

# Report for SDM274 Assgn05 MLP

---

Author: Ziyuan Li, 12211225

This document is created as a report for SDM274, Sustech.

## Background

---

Multiple Layer Perceptron(MLP) is a type of neural network, primarily used in supervised learning. It's made up of multiple layers of nodes, or neurons, organized into an input layer, one or more hidden layers, and an output layer. Each neuron is connected to every neuron in the previous and following layer (fully connected), and each connection has an associated weight, which is learned during training.

It was initially invented in the 1980s, aimed at addressing problems where linear models were insufficient. By introducing non-linear activation functions and organizing neurons into multiple fully connected layers, MLPs allowed for more sophisticated feature extraction and data representation. As deep learning evolved, MLPs were largely surpassed by more specialized architectures like CNN for image and spatial data, and RNNs for sequential tasks such as language processing. Despite this, MLPs are still widely used in scenarios where simplicity and computational efficiency are prioritized, or as subcomponents in complex architectures.

Furthermore, the concept of the fully connected layer remains fundamental across various modern deep learning models, including CNNs, RNNs, and transformers, where they play a crucial role in transforming high-dimensional data, interpreting features, and synthesizing final outputs, illustrating their enduring importance in deep learning.

## Code Design

---

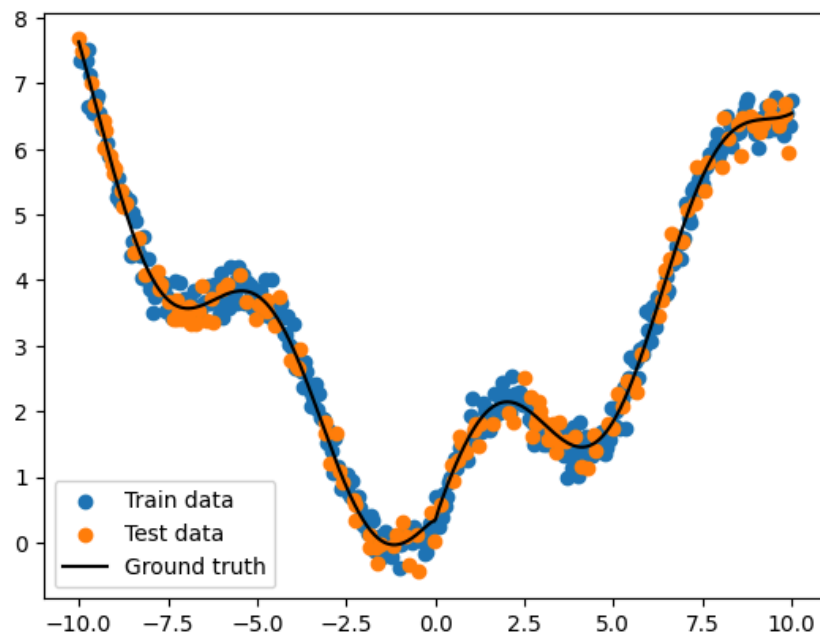
### Designing the nonlinear function

---

```
1 def complex_nonlin_func(x):  
2     return (  
3         np.sin(x) +  
4         0.3 * np.exp( np.abs(x)**0.5 )  
5     )
```

Here I introduced a nonlinear function composed of a sine and an exponential part. It is nonlinear, and not symmetrical, which is hard to approximate using a linear model.

I generated 500 points on the line, spanning from -10 to 10, and also split the data into (0.7, 0.3) for training and testing.



```
1 y = y + 0.2 * np.random.normal(size=x.size)
```

I added a random noise with the amplitude of 0.2 to the signal.

## Modules.py resembling torch.nn

I designed a family of class that **resembles the Pytorch Torch.nn library**. This library is well designed for its modularity, ease of use, and flexibility. User can design a network by simply declaring and stacking layers in a list. An **example code snippet of torch.nn code** is provided below:

```
1 net = nn.Sequential()
2 net.add(nn.Dense(256, activation='relu'),
3         nn.Dense(10))
4 net.initialize(init.Normal(sigma=0.01))
```

## modularity

In my designed class, I implement the modularity with **a base class Module**. The `Module` class contains shared attributes and methods, like `parameters` (to manage learnable parameters) and caches (`output_cache`, `input_cache`) for forward and backward passes.

```
1 class Module:
2
3     def __init__(self):
4         self.parameters = []
5         self.output_cache = None
6         self.input_cache = None
7         self.isInference = False
8
9     def forward(self):
10        raise NotImplementedError
11
12    def backward(self):
13        raise NotImplementedError
14
15    def parameters(self):
16        return self.parameters
```

Adding new types of layers or operations is straightforward—I can inherit from `Module` and override the `forward` and `backward` methods as needed.

## caching intermediate values

Implemented layers **cache intermediate values** needed for backpropagation.

For example, `Linear` caches `input_cache` to store inputs during the forward pass. This allows the layer to reuse these inputs during the backward pass for gradient calculations.

## forward and backward methods

Each module implements its own **forward and backward methods**, essential for differentiable operations in neural networks.

