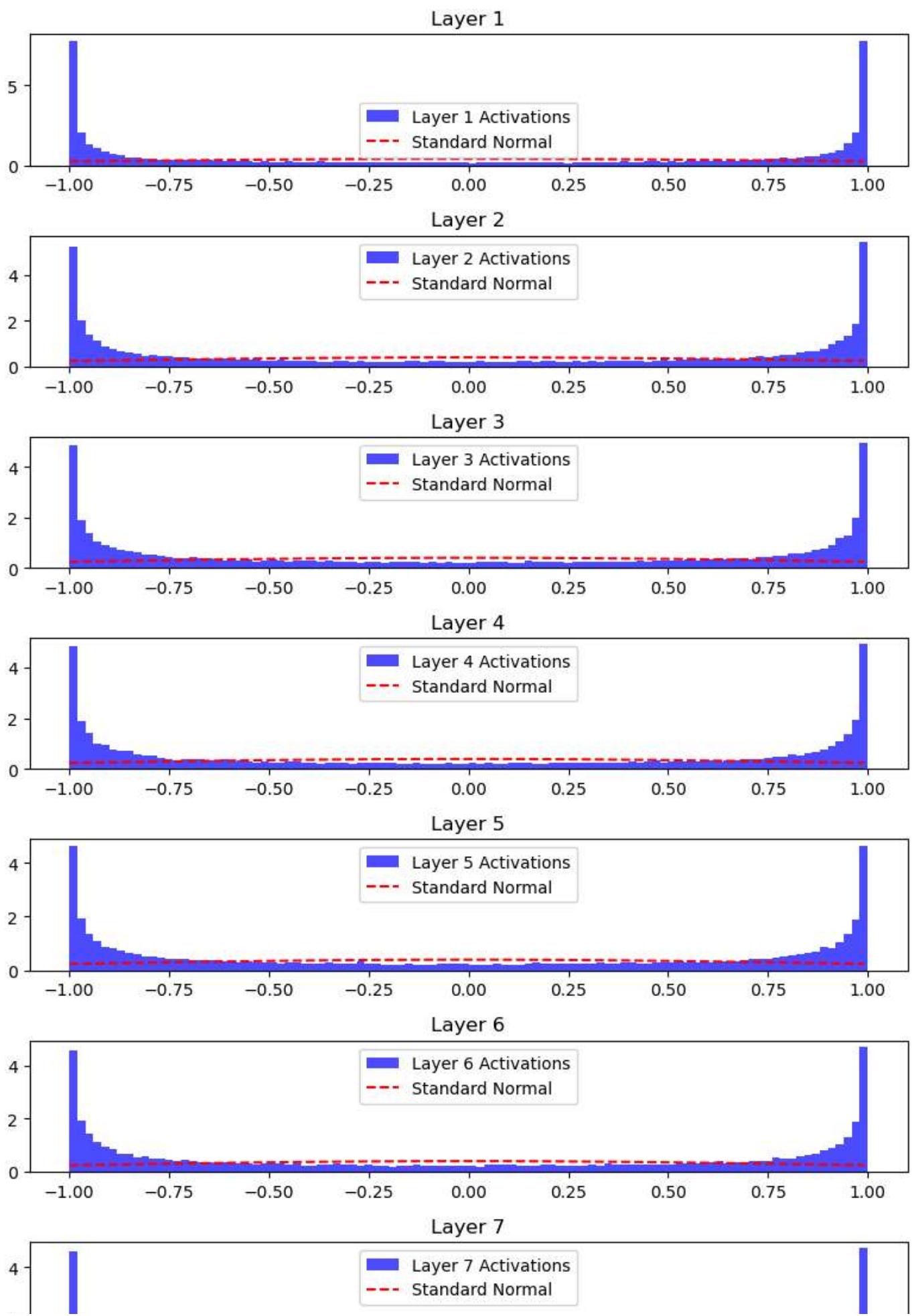
# Create a figure with subplots for each Layer's activation distribution
fig, axes = plt.subplots(num_layers, 1, figsize=(8, 12))

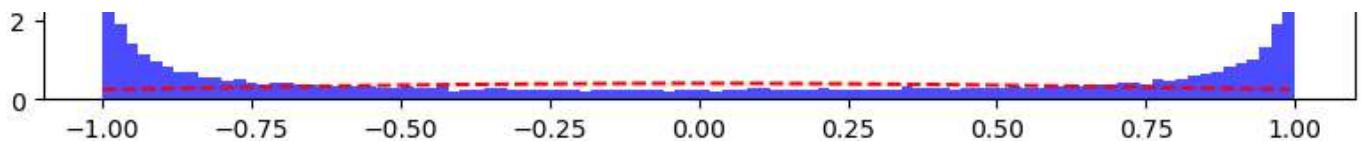
```

```
x_axis = np.arange(-1, 1, 0.01)

for i in range(num_layers):
    ax = axes[i]
    ax.hist(
        activations[i].flatten(),
        bins=100,
        density=True,
        alpha=0.7,
        color="blue",
        label=f"Layer {i+1} Activations",
    )
    ax.plot(
        x_axis,
        norm.pdf(x_axis, 0, 1),
        color="red",
        linestyle="--",
        label="Standard Normal",
    )
    ax.set_title(f"Layer {i+1}")
    ax.legend()

plt.tight_layout()
plt.show()
```





4.C) Inline Question (2 points): All activations saturate due to big weights. What can be expected regarding the gradients dL/dW , and is there still potential for learning?

The saturated activations lead to extremely small gradients, slow down or even stall the whole learning process. It will be almost impossible to learn.

3. Xavier Initialization

Xavier Initialization, also called Glorot Initialization, is a weight initialization method for deep neural networks. It sets initial weights to prevent vanishing and exploding gradients by controlling the variance of activations. This technique stabilizes training and is widely used in practice.

The weights are initialized from a Gaussian distribution with a mean of 0 and a variance of $(1/n_{in})$:-

$$\mathbf{W} = \mathcal{N} \left(0, \frac{1}{n_{in}} \right) \quad (3)$$

The below code visualizes layer-wise activation distributions for a Xavier Weight Initialization in a neural network. It offers insights into the network's learning and can reveal potential issues.

Task 4.4 - Initialize weights using Xavier method for each layer (1 point)

```
In [36]: num_layers = 7
layer_dims = [2048] * num_layers

input_data = np.random.randn(16, layer_dims[0])
activations = []
weights = []

for i in range(num_layers - 1):
    #####
    # Initialize weights using Xavier method for each Layer (except the
    # Last one) and store in variable 'W'
    #####
    W = np.random.randn(layer_dims[i], layer_dims[i + 1]) * np.sqrt(1.0 / layer_dims[i])
    #####
    # END OF YOUR CODE
    #####
    weights.append(W)

