

# Artificial Intelligence

Neural Networks

Fall 2024

**Parsa Darban 810100141**

This file contains the report and results of the simulations conducted.

[Abstract](#)

[Layer of Neural Network](#)

[Regularization](#)

[Multi-Layer Fully Connected Network](#)

[Classification](#)

[Conclusion](#)

---



## Abstract

In this project, we will explore deep neural networks. First, we will implement neural networks and analyze their performance. Then, we will combine these networks to create a deep neural network. Finally, we will use this network for classification and regression on two datasets: MNIST and California housing dataset, and we will train and evaluate the model.

## Layer of a neural network

### Affine layer

An affine layer, also called a fully connected layer or dense layer, is a layer in which the input signal of the neuron is multiplied by the weight, summed, and biased. An affine layer can be a layer of an artificial neural network in which all contained nodes connect to all nodes of the subsequent layer. Affine layers are commonly used in convolutional neural networks. An affine layer is consist of forward and backward pass. We can get the output of forward pass by dot product inputs and weights with the addition of bias.

$$output = xW + b$$

$x \in R^{n \times D}$  is input data with n sample and D features and  $w \in R^{D \times M}$  is Weight matrix, where D is the number of input features and M is the number of output features.  $b$  is bias vector.

For the backward pass, we use the chain rule of derivatives. For gradient of the loss with respect to x:

$$\frac{\partial Loss}{\partial x} = \frac{\partial Loss}{\partial out} \frac{\partial out}{\partial x}$$

$\frac{\partial Loss}{\partial out}$  contains the gradients of the loss with respect to the output of the affine layer and  $\frac{\partial out}{\partial x} = w^T$ . This gives a tensor of shape N×D, which is the gradient with respect to the flattened input.

For gradient of the loss with respect to w, like the x we use chain rule and we have:

$$\frac{\partial Loss}{\partial w} = \frac{\partial Loss}{\partial out} \frac{\partial out}{\partial w} = x^T \frac{\partial Loss}{\partial out}$$

For gradient of the loss with respect to b:

$$\frac{\partial Loss}{\partial b} = \sum_{i=1}^n \frac{\partial Loss}{\partial out_i}$$

These gradients are used during the training process to update the parameters W and b via optimization algorithms like gradient descent.

The result of comparing the output of the written code and the values obtained from the provided functions is as follows:

The error should be around e-9 or less.  
Testing affine\_forward function:  
difference: 9.769849468192957e-10

The error should be around e-10 or less.  
Testing affine\_backward function:  
dx error: 5.399100368651805e-11  
dw error: 9.904211865398145e-11  
db error: 2.4122867568119087e-11

## ReLU activation

The ReLU function defined mathematically as:

$$ReLU = \max(0, x)$$

This activation function introduces non-linearity to the model and is widely used because efficiency and mitigates the vanishing gradient problem in deep networks (compared to sigmoid/tanh).

So, for the forward pass, we use the mentioned function, and the result of comparing it with the provided functions is as follows:

```
The error should be on the order of e-8.  
Testing relu_forward function:  
difference: 4.999999798022158e-08
```

Its derivative is as follows.

$$\frac{\partial ReLU}{\partial x} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

And because of the chain rule of derivative we have:

$$\frac{\partial Loss}{\partial x} = \frac{\partial Loss}{\partial ReLU} \frac{\partial ReLU}{\partial x}$$
$$\frac{\partial Loss}{\partial x} = \begin{cases} \frac{\partial Loss}{\partial ReLU} & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

The result of comparing the backward pass output of the written code and the values obtained from the provided functions is as follows:

```
The error should be on the order of e-12  
Testing relu_backward function:  
dx error: 3.2756349136310288e-12
```

## Sigmoid activation

The sigmoid function is defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

This function squashes the input  $x$  into a range between 0 and 1, making it useful for binary classification and as an activation function in neural networks.

