



Lecture 01: Introduction, Overview and DNN Basics

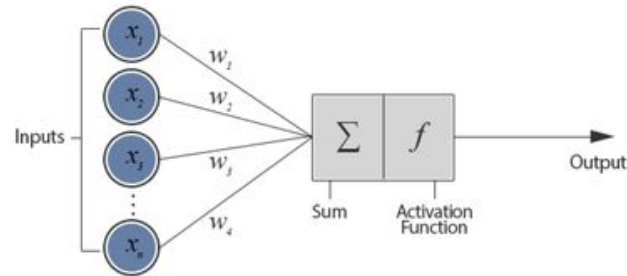
ECE-GY 9483/CSCI-GA 3033
Special Topics in Electrical Engineering EFFICIENT AI AND
HARDWARE ACCELERATOR DESIGN

Basics of Deep Neural Networks

- Multi-layer Perceptrons (MLPs)
 - Fully-connected layers
 - Activation functions
 - Loss function
 - Backpropagation
- How forward and backward propagation is performed?
- How to compute the gradient?
- How to update the weight?
- How to initialize the weight before training?

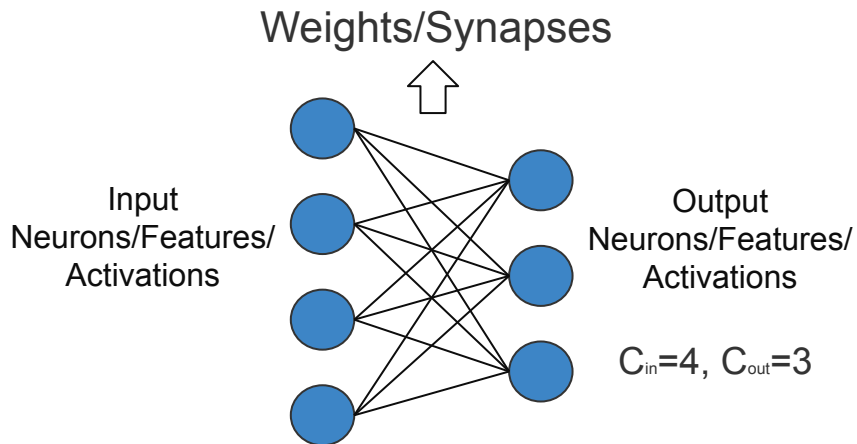
Multi-layer Perceptrons

- Usually consists of fully-connected layers with nonlinear activation functions.



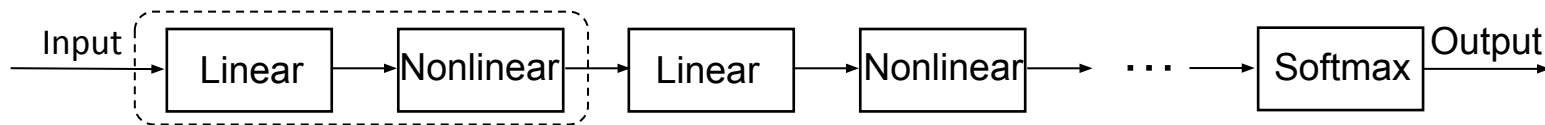
- A neural network consists of interconnected nodes, called neurons, organized into layers.
- Each neuron receives input signals (activations), performs a computation on them, and produces an output signal that may be passed to other neurons in the network.

Fully-connected layers (Linear layers)

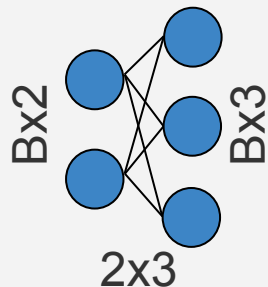
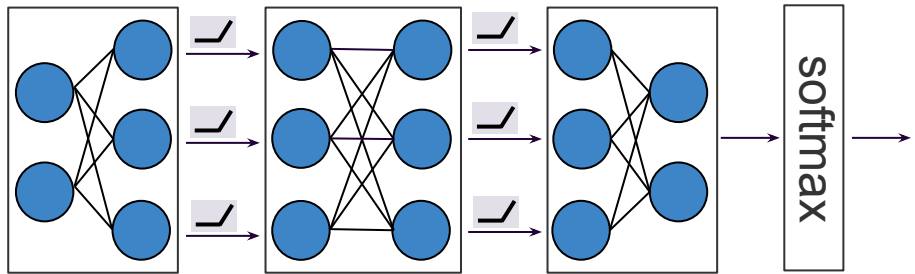
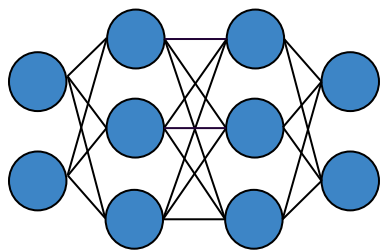


$$Y = XW + b$$

- X (input activations): $B \times C_{in}$
- Y (output activations): $B \times C_{out}$
- W (weights): $C_{in} \times C_{out}$
- b (bias): $1 \times C_{out}$
- C_{in} : Number of input activations
- C_{out} : Number of output activations
- B : batch size