# Forward pass through the network
for i in range(num_layers):
    W = weights[i] if i < num_layers - 1 else None
```

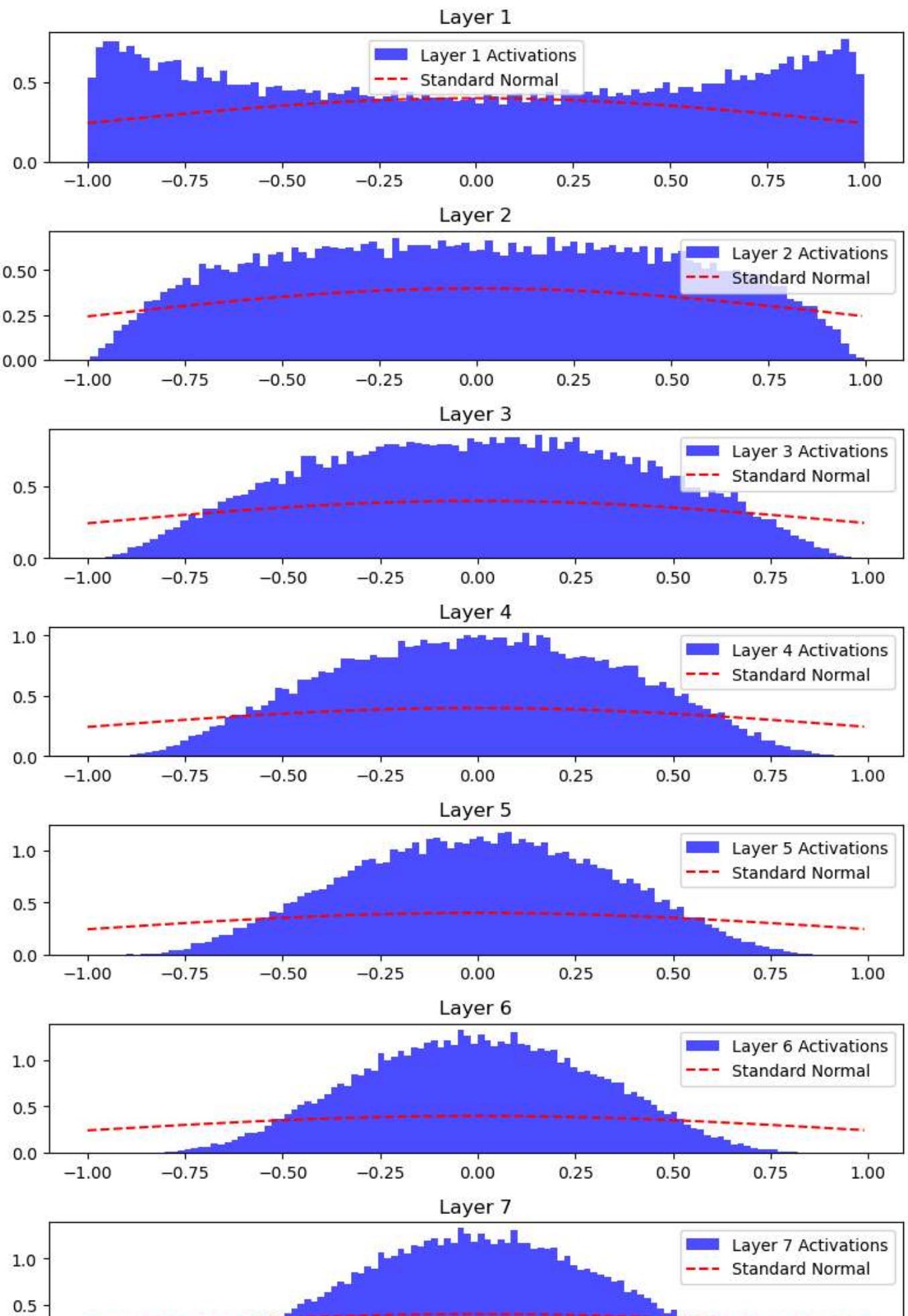
```
layer_input = input_data if i == 0 else activations[i - 1]
if W is not None:
    layer_output = np.tanh(np.dot(layer_input, W))
else:
    layer_output = layer_input
activations.append(layer_output)

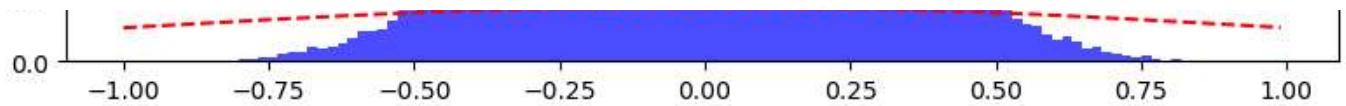
# Create a figure with subplots for each Layer's activation distribution
fig, axes = plt.subplots(num_layers, 1, figsize=(8, 12))

x_axis = np.arange(-1, 1, 0.01)

for i in range(num_layers):
    ax = axes[i]
    ax.hist(
        activations[i].flatten(),
        bins=100,
        density=True,
        alpha=0.7,
        color="blue",
        label=f"Layer {i+1} Activations",
    )
    ax.plot(
        x_axis,
        norm.pdf(x_axis, 0, 1),
        color="red",
        linestyle="--",
        label="Standard Normal",
    )
    ax.set_title(f"Layer {i+1}")
    ax.legend()

plt.tight_layout()
plt.show()
```





4. He/ MSRA Initialization

It is a weight initialization technique commonly used in deep neural networks. It is designed to address the vanishing gradient problem and is particularly effective when Rectified Linear Unit (ReLU) activation functions are used.

For a layer with n_{in} input units, He Initialization initializes the weights by sampling them from a Gaussian distribution with a mean of 0 and a variance of $2 / n_{in}$. The choice of variance (2) is specific to the ReLU activation function and ensures that the weights are set to values that allow ReLU units to activate in a desirable range.

The formula for He Initialization can be expressed as follows:

$$\mathbf{W} = \mathcal{N}\left(0, \frac{2}{n_{in}}\right) \quad (4)$$

The below code visualizes layer-wise activation distributions for a HE/ MSRA Weight Initialization in a neural network. It offers insights into the network's learning and can reveal potential issues.

Task 4.5- Initialize weights using Kaiming He's method for each layer (1 point)

```
In [37]: num_layers = 7
layer_dims = [2048] * num_layers

input_data = np.random.randn(16, layer_dims[0])

activations = []
weights = []

for i in range(num_layers - 1):
    #####
    # Initialize weights using He method for each Layer (except the
    # Last one) and store in variable 'W'
    #####
    W = np.random.randn(layer_dims[i], layer_dims[i + 1]) * np.sqrt(2.0 / layer_dims[i])
    #####
    #           END OF YOUR CODE
    #####
    weights.append(W)

# Forward pass through the network
for i in range(num_layers):
    W = weights[i] if i < num_layers - 1 else None
    layer_input = input_data if i == 0 else activations[i - 1]
    if W is not None:
        layer_output = np.maximum(0, layer_input.dot(W))
```

```
    else:
        layer_output = layer_input
        activations.append(layer_output)

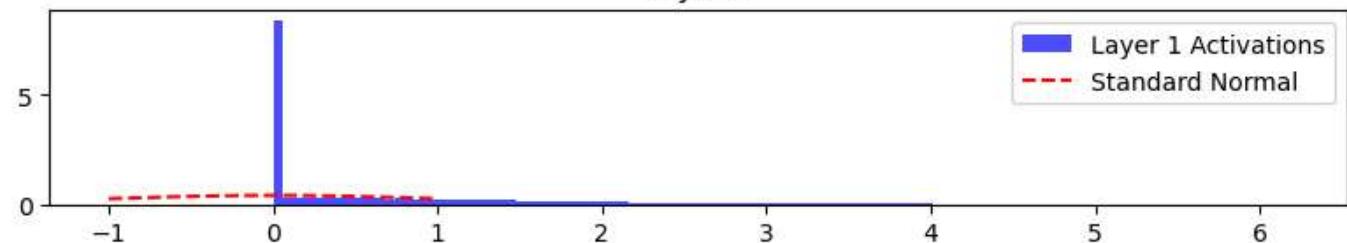
# Create a figure with subplots for each layer's activation distribution
fig, axes = plt.subplots(num_layers, 1, figsize=(8, 12))

x_axis = np.arange(-1, 1, 0.01)