- In `forward`, each layer computes its output based on the input, caching intermediate results if necessary.
- In `backward`, the layers calculate the gradients with respect to their inputs and parameters. For example, `Linear` calculates gradients for weights (`dw`), biases (`db`), and the input (`dx`) to propagate gradients to previous layers.

## How parameters are stored

Layers with learnable parameters, such as `Linear`, define and manage their parameters (weights `w` and biases `b`). These parameters are stored in a dictionary (`self.parameters`).

However this design may lead to unwanted re-train if lack proper management

## managing parameters

- The library introduces an `isInference` flag to distinguish between **training and inference modes**. When `isInference` is set to `True`, the forward pass bypasses caching and specific training-only operations.
- The `reset` method in `Linear` layers reinitializes the weights and biases. This can be useful when experimenting with different initializations and hyperparameters.

*The source code is provided in the appendix*

# Nonlinear Regression with MLP

## Design of the Network

I defined a MLP class which inherits from the Module class (which might not be the best choice but it works).

MLP
-layers -epochs -lr -loss -input_shape -ouatput_shape
+__init__(layers, epochs=1000, lr=0.01, input_shape=1, output_shape=1) +forward(X) +backward(dA) +inference(X) +update_params(lr) +predict(X) +reset() +train_BGD(X, Y, epochs=None, lr=None)

## Initialization

First, initialize the `MLP` class with the desired layers, number of epochs, learning rate, input shape, and output shape. Here is the structure used for training.

```

1  layers = [
2      Linear(input_size=1, output_size=64),
3      ReLU(),
4      Linear(input_size=64, output_size=128),
5      ReLU(),
6      Linear(input_size=128, output_size=256),
7      ReLU(),
8      Linear(input_size=256, output_size=1)
9  ]
10
11 mlp = MLP(layers=layers, epochs=1000, lr=0.01, input_shape=1, output_shape=1)

```

To train the model, call the `train_BGD` method with the input data `X` and target data `Y`.

```

1  x = np.array([...]) # Input data
2  y = np.array([...]) # Target data
3
4  mlp.train_BGD(X, Y)

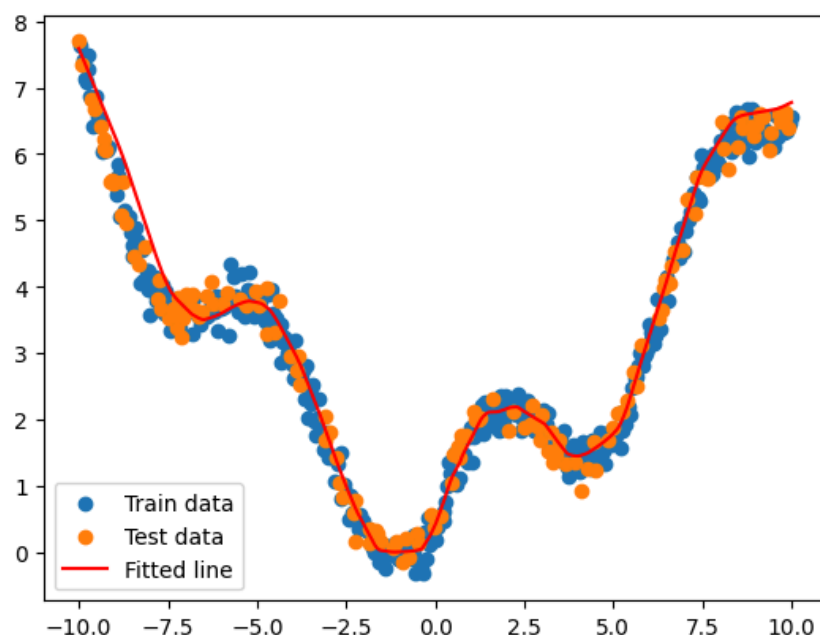
```

# Training Steps

1. **Forward Pass:** The `forward` method is called to compute the output of the model for the given input  $X$ .
2. **Loss Calculation:** The loss is computed using Mean Squared Error (MSE) between the predicted output and the target  $Y$ .
3. **Backward Pass:** The `backward` method is called to compute the gradients of the loss with respect to the model parameters.
4. **Parameter Update:** The `update_params` method updates the model parameters using the computed gradients and the learning rate.
5. **Epoch Loop:** Steps 1-4 are repeated for the specified number of epochs.