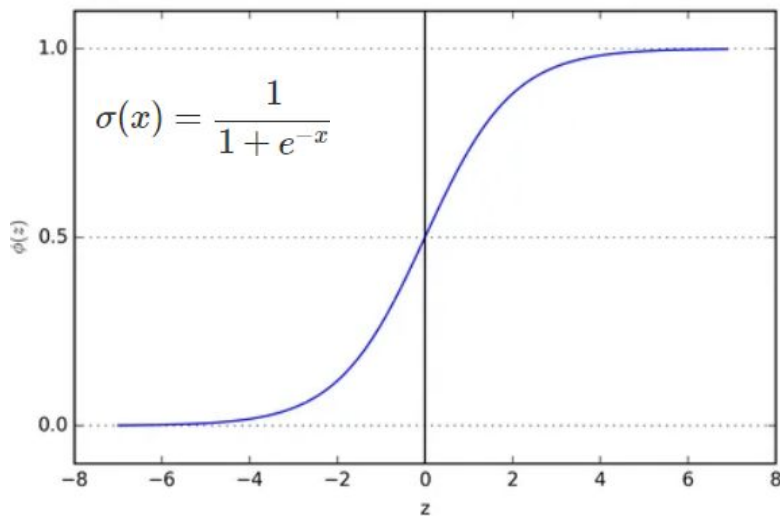


Computational Cost for MLP



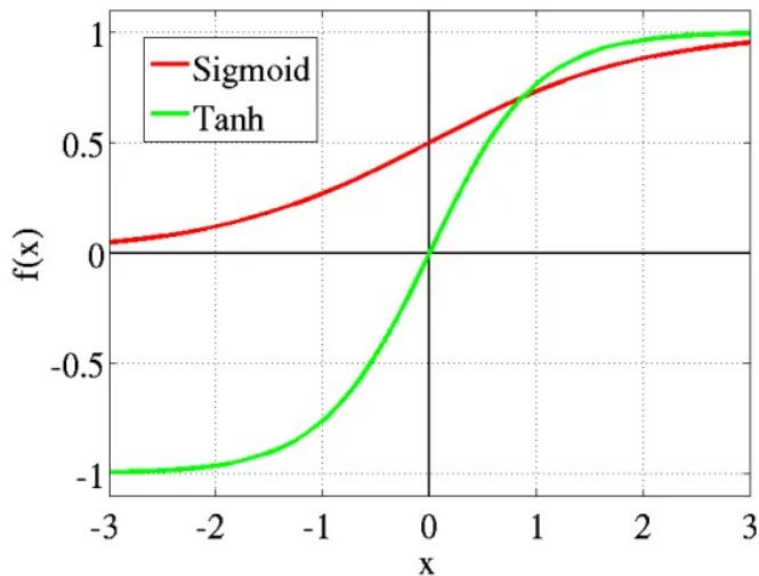
- Number of MACs:
 - $B \times 2 \times 3 = 6B$
- Storage cost:
 - $6 \times 32 = 192$ bits (Weights)
 - $(2B + 3B) \times 32$ bits (Activation)

Sigmoid



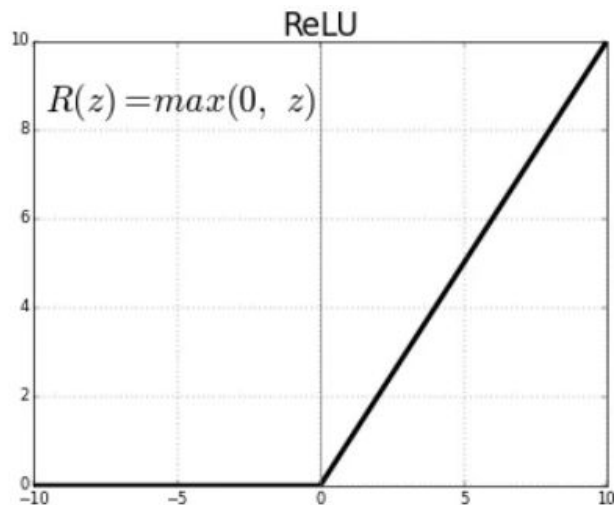
- Function: $\sigma(x) = \frac{1}{1 + e^{-x}}$
- Domain: $(-\infty, \infty)$
- Range: $[0, 1]$
- Differentiable everywhere
- Derivative: $\delta(x)(1 - \delta(x))$

Tanh



- Function: $\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$
- Domain: $(-\infty, \infty)$
- Range: $[-1, 1]$
- Differentiable everywhere
- Derivative: $1 - \tanh^2(x)$