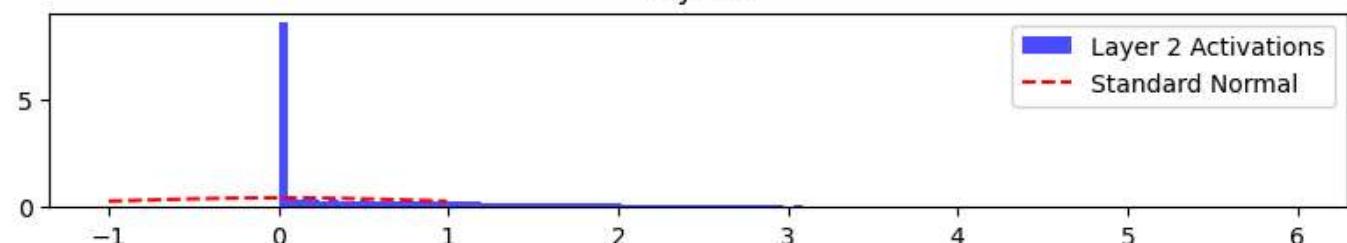
for i in range(num_layers):
    ax = axes[i]
    ax.hist(
        activations[i].flatten(),
        bins=100,
        density=True,
        alpha=0.7,
        color="blue",
        label=f"Layer {i+1} Activations",
    )
    ax.plot(
        x_axis,
        norm.pdf(x_axis, 0, 1),
        color="red",
        linestyle="--",
        label="Standard Normal",
    )
    ax.set_title(f"Layer {i+1}")
    ax.legend()

plt.tight_layout()
plt.show()
```

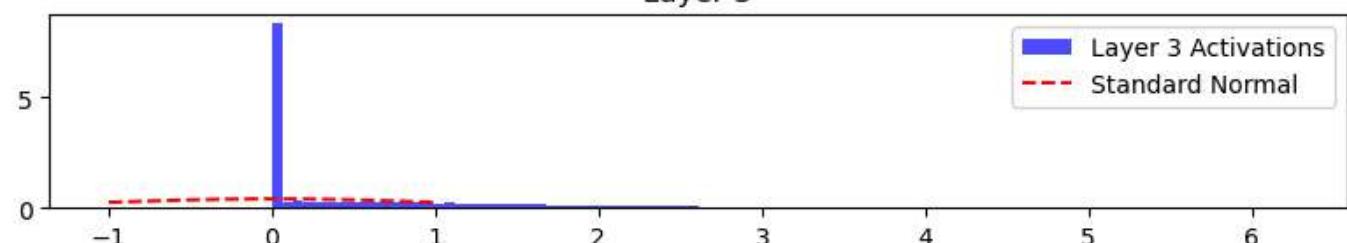
Layer 1



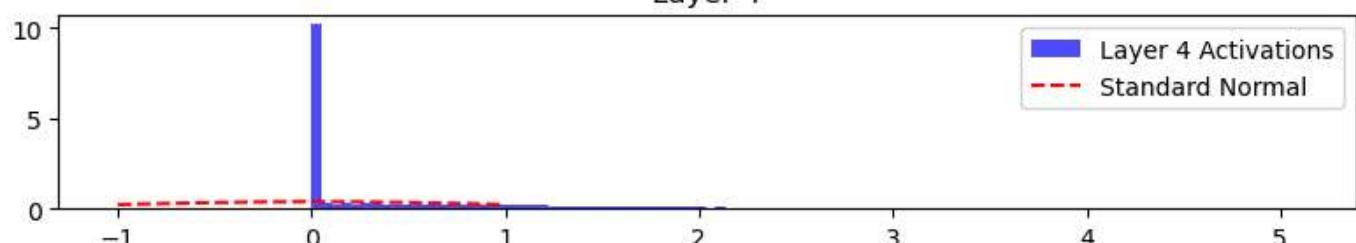
Layer 2



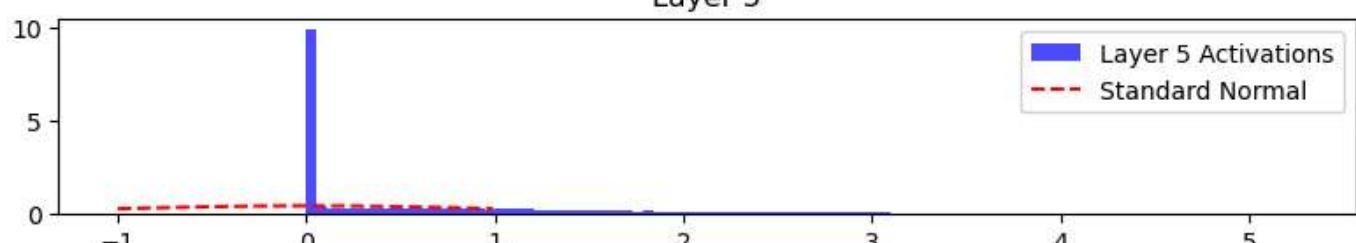
Layer 3



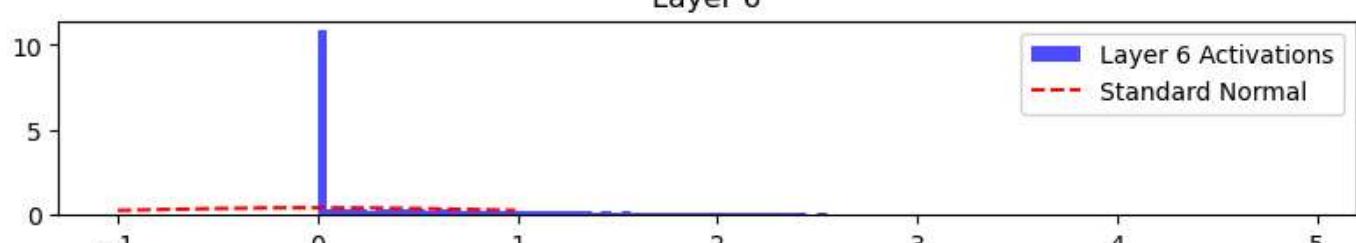
Layer 4



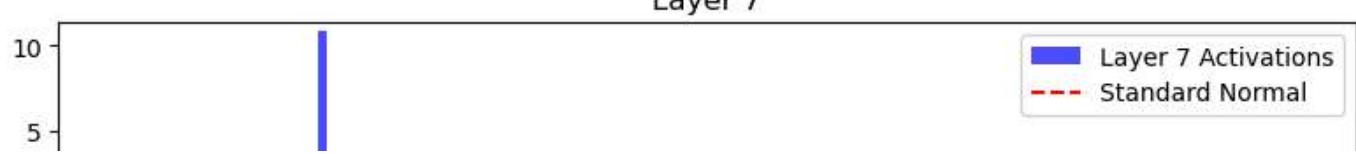
Layer 5

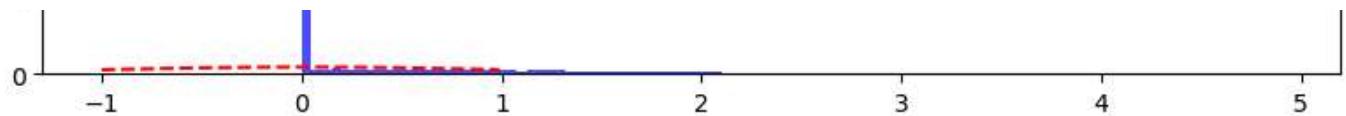


Layer 6



Layer 7





5. Custom Weight Initialization in Pytorch

Custom weight initialization in PyTorch involves setting the weights of a neural network to user-defined values. This can serve different purposes, including enhancing model performance or preventing overfitting.

In PyTorch, you can achieve custom weight initialization using the `init` module, which offers various weight initialization techniques like `normal_`, `uniform_`, and `kaiming_normal_`.

Read Documentation for more information - <https://pytorch.org/docs/stable/nn.init.html>

An illustration of custom weight initialization for a neural network model using values sampled from a normal distribution.