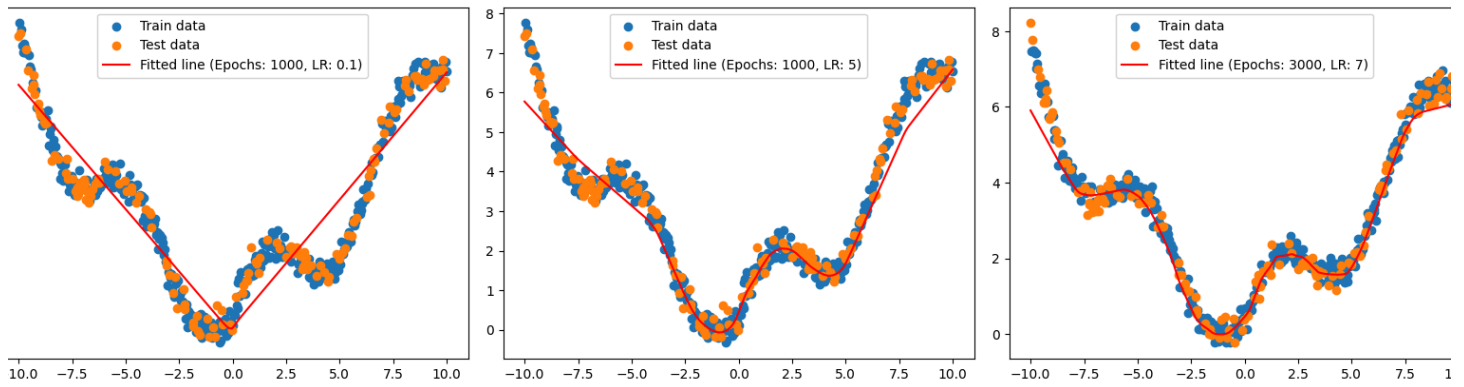
*Code provided in the appendix.*

## Fitting



```
1 # Creates the model
2 mlp_model = MLP(layers, epochs=5000, lr=7)
3 if not trained:
4     # Train the model with the generated data
5     mlp_model.train_BGD(X_train, y_train)
6     trained = True
7 elif trained:
8     mlp_model.reset()
9     mlp_model.train_BGD(X_train, y_train)
```

This image is the result of successful fitting. For the initial version of the Multi-Layer Perceptron (MLP), which lacked any form of regularization, I had to set hyperparameters to rather extreme values in order to achieve satisfactory training performance. This process significantly increased the training time, requiring approximately 2 minutes to complete, which is notably longer compared to the simpler single-neuron models used previously.



Some comparisons. The MLPs need significantly more resources to train.

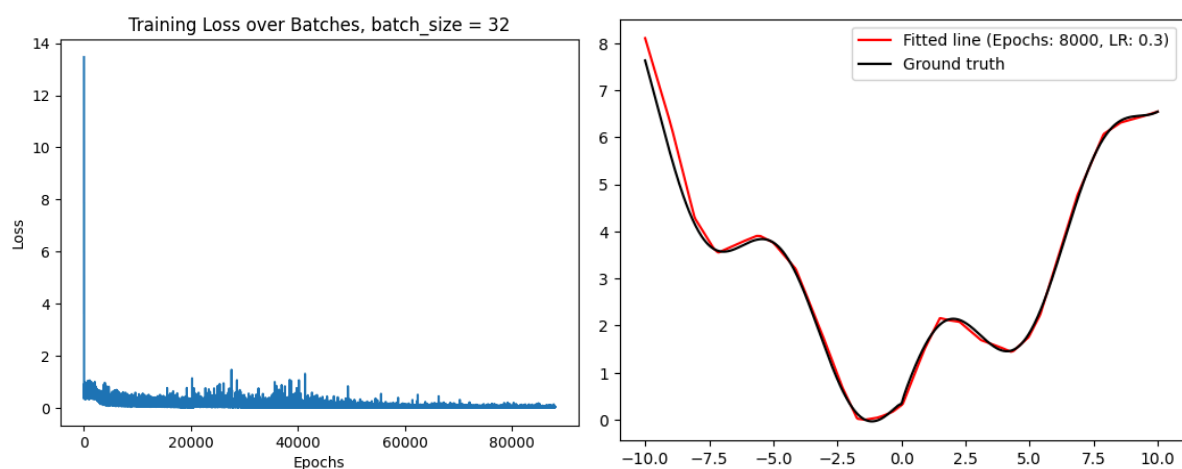
## Optimization of structure and training strategy

Empirically, I discovered that a shallower layer can achieve approximately the same result, with lesser demand on training resources.