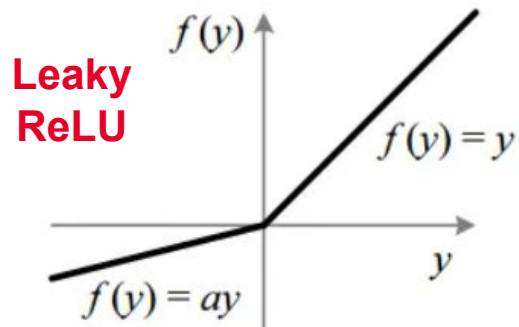
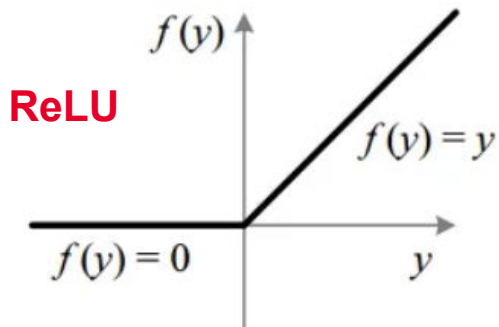
ReLU



$$\text{ReLU}(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

- Domain: $(-\infty, \infty)$
- Range: $[0, \infty]$
- Differentiable everywhere $\begin{cases} 1, & x > 0 \\ 0, & x < 0 \end{cases}$

Leaky ReLU



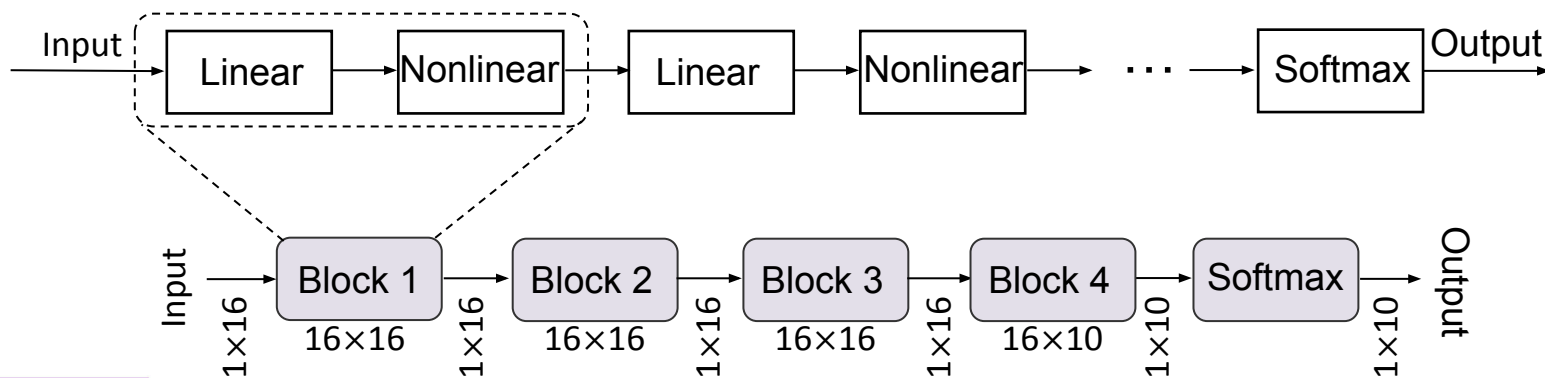
$$\text{Leaky_ReLU}(x) = \begin{cases} x, & \text{if } x \geq 0 \\ ax, & \text{otherwise} \end{cases}$$

- Domain: $(-\infty, \infty)$
- Range: $(-\infty, \infty)$

Softmax

$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

- Domain: $[-\infty, \infty]^N$
- Range: $[0, 1]^N$
- It is a multivariate function



Loss Functions

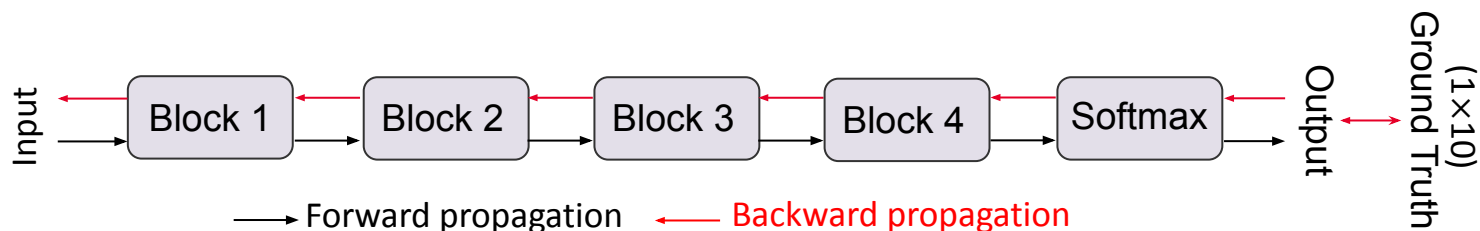
- Loss functions quantify the difference between the DNN output and the ground truth output in the training dataset.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

L2 loss

$$L = -\frac{1}{m} \sum_{i=1}^m (y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i))$$