So, for the forward pass, we use the mentioned function, and the result of comparing it with the provided functions is as follows:

The error should be on the order of  $e^{-7}$   
 Testing `sigmoid_forward` function:  
 difference:  $6.383174040859927e-07$

The derivative of the sigmoid function is:

$$\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$$

So, for the backward pass, we use this formula and we have:

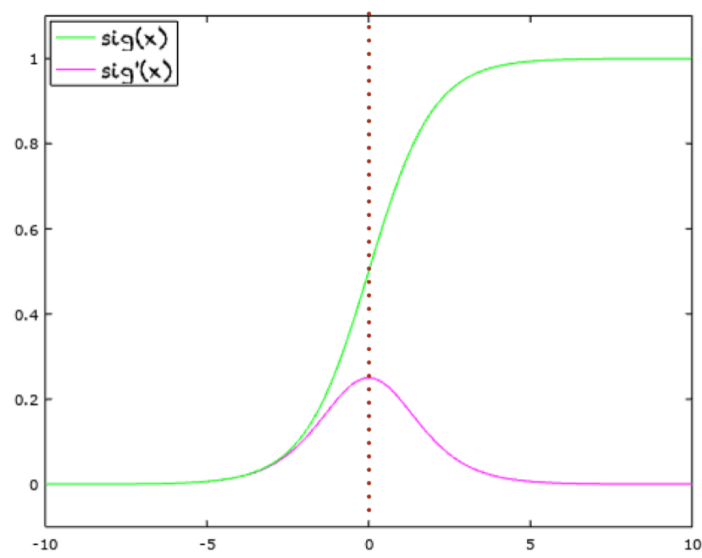
$$\frac{\partial Loss}{\partial x} = \frac{\partial Loss}{\partial \sigma} \frac{\partial \sigma}{\partial x}$$

$$\frac{\partial Loss}{\partial x} = \frac{\partial Loss}{\partial \sigma} (\sigma(x) \cdot (1 - \sigma(x)))$$

The result of comparing the backward pass output of the written code and the values obtained from the provided functions is as follows:

The error should be on the order of  $e^{-11}$   
 Testing `sigmoid_backward` function:  
 dx error:  $3.446520386706568e-11$

The following graph represents the sigmoid function and its derivative.



## Sandwich layers

A "sandwich layer" in deep neural networks refers to a functional block where a nonlinear activation is sandwiched between two linear operations, such as affine transformations or convolutions. This structure enhances the network's expressivity by alternating between linear transformations and nonlinear activations, enabling the modeling of complex patterns in data. These layers often include regularization methods like dropout and normalization techniques such as batch normalization to ensure stable gradient flow during backpropagation and reduce overfitting. Examples include fully connected layers, convolutional layers, and their variations in architectures like ResNet and DenseNet, which incorporate additional elements like residual connections.

The sandwich layer pattern is fundamental in modern deep learning because it encapsulates modular, reusable components that enhance performance and simplify design. Beyond simple linear-nonlinear sequences, it supports advanced configurations, including residual blocks and attention mechanisms in transformers. These layers contribute to stability in optimization and serve as the backbone for building powerful architectures for tasks like classification, detection, and regression.

This layer in this project is constructed by combining a fully connected layer with the ReLU activation function. Thus, it first performs calculations based on the features and weights using the affine layer formula, then passes the output to the ReLU activation function, which produces the final output.

$$y = xW + b \text{ (affine layer)}$$

$$out = ReLU(y) \text{ (ReLU activation function)}$$

So, for the forward pass, we use the mentioned function. The `affine_relu_backward` function first computes the gradient through the ReLU activation by applying the derivative  $\frac{\partial Loss}{\partial y}$ . Then it uses the chain rule to compute gradients for the affine transformation:  $dx$ ,  $dw$ , and  $db$ . For example:

$$\frac{\partial Loss}{\partial x} = \frac{\partial Loss}{\partial y} W^T$$

Another formula is similar to affine layer.

The result of comparing the forward and backward pass output of the written code and the values obtained from the provided functions is as follows:

```
Relative error should be around e-10 or less.
Testing affine_relu_forward and affine_relu_backward:
dx error:  2.2995909845854045e-11
dw error:  8.162011105764925e-11
db error:  7.826724021458994e-12
```

## Softmax Loss function

Loss functions are one of the most important aspects of neural networks, as they (along with the optimization functions) are directly responsible for fitting the model to the given training data. A loss function is a function that compares the target and predicted output values; measures how well the neural network models the training data. When training, we aim to minimize this loss between the predicted and target outputs.