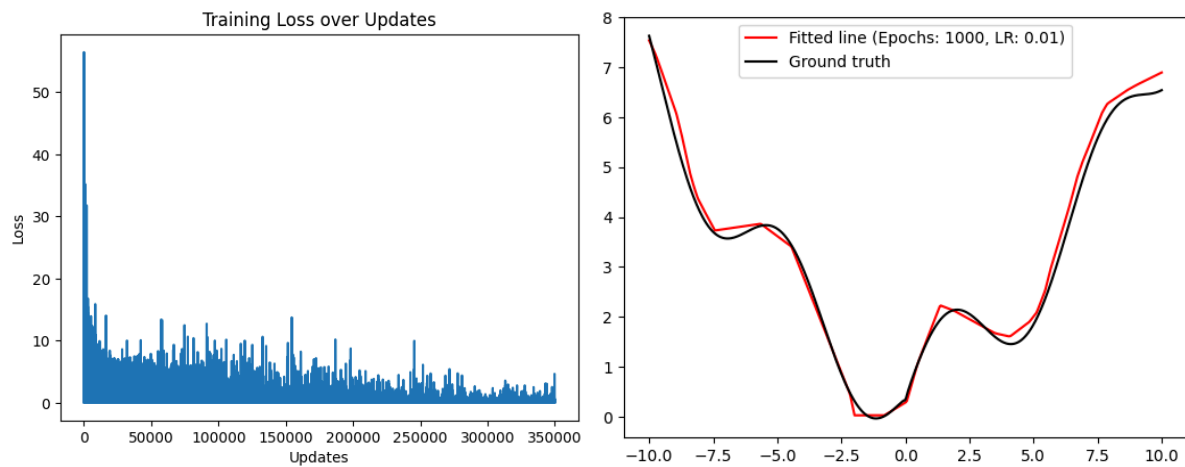
```
1 layers = [
2     Linear(input_size=1, output_size=64),
3     ReLU(),
4     Linear(input_size=64, output_size=128),
5     ReLU(),
6     Linear(input_size=128, output_size=1)
7 ]
```

On my system with a AMD Ryzen 7 7840HS@25W, it takes only 9.7 seconds to train, and achieved fairly good results on the regression task.

## Results Implemented by the Mini-Batch Gradient Descent (MBGD) Method



## Results Implemented by Stochastic Gradient Descent (SGD) Method



Empirically discovered, MBGD with  $bs=32$ ,  $epochs=8000$ ,  $lr = 0.3$  is a good set of hyperparameters. We will do the cross validation part on it.

## Cross Validation

I implemented a k-fold cross validation method to evaluate the overall model.

```
1 def train_k_fold(self, X, Y, k=5, epochs=None, lr=None):
2     if epochs is None:
3         epochs = self.epochs
4     if lr is None:
5         lr = self.lr
6
7     m = X.shape[0]
8     fold_size = m // k # integer division, i.e. 103 // 5 = 20
9     indices = np.arange(m) # 0, 1, 2, ..., m-1
10    np.random.shuffle(indices) # shuffle the indices, train randomly
11
12    fold = 1
13    for i in range(k): # iterate over the folds, totally we train for k times
14        val_start = i * fold_size
15        val_end = val_start + fold_size if i < k - 1 else m # if it's the last
16        fold, use the remaining data
17
18        val_indices = indices[val_start:val_end] # takes the validation
19        indices
20        train_indices = np.concatenate([indices[:val_start],
21        indices[val_end:]]) # takes all the indices except the validation indices
22
23        X_train = X[train_indices]
24        Y_train = Y[train_indices]
25        X_val = X[val_indices]
26        Y_val = Y[val_indices]
27
28        self.reset()
29        self.train_MBGD(X_train, Y_train, epochs=epochs, lr=lr) # we still use
30        the mini-batch gradient descent, with a default bs of 32.
31
32        A_val = self.inference(X_val)
33        loss = mse_loss(A_val, Y_val)
```

```

30         self.validation_loss.append(loss)
31
32         print(f'Fold {fold}, Loss: {loss}')
33         fold += 1
34

```

It shuffles the x and y data and split them into 5 parts. The method then iterates over the number of folds(k). For each fold, it determines the start and end indices for the validation set. The remaining indices are used for the training set. This ensures that each fold has a unique validation set and the rest of the data is used for training.

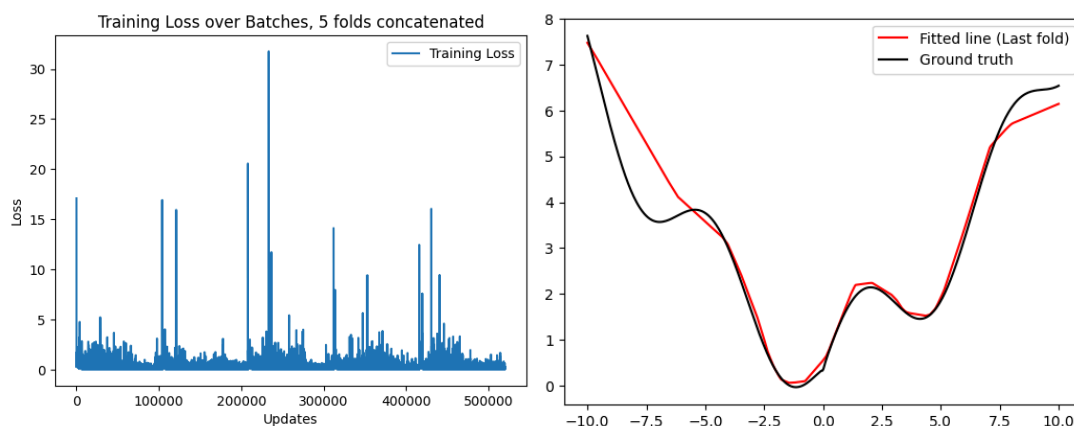
Within each fold, the method resets the model parameters to ensure that the training starts from scratch. It then trains the model using MBGD on the training set. After training, the model's performance is evaluated on the validation set using the inference method. The MSE loss is calculated for the validation set and appended to the [validation\_loss list. The loss for each fold is printed to provide feedback on the training process.

```

1  Training MBGD: 100%|██████████| 8000/8000 [00:11<00:00, 679.42it/s]
2  Fold 1, Loss: 0.135129258476114
3  Training MBGD: 100%|██████████| 8000/8000 [00:11<00:00, 704.24it/s]
4  Fold 2, Loss: 0.37937692480483626
5  Training MBGD: 100%|██████████| 8000/8000 [00:11<00:00, 703.15it/s]
6  Fold 3, Loss: 0.05067628846152356
7  Training MBGD: 100%|██████████| 8000/8000 [00:11<00:00, 716.51it/s]
8  Fold 4, Loss: 0.7495206586337292
9  Training MBGD: 100%|██████████| 8000/8000 [00:11<00:00, 704.32it/s]
10 Fold 5, Loss: 0.3724279561516599