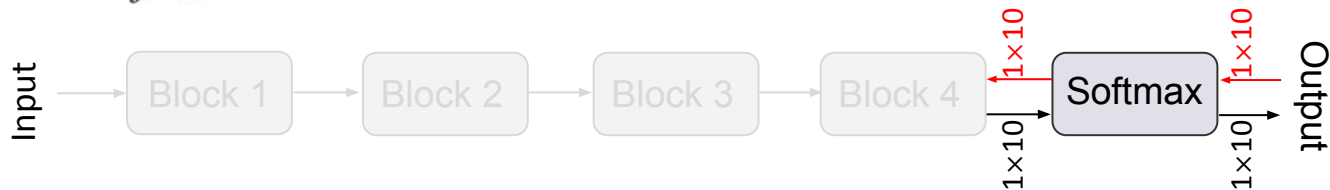
Cross-entropy loss



Softmax

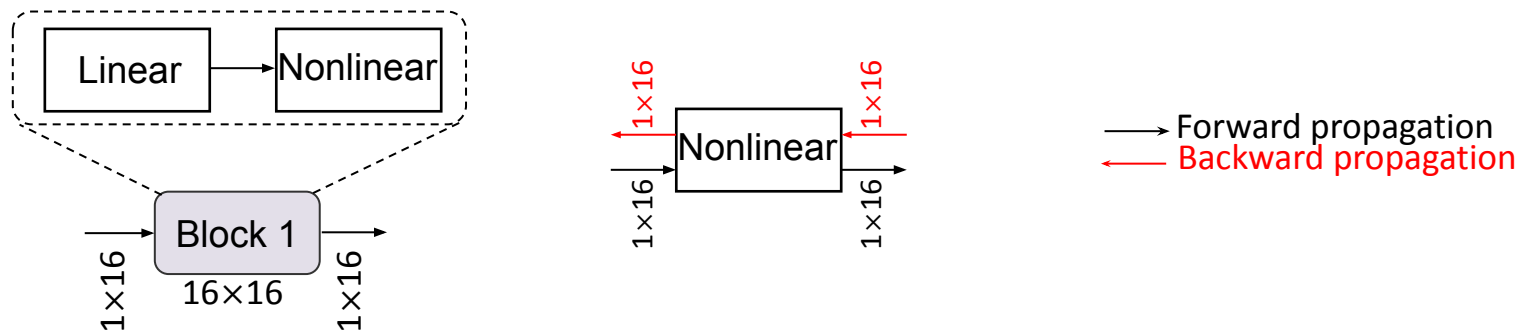
$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

- Domain: $(-\infty, \infty)$
- Range: $[0, 1]$



$$\frac{ds}{dz} = \text{diag}(s) - ss^T \quad \xrightarrow{\text{When } s \text{ has a dimension of 3}} \quad \frac{ds}{dz} = \begin{bmatrix} s_1 - s_1^2 & -s_1 \cdot s_2 & -s_1 \cdot s_3 \\ -s_2 \cdot s_1 & s_2 - s_2^2 & -s_2 \cdot s_3 \\ -s_3 \cdot s_1 & -s_3 \cdot s_2 & s_3 - s_3^2 \end{bmatrix}$$

Backpropagation for Nonlinear Layers

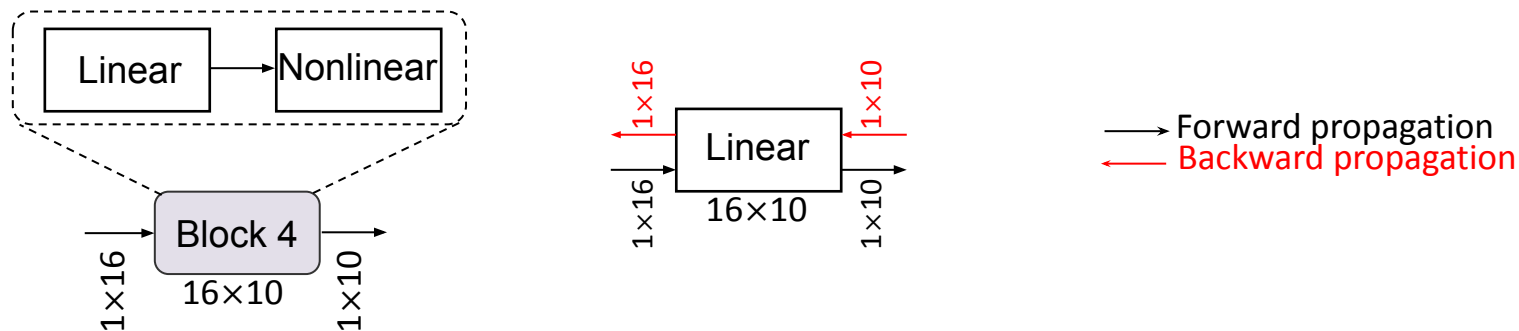


- Due to the elementwise nature, usually the nonlinear layer does not change the input and output shape during both forward and backward passes.

Backpropagation for Nonlinear Layers

- Tanh: $\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$ $1 - \tanh^2(x)$
- ReLU: $\text{ReLU}(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$ $\frac{d\text{ReLU}(x)}{dx} = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$
- Leaky_ReLU: $\text{Leaky_ReLU}(x) = \begin{cases} x, & \text{if } x \geq 0 \\ ax, & \text{otherwise} \end{cases}$ $\frac{d\text{Leaky_ReLU}(x)}{dx} = \begin{cases} 1, & \text{if } x \geq 0 \\ a, & \text{otherwise} \end{cases}$
- Sigmoid: $\sigma(x) = \frac{1}{1 + e^{-x}}$ $\sigma'(x) = \frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$

Backpropagation for Nonlinear Layers

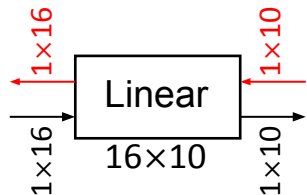


- Due to the elementwise nature, usually the nonlinear layer does not change the input and output shape during both forward and backward passes.