The Softmax function is used to convert raw logits (unscaled scores) from a model into a probability distribution. It is typically applied to the output of the final layer in a neural network for classification tasks. It produces the predicted class probabilities.

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_{j=1}^C e^{y_j}}$$

In practice, instead of separately applying Softmax, many deep learning libraries provide a combined function (e.g., `torch.nn.CrossEntropyLoss`), which directly computes both the Softmax and Cross-Entropy loss in a single step, optimizing numerical efficiency and stability.

## MSE Loss Function

The Mean Squared Error loss between the predicted values  $x$  and the true values  $y$  is given by:

$$MSE = \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2$$

$n$  is the number of examples (the length of the vectors  $x$ )

The gradient of the MSE loss with respect to the predicted values  $x$  is given by:

$$\frac{\partial MSE}{\partial x_i} = \frac{2}{n} (x_i - y_i)$$

This is the gradient for each predicted value, and we sum it across the batch to get the total gradient for the loss with respect to  $x$ .

The result of test is:

```
Loss should be close to 1.9 and dx error should be around e-9
Testing mse_loss:
loss:  1.8672282748726519
dx error:  2.8607953262121067e-09
```

This function has numerous properties that make it especially suited for calculating loss. The difference is squared, which means it does not matter whether the predicted value is above or below the target value; however, values with a large error are penalized. MSE is also a convex function with a clearly defined global minimum. This allows us to more easily utilize gradient descent optimization to set the weight values.

One disadvantage of this loss function is that it is very sensitive to outliers; if a predicted value is significantly greater than or less than its target value, this will significantly increase the loss.

## Regularization

Regularization is a crucial technique in machine learning that helps prevent overfitting. Overfitting occurs when a model performs well on training data but fails to generalize to new, unseen data. This problem arises when the model becomes too complex and is sensitive to the details and noise in the training data. Regularization techniques prevent this phenomenon by applying a penalty to the complexity of the model, helping it generalize better to new data.

There are various forms of regularization, but two of the most commonly used techniques are L2 regularization (Ridge regularization) and L1 regularization (Lasso regularization).

- L1 Regularization (Lasso Regularization): L1 regularization adds a penalty equal to the sum of the absolute values of the model parameters to the loss function.

$$Reg = \lambda \sum_i |W_i|$$

- L2 Regularization (Ridge Regularization): L2 regularization adds a penalty equal to the sum of the squared values of the model parameters (weights) to the loss function.

$$Reg = \frac{\lambda}{2} \sum_i W_i^2$$

- Elastic Net Regularization: Elastic Net is a combination of both L1 and L2 regularization. It is useful when there are many correlated features in the dataset, as it can group them together and select or shrink them jointly.

$$Reg = \lambda_1 \sum_i |W_i| + \lambda_2 \sum_i W_i^2$$

Where  $\lambda_1$  and  $\lambda_2$  control the contributions from L1 and L2 regularization, respectively.

The regularization strength  $\lambda$  is a hyperparameter that needs to be tuned. If  $\lambda$  is too large, the model will be too constrained and unable to fit the training data well, leading to underfitting. On the other hand, if  $\lambda$  is too small, the model might overfit to the training data. The optimal value of  $\lambda$  can be found using cross-validation, where the model is trained on the training set with various values of  $\lambda$  and the best-performing value is selected based on the validation set.

Regularization helps prevent overfitting by introducing a penalty for large weights, encouraging the model to focus on the most relevant features while ignoring noisy, less significant ones. This process shrinks the weights, leading to a simpler model that generalizes better to unseen data. By controlling the size of model parameters, regularization reduces sensitivity to fluctuations in the training data, improving the model's ability to perform well on new, unseen examples.