```
class CustomModel(nn.Module):
    def __init__(self):
        super(CustomModel, self).__init__()
        self.fc1 = nn.Linear(64, 32)
        self.fc2 = nn.Linear(32, 10)

        self.initialize_weights()

    #Initialize weights with values from a normal distribution
    def initialize_weights(self):
        nn.init.normal_(self.fc1.weight, mean=0, std=0.01)
        nn.init.normal_(self.fc2.weight, mean=0, std=0.01)

    def forward(self, x):
        pass

model = CustomModel()

for name, param in model.named_parameters():
    if param.requires_grad:
        print(f'Parameter name: {name}')
        print(param.data)
```

Task 4.6 - Create a class called `Supervise_random_weights` that defines weight using Random Weight Initialization method. (1 point)

```
In [38]: class Supervise_random_weights(nn.Module):
    def __init__(self, criterion, net):
        super().__init__()
        self.net = net
        self.criterion = criterion

        for module in self.net.modules():
            if isinstance(module, nn.Linear):
                ##### Set the weights to random values from a normal
                # distribution with a mean of 0 and a standard deviation of 0.01
                ##### nn.init.normal_(module.weight.data, 0, 0.01)
                #### END OF YOUR CODE #####
                #####
            def forward(self, x, y):
                out = self.net(x).squeeze()
                return self.criterion(out, y)
```

Task 4.7 - Build a Neural Network architecture as instructed below using the `Supervise_random_weights` class which defines weight Randomly with a mean =0 and std = 0.01 (2 points)

```
In [39]: torch.manual_seed(7150)
#####
# In this task, we will configure a neural network with Supervise_random_weights
# Hidden Dimension = 32
# Loss - Binary Cross Entropy
# Optimizer - Adam - (lr=0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU)-> (Linear + ReLU)->(Linear + ReLU)->
# (Linear + ReLU)->(Linear + ReLU)->(Linear + Sigmoid)
#####

hidden_dims = 32

run_test(
    Supervise_random_weights(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, hidden_dims),
            nn.ReLU(),
```

```

        nn.Linear(hidden_dims, hidden_dims),
        nn.ReLU(),
        nn.Linear(hidden_dims, output_dims),
        nn.Sigmoid(),
    ),
),
lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)

```

```

#####
#           END OF YOUR CODE
#####

```

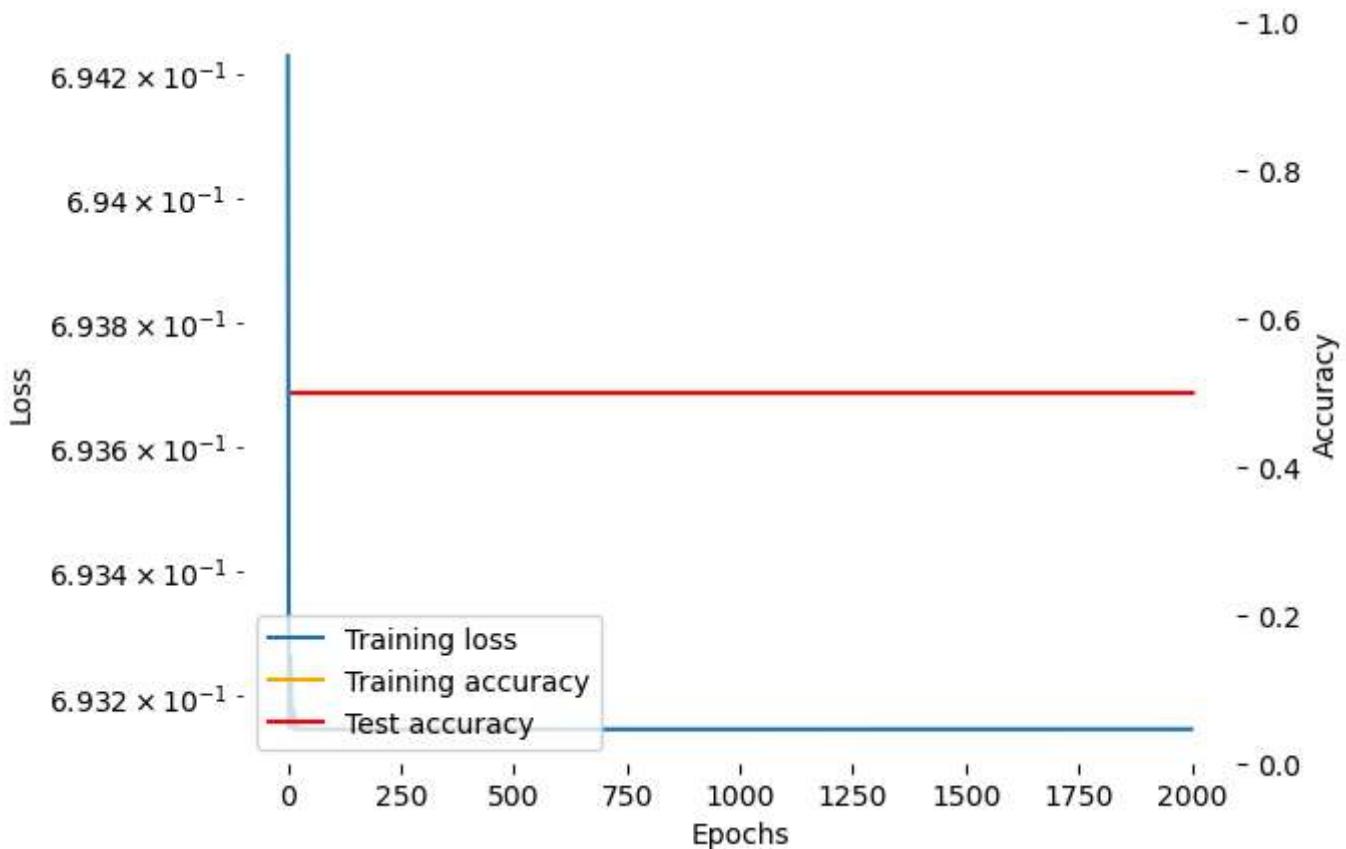
5441 parameters

Epoch 2000, Loss: 0.69315, Grad range 4.7e-10 to 2.6e-30, Train Accuracy: 0.5, Test Accuracy: 0.5

Train Accuracy: 0.50000, Test Accuracy: 0.50000

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```



Task 4.8 - Create a class called `Supervise_Kaiming_weights` that defines weight using Kaiming He's Weight Initialization method. (1 point)

```
In [40]: class Supervise_Kaiming_weights(nn.Module):
    def __init__(self, criterion, net):
        super().__init__()
        self.net = net
        self.criterion = criterion
```

```

for module in self.net.modules():
    if isinstance(module, nn.Linear):
        ##### Set the weights using He initialization #####
        # [Hint: Python's torch.nn.init.kaiming would be useful]
        ##### nn.init.kaiming_normal_(module.weight.data)
        ##### END OF YOUR CODE #####
        #####
    def forward(self, x, y):
        out = self.net(x).squeeze()
        return self.criterion(out, y)

```

Task 4.9 - Build a Neural Network architecture as instructed below using the `Supervise_Kaiming_weights` class (2 points)