```

## Training of 5 folds



## Trying different Hyperparameters

### Test 1

```

1  Average validation loss: 0.3374262173055726
2  Hyperparameters are: Epochs: 8000, LR: 0.4
3  Model structure: [<modules.Linear object at 0x7fd263aa6670>, <modules.ReLU object
  at 0x7fd263aa68b0>, <modules.Linear object at 0x7fd263aa6190>, <modules.ReLU object
  at 0x7fd263aa6070>, <modules.Linear object at 0x7fd263aa64c0>]

```



## Test 2 - Simpler Nets

```
1 Average validation loss: 0.46817570040655837
2 Hyperparameters are: Epochs: 8000, LR: 0.4
3 Model structure:
4 Layer 1: Linear
5 Input size: 1, Output size: 128
6 Layer 2: ReLU
7 Layer 3: Linear
8 Input size: 128, Output size: 1
```

## Test 2 - Simpler Nets, extensive training

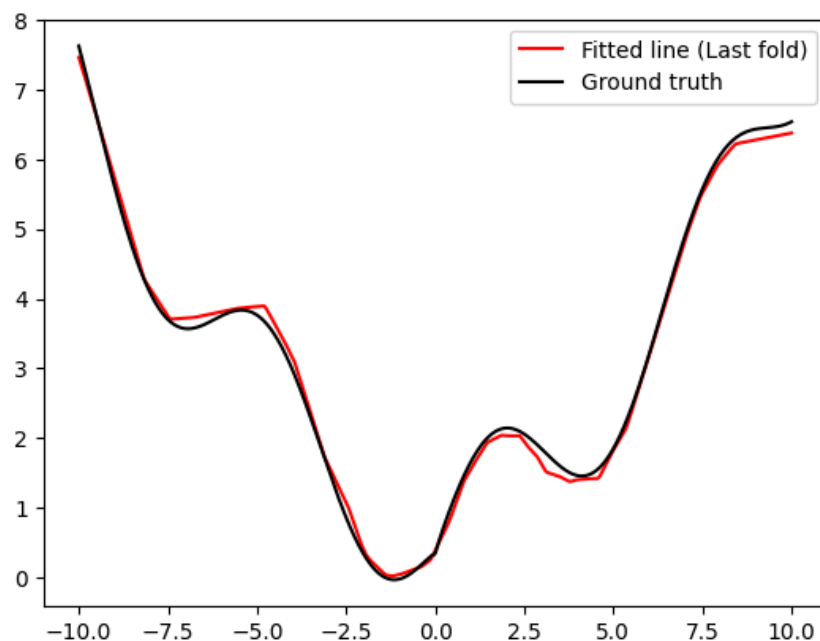
```
1 Average validation loss: 0.3243354289653034
2 Hyperparameters are: Epochs: 15000, LR: 0.4
3 Model structure:
4 Layer 1: LinearInput size: 1, Output size: 128
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 128, Output size: 1
```

## Test 3 - Deeper Nets

```
1 Average validation loss: 0.09484475132123227
2 Hyperparameters are: Epochs: 8000, LR: 0.4
3 Model structure:
4 Layer 1: LinearInput size: 1, Output size: 128
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 128, Output size: 512
7 Layer 4: ReLU
8 Layer 5: LinearInput size: 512, Output size: 1
9
```