## Multi-Layer Fully Connected Network

This code defines a class `FullyConnectedNet` which is a multi-layer neural network implementation for classification or regression tasks. The constructor `__init__` takes various parameters like the category of the task (classification or regression), hidden layer dimensions, input/output sizes, regularization strength, weight initialization scale, and data type for computation. Inside the constructor, the network's weights ( $W$ ) and biases ( $b$ ) are initialized for each layer. The layers are created sequentially, where the weights are initialized using a random normal distribution scaled by `weight_scale`, and biases are initialized to zeros. The dimensions of each layer are determined by the input and hidden layer sizes, with the final layer's output matching the output dimension (e.g., 10 for classification with MNIST). The parameters are cast to the specified data type (`np.float32` or `np.float64`).

The loss method computes both the forward and backward passes through the network. During the forward pass, the input  $X$  is propagated through each layer sequentially, applying an affine transformation followed by a ReLU activation. The output of each layer is cached for use in the backward pass. The final output scores are computed by applying a final affine transformation to the last hidden layer's output. If no labels  $y$  are provided (test mode), the function simply returns these scores, which represent the predicted class probabilities or continuous values, depending on the task.

When labels  $y$  are provided (train mode), the method computes the loss based on the category of the task. If the task is classification, a softmax loss function is used, which is suitable for multi-class classification problems like MNIST, where each class corresponds to a digit. If the task is regression, a mean squared error (MSE) loss function is used to evaluate the network's performance. In both cases, the loss is combined with an L2 regularization term to penalize large weights, controlled by the regularization parameter `reg`. The regularization term is computed as the sum of squared weights across all layers, and it is added to the final loss value.

The backward pass computes gradients for the weights and biases using the chain rule. First, the gradient of the loss with respect to the output scores is calculated, and then the gradients are propagated backward through each layer using the cached information from the forward pass. The gradients of the weights ( $W$ ) and biases ( $b$ ) are computed at each layer. Additionally, L2 regularization is applied to the gradients of the weights, ensuring that the model learns to prevent overfitting by controlling the magnitude of the weights. The loss method finally returns both the computed loss and the gradients, which will be used for optimization during training.

In the provided class, L2 regularization is implemented. After computing the main loss (which could be softmax loss for classification or MSE loss for regression), the regularization term is added to the loss. The regularization term is computed by summing the squares of all the weights across all layers.

$$L_{total} = L_{Loss} + \frac{\lambda}{2} \sum_i ||W_i||^2$$

Where  $L_{Loss}$  is the original loss, and  $||W_i||^2$  is the squared sum of weights for each layer  $i$ .

In the backward pass, the gradient of the regularization term is added to the gradient of the loss with respect to the weights. The gradient of the regularization term for each weight is simply  $\lambda W_i$ , where  $W_i$  is the weight for the layer.

$$\frac{\partial L_{total}}{\partial W_i} = \frac{\partial L_{Loss}}{\partial W_i} + \lambda W_i$$

This helps the model update the weights to not only minimize the main loss but also reduce their magnitude by a factor proportional to  $\lambda$ .

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. This is a good way to see if the initial losses seem reasonable. The result is:

```
Running check with reg = 0
Initial loss: 2.300479089768492
W1 relative error: 1.025267445229243e-07
W2 relative error: 2.2120479296856732e-05
W3 relative error: 4.562327873666505e-07
b1 relative error: 4.660094188020383e-09
b2 relative error: 2.085654200257447e-09
b3 relative error: 1.689724888469736e-10
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 3.904542008453064e-09
W2 relative error: 6.86942277940646e-08
W3 relative error: 1.0
b1 relative error: 1.4752427901105963e-08
b2 relative error: 1.4615869332918208e-09
b3 relative error: 1.3200479211447775e-10
```

## SGD

Stochastic Gradient Descent (SGD) is an optimization algorithm used to minimize a function, typically a loss or cost function in machine learning models. It is a variation of gradient descent, a widely used optimization method that aims to find the minimum (or maximum) of a function.

In the context of machine learning, the goal is to minimize the loss function, which measures how well the model is performing. The gradient descent algorithm updates the model's parameters (weights) iteratively by computing the gradient of the loss function with respect to the parameters and adjusting the weights in the direction of the negative gradient (downhill).