Fully-connected layers (Linear layers)

$$Y = XW + b$$

- X (input activations): $B \times C_{in}$
- Y (output activations): $B \times C_{out}$
- W (weights): $C_{in} \times C_{out}$
- b (bias): $1 \times C_{out}$



$$\frac{dL}{dX} = \frac{dL}{dY} \frac{dY}{dX} = \frac{dL}{dY} W^T$$

Derivative wrt data

$$\frac{dL}{db} = \frac{dL}{dY}$$

Derivative wrt bias

$$\frac{dL}{dW} = X^T \frac{dL}{dY}$$

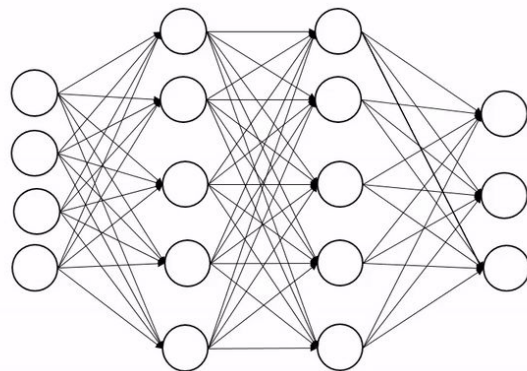
Derivative wrt weight

Weight Decay and Dropout

- The loss function is usually attached with a weight decay loss to penalize the complexity of the function and prevent the overfitting.

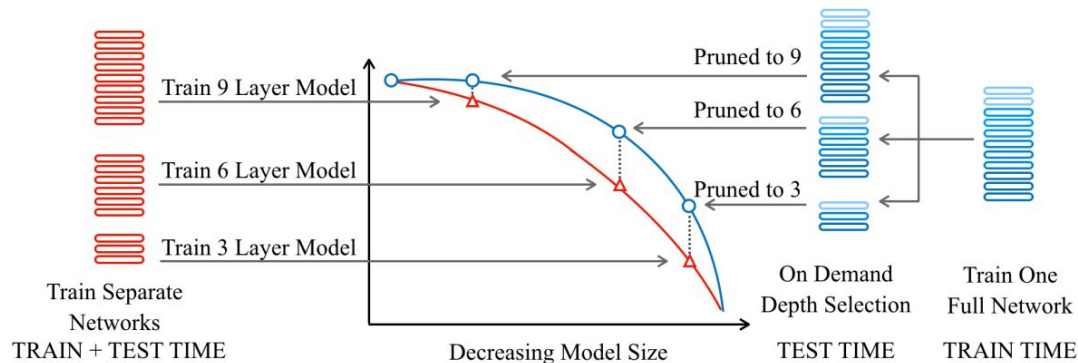
$$L = L + \lambda ||W||^2$$

- Dropout refers to the practice of disregarding certain nodes in a layer at random during training.
- All the nodes will be there during inference.
- Can be used to prevent overfitting and reduce the dependency on any one of a single node.



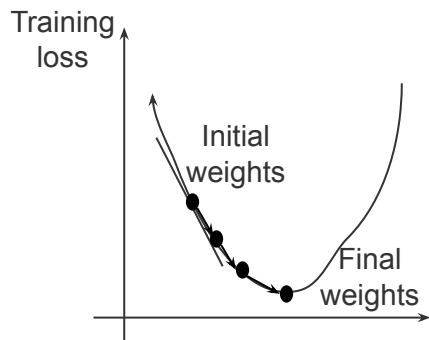
Layer Dropout

- LayerDrop, a form of structured dropout, which has a regularization effect during training and allows for efficient skipping at inference time.
- It is possible to select sub-networks of any depth from one large network without having to finetune them and with limited impact on performance.
- Usually used in transformer.



DNN Training Process

- An optimizer is a crucial element that adjusts DNN parameters during training. Its primary role is to minimize the training loss defined by the loss function.
 - Epoch: The number of times the algorithm runs on the whole training dataset.
 - Batch: The size of block of dataset that is used to update the model weights.
 - Iteration: $\text{total_training_data_size}/\text{Batch}$
 - Learning rate: It is a parameter that provides the model a scale of how much model weights should be updated.



Initialized W for each layer.

For each epoch:

Shuffle the training data.

For each batch in training datasize:

Perform the forward propagation

Compute the loss and weight gradient via backward propagation.