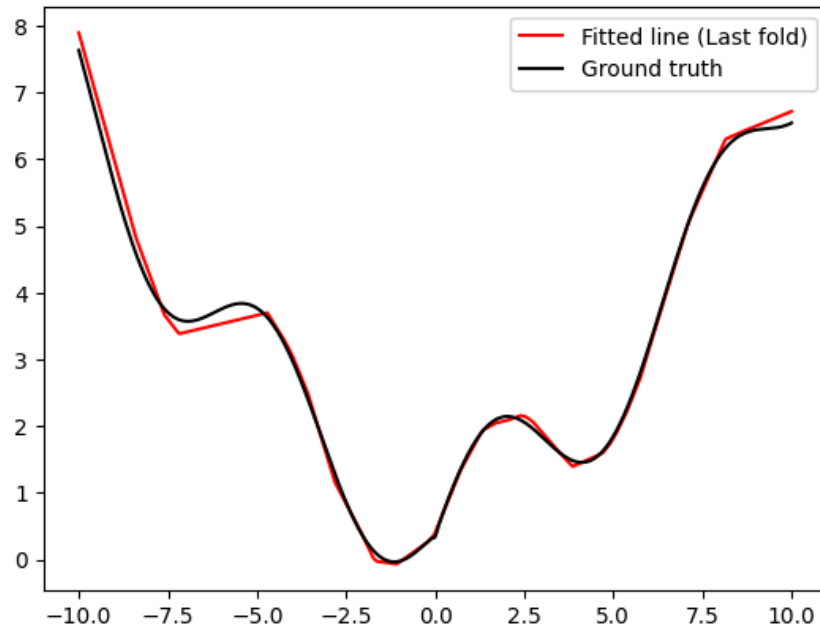
Note: this network took 23m39.7s to finish the 5-fold cross validation.

The result looks pretty good though.



## Test 4 - Simply Better Hyperparams

```
1 Average validation loss: 0.06755265032590395
2 Hyperparameters are: Epochs: 8000, LR: 0.1
3 Model structure:
4 Layer 1: LinearInput size: 1, Output size: 64
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 64, Output size: 128
7 Layer 4: ReLU
8 Layer 5: LinearInput size: 128, Output size: 1
```



Lots easier to train.

## Summary

- Deeper and wider nets are more difficult to train.

```
1 eval_layers = [
2     Linear(input_size=1, output_size=128),
3     ReLU(),
4     Linear(input_size=128, output_size=256),
5     ReLU(),
6     Linear(input_size=256, output_size=512),
7     ReLU(),
8     Linear(input_size=512, output_size=1)
9 ]
```

with a network like this, the speed is about **10~20 it/s**, while using a simpler net in simpler nets, the speed can go up to **~1700it/s**.

```

1 Training MBGD: 0%|          | 0/8000 [00:00<?, ?it/s]
2 Training MBGD: 100%|██████████| 8000/8000 [08:42<00:00, 15.32it/s]
3 Fold 1, Loss: 0.09874567588570843
4 Training MBGD: 0%|          | 0/8000 [00:00<?, ?
  it/s]/tmp/ipykernel_29773/4177186514.py:4: RuntimeWarning: overflow encountered in
  square
5     return np.mean((predicted - target) ** 2)
6 /root/sdm274/assgn05_MLP/modules.py:45: RuntimeWarning: overflow encountered in
  matmul
7     dw = (dz.T @ X) / m # dz.T ( output_size, m) @ X (m, input_size) = dw
  (output_size, input_size)
8 /root/sdm274/assgn05_MLP/modules.py:90: RuntimeWarning: invalid value encountered
  in multiply
9     dz = da * (self.output_cache > 0).astype(float)
10 /root/sdm274/assgn05_MLP/modules.py:45: RuntimeWarning: invalid value encountered
  in matmul
11     dw = (dz.T @ X) / m # dz.T ( output_size, m) @ X (m, input_size) = dw
  (output_size, input_size)
12 Training MBGD: 55%|██████    | 4396/8000 [04:29<03:41, 16.29it/s]

```

It is so resource-hungry that it might trigger **overflows**.

- Deep nets capture nonlinearities better, when properly trained.
- Shallow networks do not have enough layers to learn complex hierarchical features, making them less capable of capturing intricate nonlinear relationships. However, If you train them enough, you might have it yield a small validation loss, with the danger of overfitting.
- The improvement of accuracy is not proportionately big when you go from a normal model to a large model. The more accuracy you want, the more resources it takes to train and run a model.
- Hardware problem:

名称	状态	94% CPU	3% 磁盘	66% 内存	0% 网络
VmmemWSL		73.6%	0.1 MB/秒	3,703.5 MB	0 Mbps

What the program does is do tons of matrix multiplications for the training. It can cause the CPU to be fully loaded. But those matrix multiplications can be well accelerated by a GPU, which is suitable for computing simple calculations for millions of times.