However, batch gradient descent, which computes the gradient using the entire dataset, can be computationally expensive, especially for large datasets. Stochastic Gradient Descent (SGD) addresses this by using only a single data point or a small batch of data points to compute the gradient in each iteration.

Stochastic Gradient Descent (SGD) is an optimization algorithm used to minimize a loss function  $L(w)$  by iteratively updating model parameters  $W$ . For each data point  $(x_i, y_i)$ , the gradient of the loss  $\nabla L(W; x_i, y_i)$  is computed to determine how the parameters should change to reduce the loss. The weights are then updated using the rule  $W_t = W_{t-1} - \eta \nabla L(W_{t-1}; x_i, y_i)$ , where  $\eta$  is the learning rate. This update step moves the parameters in the direction opposite to the gradient to minimize the loss. The process is repeated over multiple iterations until the loss converges to a minimum.

Stochastic Gradient Descent (SGD) is an efficient optimization algorithm that updates model parameters using gradients computed from individual data points or mini-batches, making it faster than batch gradient descent for large datasets. While each update is noisy, this randomness helps escape local minima and explore the solution space effectively. Proper learning rate tuning is critical to ensure convergence, as overly high rates can cause oscillations, while decreasing rates stabilize the updates over time. Variants like mini-batch SGD balance efficiency and stability, while momentum-based approaches, such as SGD with Momentum and Nesterov Accelerated Gradient, smooth updates and accelerate convergence for more efficient optimization.

## Momentum

Momentum is a technique used to accelerate gradient descent algorithms by considering not just the current gradient but also the previous gradients. It can help the optimization process converge more quickly and avoid oscillations, especially in regions with noisy gradients or steep areas in the cost function. In simple terms, momentum helps smooth out the updates to the parameters by adding a "memory" of past gradients. In machine learning, momentum is often added to the update rule of gradient descent algorithms, like Stochastic Gradient Descent (SGD), to improve the optimization process.

Momentum introduces a "velocity" term that accumulates the gradient information from previous iterations, allowing for smoother and faster updates. The update rule becomes:

$$v_t = \beta v_{t-1} + (1 - \beta) \nabla L(W_{t-1})$$

$$w_t = w_{t-1} - \eta v_t$$

The key idea of momentum is to accumulate past gradients and use this accumulated velocity to update the weights. This equation shows how the current velocity is influenced by both the previous velocity and the current gradient. The term  $\beta v_{t-1}$  is the "memory" of the previous updates, and  $(1 - \beta) \nabla L(W_{t-1})$  is the new gradient contribution.

The velocity  $v_t$  is subtracted from the current weight, and the learning rate  $\eta$  controls how large the step will be. Since the velocity is a weighted combination of gradients from previous iterations, it can help the optimizer maintain a consistent direction toward the minimum, smoothing out the updates.

Momentum improves optimization by accelerating convergence, reducing oscillations, escaping local minima, and stabilizing updates. It achieves this by incorporating a fraction of the previous update's velocity into the current gradient step, enabling faster movement in consistent gradient directions and smoothing noisy updates. This makes momentum particularly effective for navigating complex loss surfaces with flat regions, sharp curvatures, or saddle points, leading to more stable and efficient training compared to standard gradient descent methods.

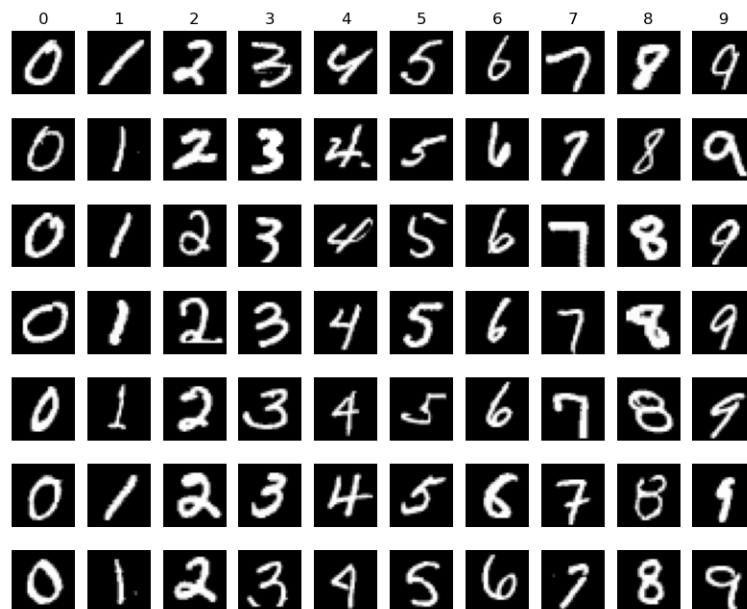
When using SGD with Momentum, the optimizer benefits from both the efficiency of gradient-based updates and the stability and speed provided by momentum, making it a popular choice in training deep neural networks.