```

In [41]: torch.manual_seed(7150)
#####
# In this task, we will configure a neural network with Supervise_He_weights
# Hidden Dimension = 32
# Loss - Binary Cross Entropy
# Optimizer - Adam - (lr=0.01, weight_decay=1e-3)
# Network Architecture - (Linear + ReLU)-> (Linear + ReLU)->(Linear + ReLU)->
# (Linear + ReLU)->(Linear + ReLU)->(Linear + Sigmoid)
#####

hidden_dims = 32

run_test(
    Supervise_Kaiming_weights(
        nn.BCELoss(),
        nn.Sequential(
            nn.Linear(input_size, hidden_dims),
            nn.ReLU(),
            nn.Linear(hidden_dims, output_dims),
            nn.Sigmoid(),
        ),
    ),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)
#####

```

#

END OF YOUR CODE

#

5441 parameters

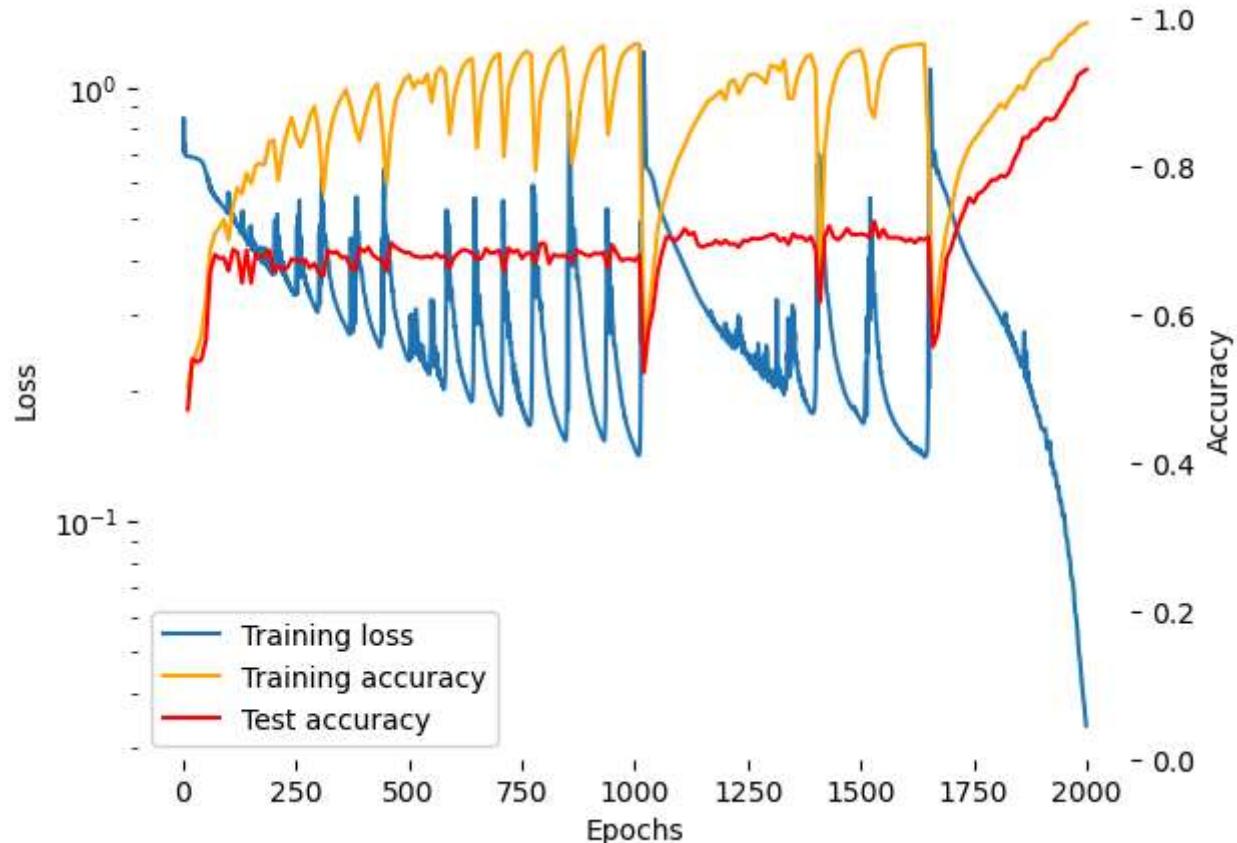
Epoch 2000, Loss: 0.03362, Grad range 1.1e-02 to 1.1e-03, Train Accuracy: 0.9934999942779541, T

est Accuracy: 0.9309999942779541

Train Accuracy: 0.99350, Test Accuracy: 0.93100

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Note - Any idea how does pytorch initialize weights and biases for a layer by default ? Read this discussion - <https://discuss.pytorch.org/t/how-are-layer-weights-and-biases-initialized-by-default/13073>

Extra Credit Questions

Extra Credit Question 1: (2 points)

Now try to train a network to classify the data in the file `hard-classification.npz` . This classification problem is very similar to the original one in `tiny-classification.npz` , with inputs that have a very similar structure. And yet the problem is harder to learn: do the same training techniques work, or is some other approach necessary? Hint: consider transfer learning or fine-tuning approaches

Answer: _____

Pretrain the model

```
In [42]: train_data, train_labels, test_data, test_labels = [
    torch.tensor(m[k]).float()
    for m in [np.load("tiny-classification.npz")]
    for k in "train_data train_labels val_data val_labels".split()
]
```

```
In [49]: torch.manual_seed(7150) # (Leave it here for deterministic behavior and easier grading)

hidden_dims = 16
fan_out_dims = 32

model = nn.Sequential(
    nn.Linear(input_size, hidden_dims),
    ResidualSequence(
        nn.BatchNorm1d(hidden_dims),
        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    nn.BatchNorm1d(hidden_dims),
    nn.Linear(hidden_dims, 1),
    nn.Sigmoid(),
)
```

```
    run_test(  
        Supervise_Kaiming_weights(nn.BCELoss(), model),  
        lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),  
    )
```

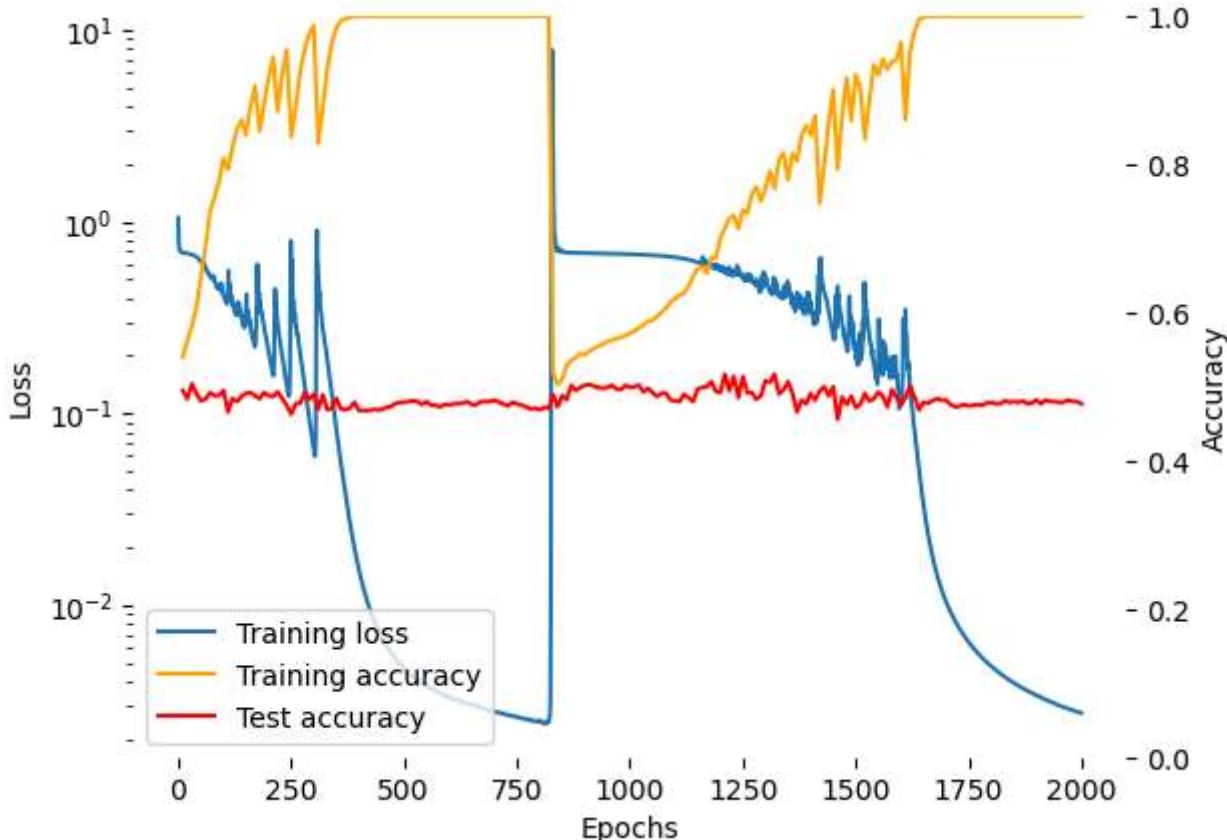
7265 parameters

Epoch 2000, Loss: 0.00274, Grad range 4.0e-03 to 2.7e-08, Train Accuracy: 1.0, Test Accuracy: 0.4769999809265137 9852180481

Train Accuracy: 1.00000, Test Accuracy: 0.47700

C:\Users\GinMa\AppData\Local\Temp\ipykernel_3496\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```