Update the weights

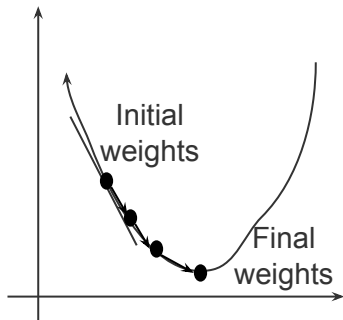
Update the learning rate (if necessary)

Batch, Iteration and Epoch

- A data batch refers to a subset of the entire training dataset used to train the network.
- Iteration refers to a single update of the model's parameters.
- An epoch represents one complete pass through the entire training dataset. Here's what typically happens during an epoch:
 - For example, if you have 1,000 training examples and you use a batch size of 100, it would take 10 iterations to complete one epoch.
- The composition of minibatches typically changes after every epoch during the training of a DNN.

DNN Training Process

- An optimizer is a crucial element that fine-tunes DNN parameters during training. Its primary role is to minimize the model's error or loss function, enhancing performance.
 - Epoch: The number of times the algorithm runs on the whole training dataset.
 - Batch: The size of block of dataset that is used to update the model weights dataset.
 - Iteration: $\text{total_trainingdata_size}/\text{Batch}$
 - Learning rate: It is a parameter that provides the model a scale of how much model weights should be updated.



Initialized W for each layer.

For each epoch:

Shuffle the training data.

For each batch in training datasize:

Perform the forward propagation

Compute the loss and weight gradient via backward propagation.

Update the weights

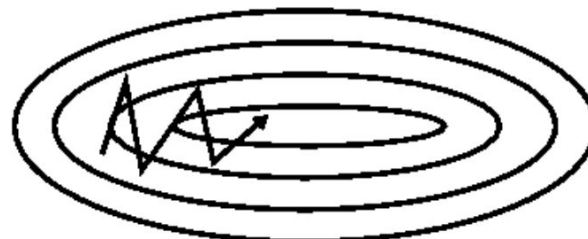
Update the learning rate (if necessary)

Stochastic Gradient Descent (with Momentum)

- $W' = W - \eta dL/dW$
- Due to the significant noise introduced during the SGD process, it is beneficial to stabilize the process.
- $W' = W - \eta g_t$ $g_t \rightarrow s g_{t-1} + (1-s)dL/dW$, s is a hyperparameter between 0 and 1, close to 1.



SGD without momentum



SGD with momentum

RMSProp

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

$$g = [0.02, -0.04, 1.6, -0.01]$$

1.6 will be scaled down with RMSProp

- All operations are elementwise operations.
- When the variance of gradients is high, we scale down the gradient as we want to be more conservative and vice versa.

Adam Optimizer

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

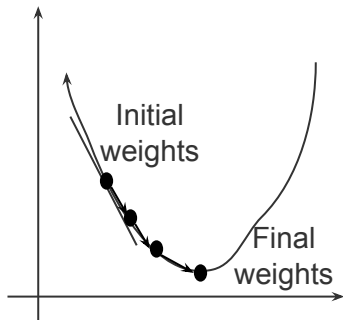
end while

return θ_t (Resulting parameters)

- Combine RMSProp with Momentum SGD.
- By adapting the learning rate during training, Adam converges much more quickly than SGD.

DNN Training Process

- An optimizer is a crucial element that fine-tunes DNN parameters during training. Its primary role is to minimize the model's error or loss function, enhancing performance.
 - Epoch: The number of times the algorithm runs on the whole training dataset.
 - Batch: The size of block of dataset that is used to update the model weights dataset.
 - Iteration: $\text{total_trainingdata_size}/\text{Batch}$
 - Learning rate: It is a parameter that provides the model a scale of how much model weights should be updated.



Initialized W for each layer.

For each epoch:

Shuffle the training data.

For each batch in training datasize:

Perform the forward propagation

Compute the loss and weight gradient via backward propagation.

Update the weights

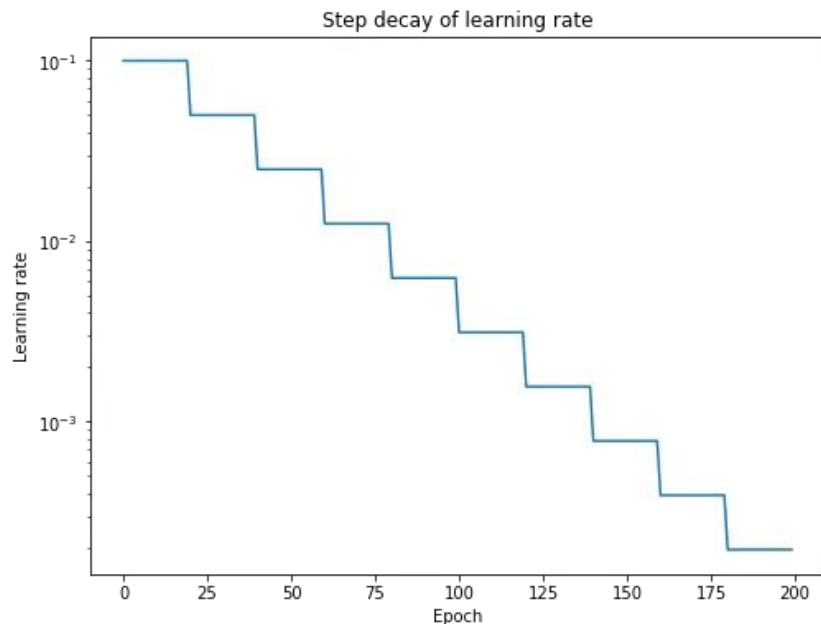
Update the learning rate (if necessary)

Learning Rate Scheduler

- Learning rate η is an important hyperparameter for training the DNNs.
- A large learning rate can help the algorithm to converge quickly. But it can also cause the algorithm to bounce around the minimum without reaching it or even jumping over it if it is too large.
- If the learning rate is too small, the optimizer may take too long to converge or get stuck in a plateau if it is too small.