When we test this on some data that we have expected value and calculate the error of this and our result, we have:

```
next_w error: 8.882347033505819e-09  
velocity error: 4.269287743278663e-09
```

## Classification

Now, we move on to the classification task using the MNIST dataset. MNIST is a widely used dataset of handwritten digits that contains 60,000 handwritten digits for training a machine learning model and 10,000 handwritten digits for testing the model. We have randomly plotted some of its data, as shown below:



Data normalization is an important step which ensures that each input parameter has a similar data distribution. This makes convergence faster while training the network.

In this section, we will build a neural network to maximize accuracy in classification.

The parameter values define and train a deep neural network for the MNIST dataset. In this model, the SGD algorithm with momentum and regularization is used for optimization. Some hidden layers with  $n$  neurons each and selected hyperparameters (such as learning rate and regularization value) are tuned to achieve a balance between accuracy and generalization. The training process is monitored by periodically reporting Loss and Accuracy values.

Several steps were taken to improve accuracy. Initially, regularization was set to 0, and the accuracy of the model was increased by increasing the number of neurons per layer. In one case, we achieved even higher accuracy by adding regularization to the model and reducing the number of neurons in each layer. The number of epochs was also found to be critical for improving model accuracy since parameters are updated during each epoch, leading to an eventual saturation point in accuracy.

The results are as follows:

```
(Epoch 100 / 100) train acc: 0.941000; val_acc: 0.948300
```

```
Input_dim = [512,512] , reg = 0
```

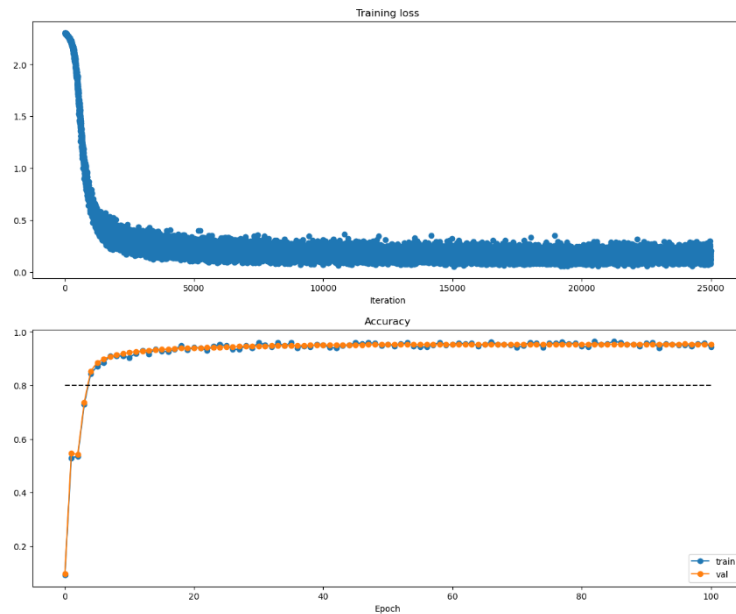
```
(Epoch 100 / 100) train acc: 0.943000; val_acc: 0.953900
```

```
Input_dim = [1024,1024] , reg = 0
```