## Classification Using MLP

### Data gathering

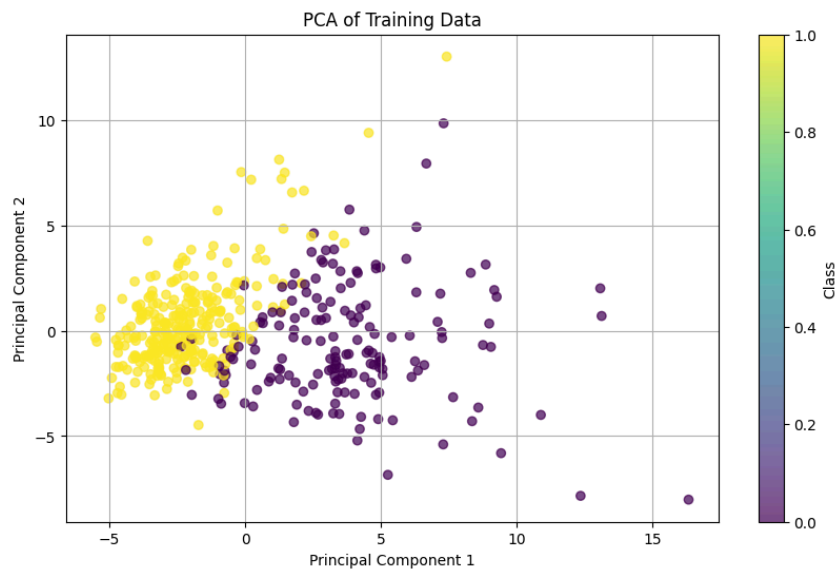
We choose the Wisconsin breast cancer dataset for classification.

```

1 Training data shape: (455, 30)
2 Testing data shape: (114, 30)

```

Here I used a **PCA to reduce the dimensions of the data to 2**. Note that we use the original 30-feature dataset for training and predicting, and we're **only using PCA for visualizing purpose**.



I defined an MLPClassifier class that inherits from the wrapped MLP class. This structure ensures reusability.

```

1  class MLPClassifier(MLP):
2
3      def __init__(self, layers, epochs=1000, lr=0.01, input_shape=1,
4      output_shape=1):
5          super().__init__(layers, epochs, lr, input_shape, output_shape)
6          self.threshold = 0.5
7
8      def get_loss(self, y_true, y_pred):
9          y_pred = y_pred.reshape(-1, 1)
10         # Clip predictions to prevent log(0)
11         y_pred = np.clip(y_pred, 1e-12, 1. - 1e-12)
12         # Compute cross-entropy loss
13         return -np.mean(y_true * np.log(y_pred) + (1 - y_true) * np.log(1 -
14         y_pred))
15
16     def get_loss_grad(self, y_pred, y_true):
17         # The loss here is the cross-entropy loss
18         # It accepts matrix with the shape (x, n_classes)
19         # Where x is the number of samples and n_classes is the number of classes
20
21         y_pred = y_pred.reshape(-1, 1)
22         # Clip predictions to prevent division by zero
23         y_pred = np.clip(y_pred, 1e-12, 1. - 1e-12)
24         # Compute gradient of cross-entropy loss
25         return (y_pred - y_true) / (y_pred * (1 - y_pred))
26
27     def predict(self, X):
28         # Forward pass
29         y_pred = self.forward(X)
30
31         # Convert the output to binary
32         return (y_pred > self.threshold).astype(int)
33

```

I reused the model structure for the nonlinear function training. It balances performance and resource.

```

1 layers = [
2     Linear(input_size=30, output_size=64),
3     ReLU(),
4     Linear(input_size=64, output_size=128),
5     ReLU(),
6     Linear(input_size=128, output_size=1),
7     Sigmoid()
8 ]
9

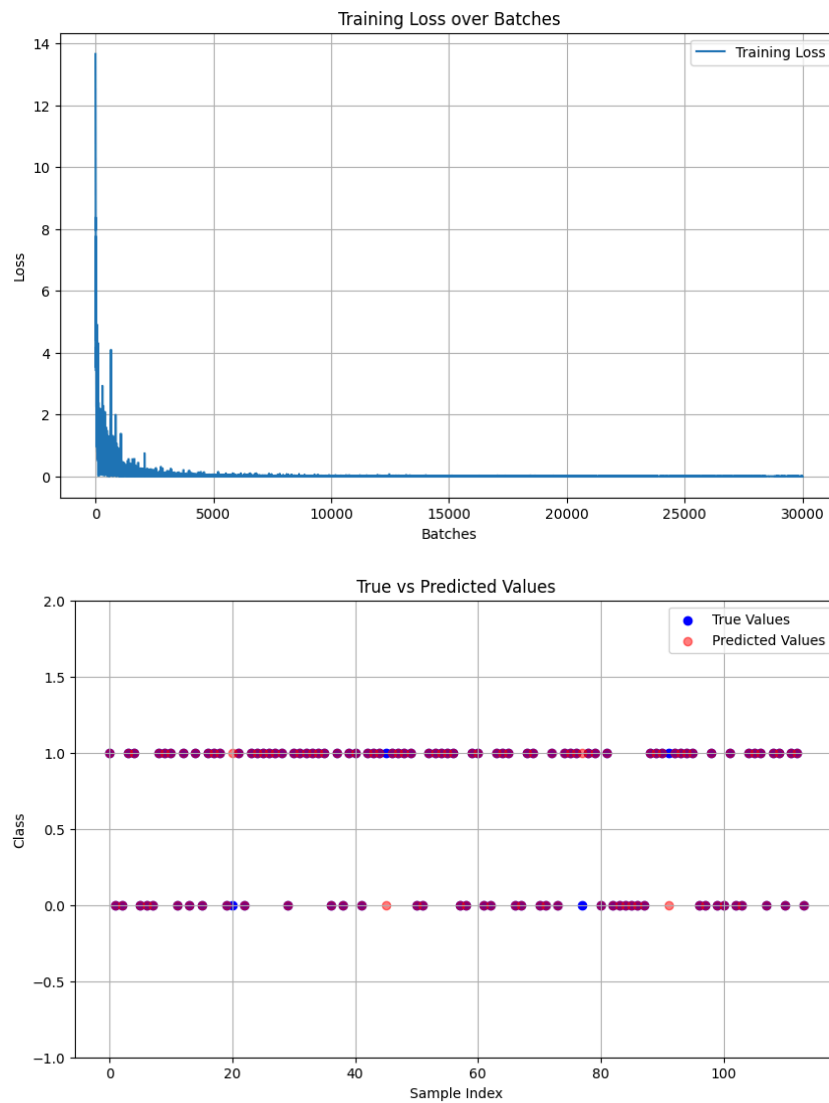
```