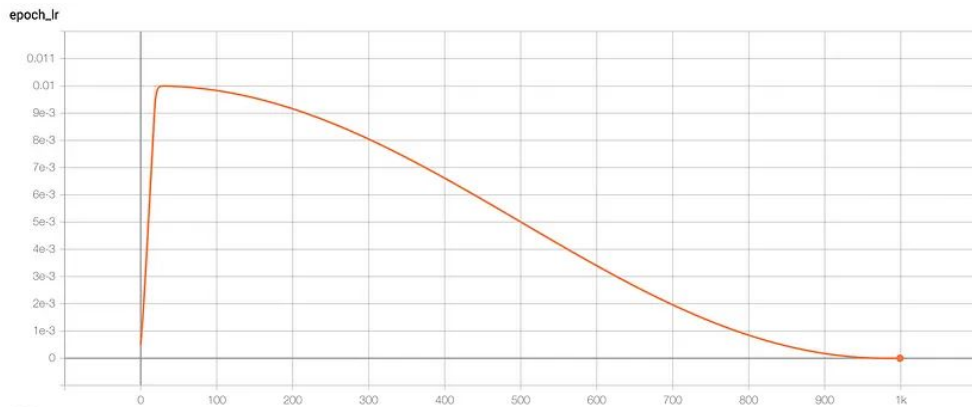
$$W' = W - \eta g_t$$

Multistage Learning Rate



- The learning rate is reduced by a fixed amount after every T epochs.
- Typically, the learning rate is reduced to 10% of its value after every T epochs.
- Widely used in image classification task.

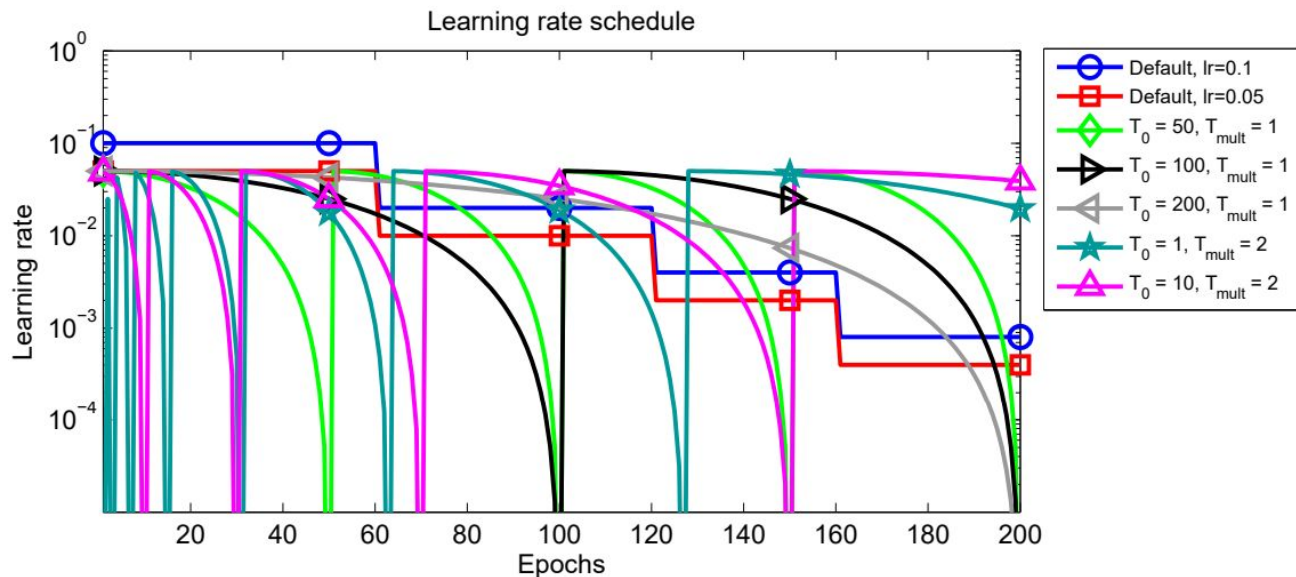
Cosine Learning Rate



- We propose to periodically simulate warm restarts of SGD, where in each restart the learning rate is initialized to some value and is scheduled to decrease.
- Periodic restart can effectively avoid local minima and saddle points during the training.

$$\eta_t = \eta_{min}^i + \frac{1}{2}(\eta_{max}^i - \eta_{min}^i)(1 + \cos(\frac{T_{cur}}{T_i}\pi)),$$