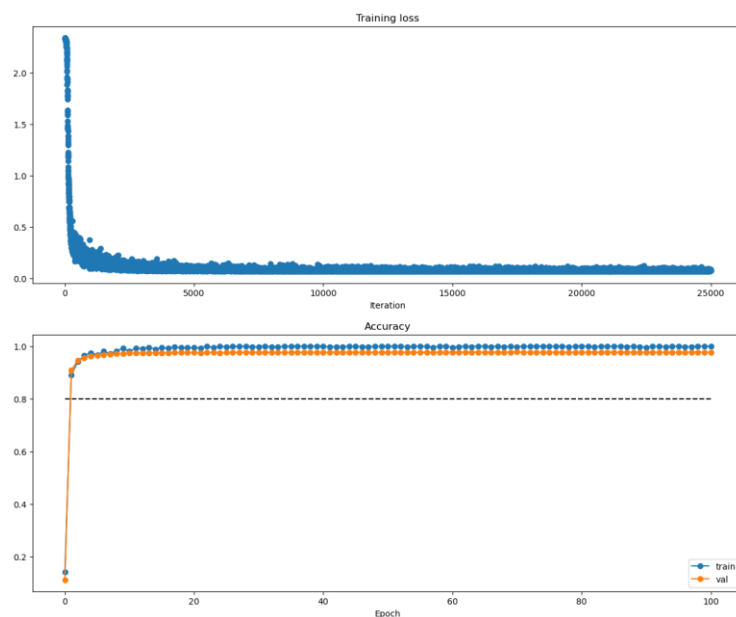
```
(Epoch 100 / 100) train acc: 0.999000; val_acc: 0.977100
```

```
Input_dim = [512,512] , reg = 10-3
```

Additionally, for the last two cases, we have the growth charts showing the changes in Loss and Accuracy over the training process:



Input\_dim = [1024,1024] , reg = 0



Input\_dim = [512,512] , reg =  $10^{-3}$

As we can see, convergence is better in the second case. Additionally, the test results are as follows:

Validation set accuracy: 0.9543  
Test set accuracy: 0.9487

Input\_dim = [1024,1024] , reg = 0

Validation set accuracy: 97.75%  
Test set accuracy: 97.70%

Input\_dim = [512,512] , reg =  $10^{-3}$

**Hidden Dimensions (hidden\_dims):** Increasing dimensions improves model capacity but risks overfitting, while decreasing simplifies the model and may underfit.

**Regularization Strength (reg):** Higher values reduce overfitting by penalizing large weights; lower values increase model flexibility but may overfit.

**Learning Rate (learning\_rate):** Higher rates speed up convergence but risk instability; lower rates stabilize training but slow convergence.

**Number of Epochs (num\_epochs):** More epochs improve training but risk overfitting; fewer epochs may lead to undertraining.

**Batch Size (batch\_size):** Larger batches smooth gradients for stability; smaller batches add noise, potentially aiding escape from local minima.

**Learning Rate Decay (lr\_decay):** Faster decay stabilizes training late, while slower decay maintains learning speed but risks instability.

By carefully tuning these parameters, a balance between training speed, accuracy, and generalization can be achieved, depending on the specific dataset and task requirements.

## Conclusion

Activation functions in neural networks are used to introduce non-linearity into models. These functions help the model learn more complex patterns and prevent it from being limited to learning linear features. Without activation functions, neural networks would not be able to properly learn the non-linear features of data.

The ReLU (Rectified Linear Unit) function transforms all negative values to zero while leaving positive values unchanged. This characteristic makes it easy for the network to model non-linear features. One of the advantages of ReLU is that it allows faster computations compared to other activation functions. It's simple and fast, reduces issues like "zero gradients" that exist in functions like Sigmoid.

$$ReLU = \max(0, x)$$

The Sigmoid function maps input values to a range between 0 and 1. This makes it particularly suitable for binary classification tasks, especially when predicting the probability of an event occurring. It's suitable for modeling probabilities and binary classification outputs. However, a downside is that for large or small input values, the gradients approach zero (the "vanishing gradient" problem), which can make learning more difficult.

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

The cost function (or loss function) plays an important role in training machine learning models by measuring the quality of the model's predictions. This function represents the difference or error between the model's predictions and the actual values (labels). The goal of training a model is to minimize the cost function and improve the model's performance. In other words, the model should be adjusted in such a way that the cost function is minimized for the training data.