I used MBGD in training, reaching about 400it/s with a good accuracy.

```

1 Training MBGD: 100%|██████████| 2000/2000 [00:04<00:00, 412.34it/s]

```

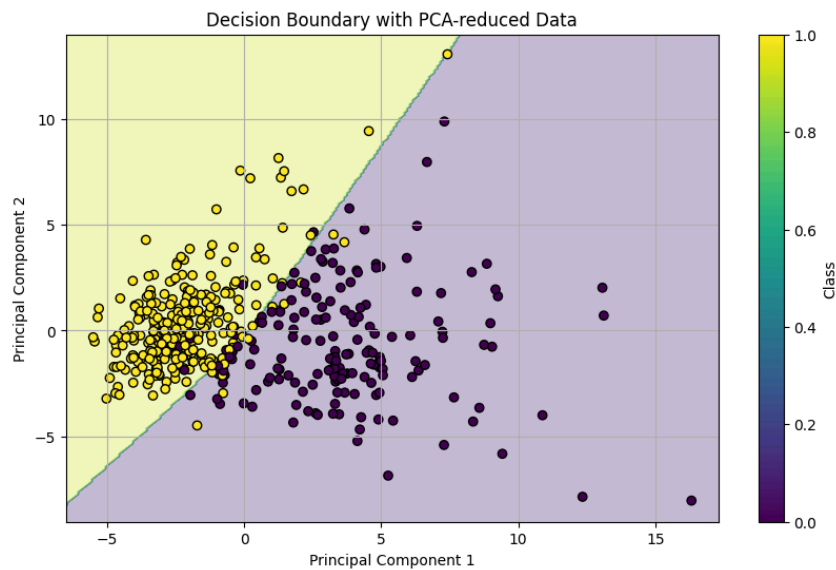


The model achieved high F1 score on this dataset.

```

1 Accuracy: 0.9649122807017544
2 Recall: 0.971830985915493
3 Precision: 0.971830985915493
4 F1 Score: 0.971830985915493

```



Above is the decision boundary of the trained MLP. As shown in the graph the boundary is **a folded line** instead of a straight line. This indicates that the MLP model has learned a **nonlinear decision boundary**. Unlike linear classifiers, which can only separate classes with straight lines or planes, an MLP with at least one hidden layer can model complex, nonlinear relationships between features.

The dataset we choose here is largely linearly separable. Most parts of it can be approximately separated by a straight line, therefore the outliers won't punish the MLP to behave very non-linearly. Theoretically we can push the training and the parameters to extremes to let it show more nonlinear characteristics, but that requires much more training resources and is somewhat beyond the scope of this assignment.

Shown below is some sets cross-validation result.

```
1 Average validation loss: 0.8235056154333691
2 Hyperparameters are: Epochs: 3000, LR: 0.1
3 Model structure:
4 Layer 1: LinearInput size: 30, Output size: 64
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 64, Output size: 128
7 Layer 4: ReLU
8 Layer 5: LinearInput size: 128, Output size: 1
9 Layer 6: Sigmoid
```

```
1 Average validation loss: 0.7517043588174628
2 Hyperparameters are: Epochs: 2000, LR: 0.1
3 Model structure:
4 Layer 1: LinearInput size: 30, Output size: 64
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 64, Output size: 128
7 Layer 4: ReLU
8 Layer 5: LinearInput size: 128, Output size: 1
9 Layer 6: Sigmoid
10
11 F1 Score: 0.9859154929577465
```

```
1 Average validation loss: 0.8006151770380503
2 Hyperparameters are: Epochs: 2000, LR: 0.01
3 Model structure:
4 Layer 1: LinearInput size: 30, Output size: 64
5 Layer 2: ReLU
6 Layer 3: LinearInput size: 64, Output size: 128
7 Layer 4: ReLU
8 Layer 5: LinearInput size: 128, Output size: 1
9 Layer 6: Sigmoid
10
11 F1 Score: 0.9929078014184397
```

We have almost come to the best of this structure, achieving very high F1 scores.

## Summary

---

The results and findings provide insights into the MLP's representational power, demonstrating the ability of them to approximate nonlinear functions and capture complex patterns in data.