```
In [50]: torch.save(model.state_dict(), "pretrained_model.pth")
```

Load hard dataset & pretrained_model

```
In [51]: train_data, train_labels, test_data, test_labels = [  
    torch.tensor(m[k]).float()  
    for m in [np.load("hard-classification.npz")]  
    for k in "train_data train_labels val_data val_labels".split()  
]  
  
print(  
    f"The training data has {train_data.size(0)} samples, each a vector of  
    {train_data.size(1)} numbers along with"  
)  
print(
```

```

        f'a corresponding set of {train_labels.size(0)} labels, assigning {train_labels.min()}'
    or {train_labels.max()} to each sample."
)

print(
    f'The test data has {test_data.size(0)} samples and labels that are disjoint from the
training data.'
)

```

The training data has 8000 samples, each a vector of 36 numbers along with
a corresponding set of 8000 labels, assigning 0.0 or 1.0 to each sample.
The test data has 1000 samples and labels that are disjoint from the training data.

```
In [52]: pretrained_model = nn.Sequential(
    nn.Linear(input_size, hidden_dims),
    ResidualSequence(
        nn.BatchNorm1d(hidden_dims),
        nn.Linear(hidden_dims, fan_out_dims),
        nn.ReLU(),
        nn.Linear(fan_out_dims, hidden_dims),
    ),
    nn.BatchNorm1d(hidden_dims),
    nn.Linear(hidden_dims, 1),
    nn.Sigmoid(),
)
```

```
In [53]: pretrained_model.load_state_dict(torch.load("pretrained_model.pth"))
pretrained_model.train()
```

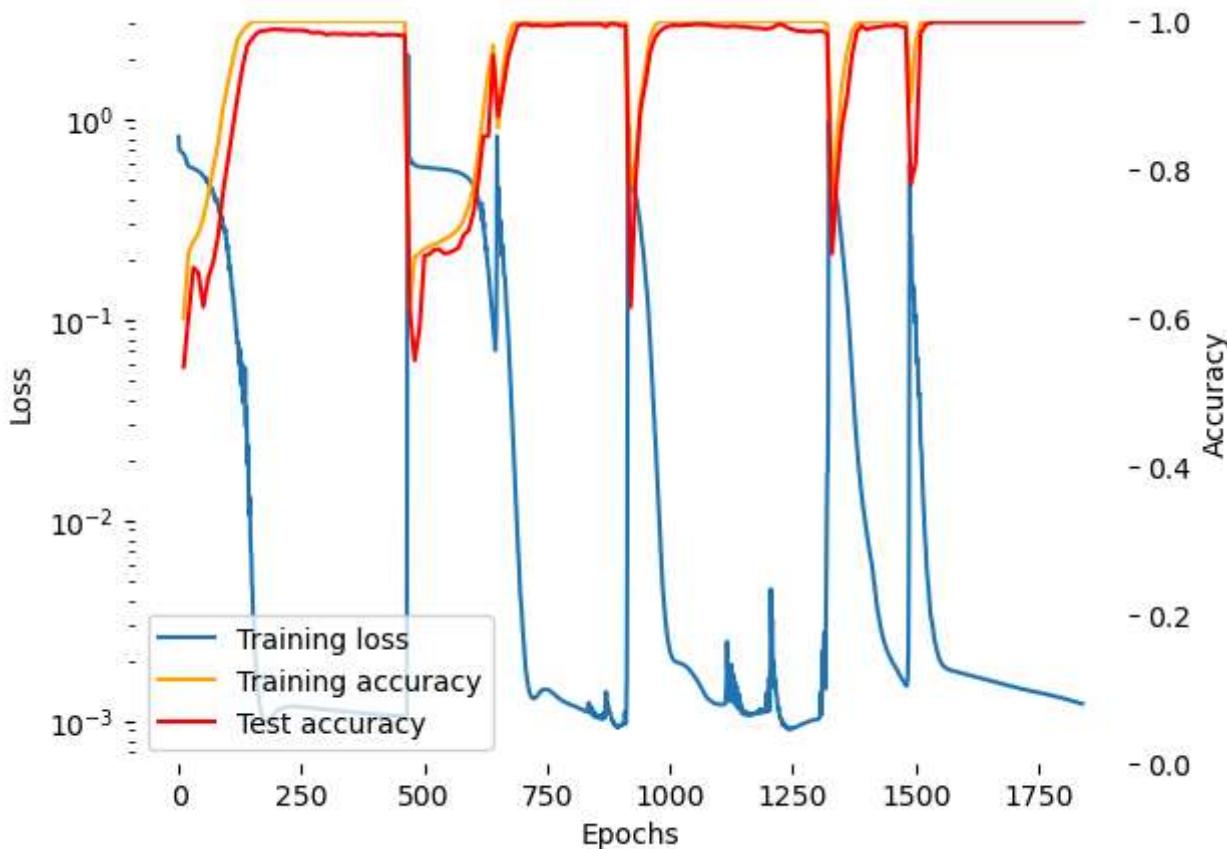
```
Out[53]: Sequential(
    (0): Linear(in_features=36, out_features=16, bias=True)
    (1): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (2): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (3): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (4): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (5): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (6): ResidualSequence(
        (0): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (1): Linear(in_features=16, out_features=32, bias=True)
        (2): ReLU()
        (3): Linear(in_features=32, out_features=16, bias=True)
    )
    (7): BatchNorm1d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (8): Linear(in_features=16, out_features=1, bias=True)
    (9): Sigmoid()
)
```

```
In [85]: run_test(
    Supervise_Kaiming_weights(nn.BCELoss(), pretrained_model),
    lambda p: ADAMOptimizer(p, lr=0.01, weight_decay=1e-3),
)
```

7265 parameters
Epoch 1840, Loss: 0.00122, Grad range 1.8e-03 to 4.0e-10, Train Accuracy: 1.0, Test Accuracy:
1.0 000128746033 000047683716
Train Accuracy: 1.00000, Test Accuracy: 0.99900