The Mean Squared Error (MSE) cost function is commonly used for regression problems. Its formula is as follows:

$$MSE = \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2$$

This cost function squares the difference between the predicted and actual values to penalize larger differences more heavily. In other words, the greater the model's prediction error, the higher the cost. The MSE function is used when the goal of the model is to predict continuous values, such as predicting house prices, temperature, or other numerical features.

This Cross-Entropy Loss is used for classification problems, especially for multi-class classification. The Cross-Entropy Loss function is used in classification tasks, particularly when the model provides probability predictions for multiple categories, such as classifying images into different categories. This cost function is specifically designed to evaluate the model's performance in classification problems. Cross-Entropy measures the mismatch between the model's predictions and the actual values. The function penalizes the model more when it makes incorrect predictions and applies a lower cost when the model predicts correctly. The greater the difference between predicted and actual probabilities, the higher the loss, making it more sensitive to wrong predictions.

Batch Normalization (BN) is a technique to normalize the inputs of each layer in a neural network. It stabilizes and accelerates the training process by ensuring that the inputs to each layer maintain a consistent distribution. BN reduces the problem of internal covariate shift, which occurs when the distribution of inputs to a layer changes due to updates in preceding layers during training.

It's first compute mean and variance for a given mini-batch of data with activations  $x$ :

$$\mu_B = 1/m \sum_{i=1}^m x_i \text{ and } \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

Here,  $m$  is the batch size,  $\mu_B$  is the mean, and  $\sigma_B^2$  is the variance of the batch. After that normalize the input by the computed value and then scale and shift to make the next layer input:

$$\bar{x}_i = \frac{x_i - \mu_B}{\sigma_B}$$

$$x_i = \gamma \bar{x}_i + \beta$$

BN learns to optimize the normalization parameters ( $\gamma$  and  $\beta$ ) during training, allowing the network to retain representational power.

Batch Normalization (BN) accelerates training by stabilizing the distribution of inputs to layers, allowing higher learning rates and reducing sensitivity to initialization. It introduces a regularization effect, slightly reducing overfitting due to batch-wise statistics. BN also mitigates vanishing /exploding gradients by normalizing intermediate activations, ensuring gradients remain within a manageable range during backpropagation. These factors contribute to improved generalization, helping the network perform better on unseen data, making BN a crucial technique for enhancing the efficiency and performance of deep neural networks. It is particularly useful in deep networks, where gradient instability such as vanishing or exploding gradients can be a problem. It is widely applied in convolutional layers, especially for tasks involving datasets like CIFAR-10 and ImageNet, to improve network stability. BN can also be used before fully connected layers to stabilize dense layers, although Dropout may suffice in smaller networks. Additionally, BN enables the use of higher learning rates, facilitating faster convergence during training. This makes BN an essential tool in improving the performance and efficiency of deep learning models.

Batch Normalization (BN) has some limitations, including its dependence on batch size, as performance may degrade with very small batches due to noisy statistics. It also introduces additional computational overhead and memory usage for maintaining the mean and variance statistics. Furthermore, BN is not always the best choice for sequence models like RNNs, where alternative normalization techniques such as Layer Normalization or Group Normalization may perform better in certain cases. These limitations should be considered when choosing the appropriate normalization technique for specific tasks.



Deep neural networks outperform classical machine learning methods when working with complex, large, and unstructured data such as images, videos, and time-series data. These networks are capable of automatically extracting features from the data and learning more complex patterns, which would require manual feature extraction and intricate tuning in classical methods. Particularly in tasks such as image recognition, natural language processing, and speech recognition, deep neural networks perform significantly better and more effectively.

If a neural network with 1000 parameters is trained on only 10 data points, it is likely to face overfitting. In this case, the model may perform very well on the training data but will likely perform poorly on new or test data because it has memorized the training features and lacks the ability to generalize to new data. This occurs because the model has far more parameters than the number of training data points, making it overfit the training data and lose its generalization capability. To prevent overfitting, regularization techniques such as Dropout, L2 regularization, or other methods should be used, and the amount of training data should be increased appropriately.