```
C:\Users\GinMa\AppData\Local\Temp\ipykernel_28232\2795418505.py:65: UserWarning: Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so can't show the figure.
```

```
fig.show()
```



```
In [54]: torch.save(pretrained_model.state_dict(), "finetuned_model.pth")
```

Extra Credit Question 2: (4 points)

One of the most serious drawbacks of deep networks is that, even if they can learn to solve a problem and recognize patterns in the data, they might not give us humans much insight about those solutions. But if we can create a network that solves a problem, it should be possible to understand that solution. As extra credit, figure out: what classification rule did the neural network learn in the above exercises when the network achieves 100% hold-out accuracy? Can you extract from the network a succinct set of rules that it implements, for example, can you decompile the network into a short python program, that can correctly assign a class to a sample?

Answer:

I didn't find a way to understand the model. I had experiments with local explanation methods (incl. LIME, Gradient-based Attribution, and SHAP), but the explanations didn't really help to understand the model on a global scale. Distillation with decision tree and logistics regression also didn't seem to work since it's a binary classification problem.

In my opinion it will be really hard to understand what the model had learned during training without mechanistically understanding each component of the model.

Reference

- Shap/shap. (2025). [Jupyter Notebook]. shap. <https://github.com/shap/shap> (Original work published 2016)
- ChatGPT. (n.d.). Retrieved January 17, 2025, from <https://chatgpt.com>
- DeepSeek. (n.d.). Retrieved February 3, 2025, from <https://chat.deepseek.com>
- Ribeiro, M. T. C. (2025). Marcotcr/lime [JavaScript]. <https://github.com/marcotcr/lime> (Original work published 2016)

```
In [ ]: from sklearn.tree import DecisionTreeClassifier, export_text
from sklearn.metrics import accuracy_score

# Load the tiny (simple) dataset
tiny_data = np.load("tiny-classification.npz")
tiny_train_data = torch.tensor(tiny_data["train_data"]).float()
tiny_train_labels = torch.tensor(
    tiny_data["train_labels"])
).float() # not used in distillation
tiny_test_data = torch.tensor(tiny_data["val_data"]).float()
tiny_test_labels = torch.tensor(
    tiny_data["val_labels"])
).float() # not used in distillation

# Load the hard (challenging) dataset
hard_data = np.load("hard-classification.npz")
hard_train_data = torch.tensor(hard_data["train_data"]).float()
hard_train_labels = torch.tensor(
    hard_data["train_labels"])
).float() # not used in distillation
hard_test_data = torch.tensor(hard_data["val_data"]).float()
hard_test_labels = torch.tensor(
    hard_data["val_labels"])
).float() # not used in distillation

#####
# 4. Generate Teacher Predictions for Both Training Datasets #
#####

# For the tiny dataset
with torch.no_grad():
    teacher_preds_tiny = pretrained_model(tiny_train_data)
teacher_preds_tiny = teacher_preds_tiny.numpy().ravel()
# Threshold at 0.5 to produce binary labels
teacher_labels_tiny = (teacher_preds_tiny > 0.5).astype(int)

# For the hard dataset
with torch.no_grad():
    teacher_preds_hard = pretrained_model(hard_train_data)
teacher_preds_hard = teacher_preds_hard.numpy().ravel()
teacher_labels_hard = (teacher_preds_hard > 0.5).astype(int)

#####
```

```

# 5. Combine the Two Datasets Using Sample Weights    #
#####
# Convert input tensors to NumPy arrays
X_tiny = tiny_train_data.numpy()
X_hard = hard_train_data.numpy()

# Concatenate the inputs and teacher-produced Labels
X_combined = np.concatenate([X_tiny, X_hard], axis=0)
y_combined = np.concatenate([teacher_labels_tiny, teacher_labels_hard], axis=0)

# To simulate the two-stage (pre-training then fine-tuning) approach,
# we assign a lower weight to the tiny dataset and a higher weight to the hard dataset.
# (These weight factors can be tuned as needed.)
weights_tiny = np.ones(len(teacher_labels_tiny)) * 1000.0
weights_hard = np.ones(len(teacher_labels_hard)) * 0.000001
sample_weights = np.concatenate([weights_tiny, weights_hard])

#####
# 6. Train the Decision Tree on the Combined Data  #
#####
tree_clf = DecisionTreeClassifier(max_depth=5, random_state=42)
tree_clf.fit(X_combined, y_combined, sample_weight=sample_weights)

#####
# 7. Evaluate the Distilled Decision Tree on Hard Test Set
#####
X_hard_test = hard_test_data.numpy()
with torch.no_grad():
    teacher_preds_hard_test = pretrained_model(hard_test_data)
teacher_labels_hard_test = (teacher_preds_hard_test.numpy().ravel() > 0.5).astype(int)

# Get the decision tree predictions on the hard test set
tree_preds_hard_test = tree_clf.predict(X_hard_test)

# Compare the tree's predictions with the teacher's outputs
accuracy = accuracy_score(hard_test_labels, tree_preds_hard_test)
print(
    "Decision Tree (distilled) vs. Teacher Accuracy on Hard Test Set: {:.2f}%".format(
        accuracy * 100
    )
)

#####

# 8. Export the Decision Tree Rules
#####
rules = export_text(
    tree_clf, feature_names=[f"f{i}" for i in range(X_combined.shape[1])]
)
print("\nExtracted Decision Tree Rules:")
print(rules)

```

Decision Tree (distilled) Accuracy on Hard Test Set: 47.60%

Extracted Decision Tree Rules:

```
|--- f21 <= 0.50
|   |--- f7 <= 0.50
|   |   |--- f17 <= 0.50
|   |   |   |--- f16 <= 0.50
|   |   |   |   |--- f12 <= 0.50
|   |   |   |   |   |--- class: 1
|   |   |   |   |--- f12 >  0.50
|   |   |   |   |   |--- class: 1
|   |   |   |--- f16 >  0.50
|   |   |   |   |--- f26 <= 0.50
|   |   |   |   |   |--- class: 1
|   |   |   |   |--- f26 >  0.50
|   |   |   |   |   |--- class: 0
|   |--- f17 >  0.50
|   |   |--- f24 <= 0.50
|   |   |   |--- f18 <= 0.50
|   |   |   |   |--- class: 0
|   |   |   |--- f18 >  0.50
|   |   |   |   |--- class: 0
|   |--- f24 >  0.50
|   |   |--- f23 <= 0.50
|   |   |   |--- class: 1
|   |   |--- f23 >  0.50
|   |   |   |--- class: 1
|--- f7 >  0.50
|   |--- f22 <= 0.50
|   |   |--- f8 <= 0.50
|   |   |   |--- f12 <= 0.50
|   |   |   |   |--- class: 0
|   |   |   |--- f12 >  0.50
|   |   |   |   |--- class: 1
|   |--- f8 >  0.50
|   |   |--- f25 <= 0.50
|   |   |   |--- class: 1
|   |   |--- f25 >  0.50
|   |   |   |--- class: 1
|--- f22 >  0.50
|   |--- f1 <= 0.50
|   |   |--- f0 <= 0.50
|   |   |   |--- class: 0
|   |   |--- f0 >  0.50
|   |   |   |--- class: 1
|   |--- f1 >  0.50
|   |   |--- f20 <= 0.50
|   |   |   |--- class: 0
|   |   |--- f20 >  0.50
|   |   |   |--- class: 0
|--- f21 >  0.50
|   |--- f20 <= 0.50
|   |   |--- f32 <= 0.50
|   |   |   |--- f19 <= 0.50
|   |   |   |   |--- f17 <= 0.50
|   |   |   |   |   |--- class: 1
|   |   |   |   |--- f17 >  0.50
```

```
| | | | |--- class: 0
| | | | --- f19 >  0.50
| | | | |--- f2 <= 0.50
| | | | |--- class: 0
| | | | |--- f2 >  0.50
| | | | |--- class: 1
| | | --- f32 >  0.50
| | | |--- f4 <= 0.50
| | | | |--- f34 <= 0.50
| | | | |--- class: 0
| | | | |--- f34 >  0.50
| | | | |--- class: 1
| | | | --- f4 >  0.50
| | | | |--- f23 <= 0.50
| | | | |--- class: 0
| | | | |--- f23 >  0.50
| | | | |--- class: 0
| | | --- f20 >  0.50
| | | |--- f30 <= 0.50
| | | | |--- f33 <= 0.50
| | | | |--- f16 <= 0.50
| | | | |--- class: 0
| | | | |--- f16 >  0.50
| | | | |--- class: 0
| | | | --- f33 >  0.50
| | | | |--- f18 <= 0.50
| | | | |--- class: 0
| | | | |--- f18 >  0.50
| | | | |--- class: 0
| | | | --- f30 >  0.50
| | | | |--- f14 <= 0.50
| | | | | |--- f3 <= 0.50
| | | | | |--- class: 0
| | | | | |--- f3 >  0.50
| | | | | |--- class: 0
| | | | |--- f14 >  0.50
| | | | | |--- f12 <= 0.50
| | | | | |--- class: 0
| | | | | |--- f12 >  0.50
| | | | | |--- class: 1
```

```
In [121... counts = hard_train_data.sum(dim=1)
tree_clf.fit(counts.reshape(-1, 1), hard_train_labels)
```

```
Out[121... ▾ DecisionTreeClassifier ⓘ ??
DecisionTreeClassifier(max_depth=5, random_state=42)
```

```
In [124... # Compare the tree's predictions with the teacher's outputs
accuracy = accuracy_score(hard_test_labels, tree_preds_hard_test)
print(
    "Decision Tree (distilled) vs. Teacher Accuracy on Hard Test Set: {:.2f}%".format(
        accuracy * 100
```

```
)
```

Decision Tree (distilled) vs. Teacher Accuracy on Hard Test Set: 50.30%

In [162...]

```
import torch
from captum.attr import IntegratedGradients
import matplotlib.pyplot as plt

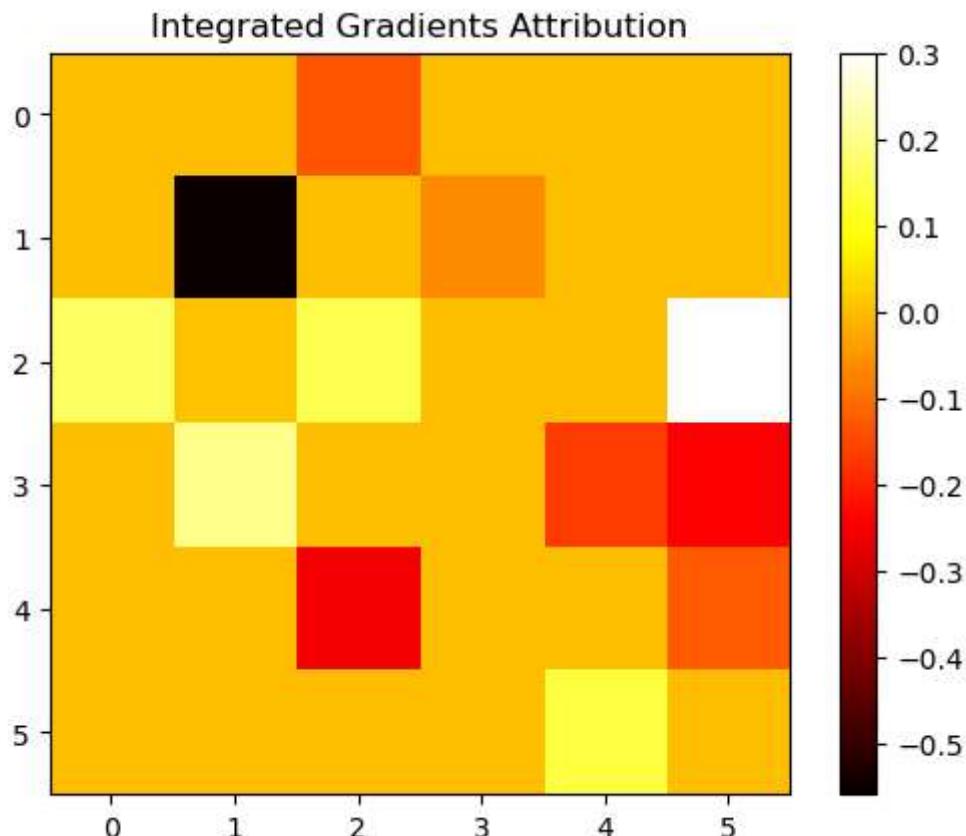
test_sample = hard_test_data[3]

# Assuming you have a model and an input tensor:
pretrained_model.eval()
input_tensor = test_sample.unsqueeze(0) # Add batch dimension
baseline = torch.zeros_like(input_tensor)

# Create an instance of IntegratedGradients
ig = IntegratedGradients(pretrained_model)

# Compute attributions
attributions, delta = ig.attribute(
    input_tensor, baseline, target=0, return_convergence_delta=True
)
attributions = attributions.reshape(6, 6)

# Visualize the attributions (for image data, for example)
plt.imshow(attributions.squeeze().cpu().detach().numpy(), cmap="hot")
plt.colorbar()
plt.title("Integrated Gradients Attribution")
plt.show()
```



In [157...]

```
import shap

# Select background data for SHAP (use first 100 training samples)
background = train_data[:100]

# Initialize SHAP DeepExplainer
explainer = shap.DeepExplainer(model=pretrained_model, data=background)

# Calculate SHAP values for test samples (using first 20 test examples)
# test_samples = test_data[:20]
test_samples = test_data
shap_values = explainer.shap_values(test_samples)

# shap.image_plot(shap_values, test_samples.detach().numpy())
```



In [136...]

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split

# Step 1: Compute average absolute SHAP values per feature.
avg_shap = np.mean(np.abs(shap_values), axis=0) # shape: (num_features,)
important_feature_indices = np.argsort(avg_shap)[-10:] # select top 10 features

# Step 2: Prepare a reduced dataset.
X_reduced = hard_train_data[:, important_feature_indices].squeeze()

simple_model = LogisticRegression()
simple_model.fit(X_reduced, hard_train_labels)
```

Out[136...]

```
▼ LogisticRegression ⓘ ?  
LogisticRegression()
```

In [153...]

```
with torch.no_grad():
    # Forward pass: compute the predicted outputs by passing hard_test_data to the model
    outputs = simple_model.predict(
        hard_test_data[:, important_feature_indices].squeeze()
    )

    # Since the model ends with a Sigmoid, outputs are probabilities in the range [0, 1].
```

```
# Threshold at 0.5 to obtain predicted classes (0 or 1).
predicted_classes = outputs

# If the outputs have shape (N, 1), we can squeeze them to compare with labels.
predicted_classes = predicted_classes.squeeze()

# Compare predictions to the true Labels
correct_predictions = predicted_classes == hard_test_labels.numpy()

# Compute accuracy as the mean of correct predictions
accuracy = correct_predictions.mean()

print(accuracy)
```

```
0.512
```