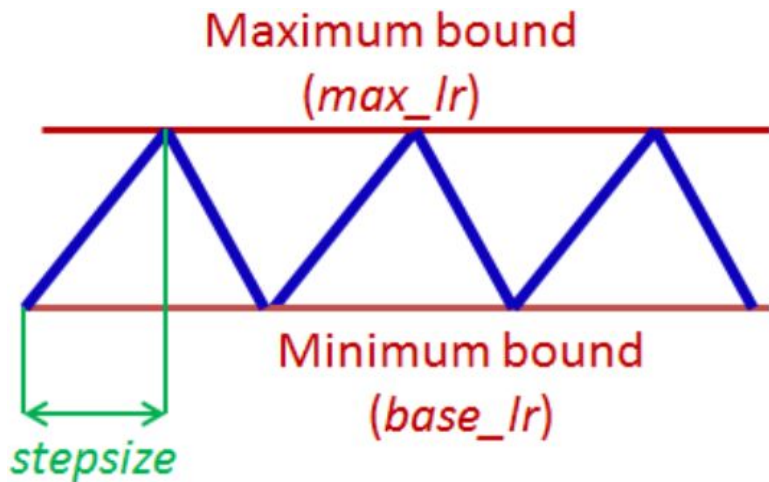
Cosine Learning Rate



- T_{cur} accounts for how many iterations have been performed since the last restart.
- T_{cur} is updated at each iteration t .
- The SGD is restarted once T_i epochs are performed, where i is the index of the run.
- T_i may increase with i .

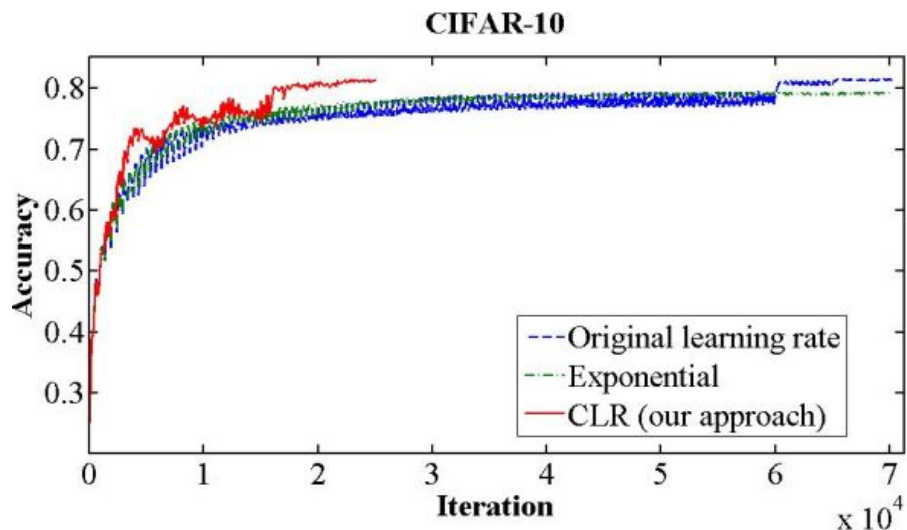
$$\eta_t = \eta_{min}^i + \frac{1}{2}(\eta_{max}^i - \eta_{min}^i)(1 + \cos(\frac{T_{cur}}{T_i}\pi)),$$

Cyclical Learning Rate



- Increasing the learning rate might have a short term negative effect and yet achieve a longer term beneficial effect.

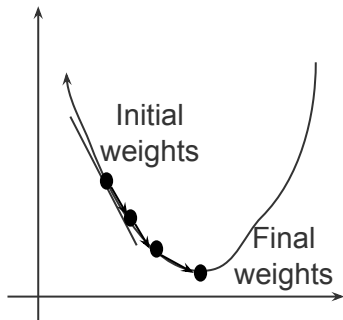
Cyclical Learning Rate



- The red curve shows the result of training with cyclical learning rate achieves the shortest convergence time.

DNN Training Process

- An optimizer is a crucial element that fine-tunes DNN parameters during training. Its primary role is to minimize the model's error or loss function, enhancing performance.
 - Epoch: The number of times the algorithm runs on the whole training dataset.
 - Batch: The size of block of dataset that is used to update the model weights dataset.
 - Iteration: $\text{total_trainingdata_size}/\text{Batch}$
 - Learning rate: It is a parameter that provides the model a scale of how much model weights should be updated.



Initialized W for each layer.

For each epoch:

Shuffle the training data.

For each batch in training datasize:

Perform the forward propagation

Compute the loss and weight gradient via backward propagation.

Update the weights

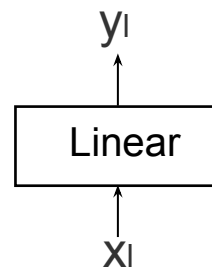
Update the learning rate (if necessary)

DNN Initialization: Kaiming Initialization

- Kaiming initialization is designed for modern DNN that uses ReLU.

$$W \sim \mathcal{N}\left(0, \frac{2}{n^l}\right)$$

- Target: ensure the activation variance is the same across different layers.
- Assumption:
 - ReLU activation.
 - Weight is normally distributed with mean of zero.
 - Weight and activations are independent.



Derivation

$$\mathbf{y}_l = \mathbf{W}_l \mathbf{x}_l + \mathbf{b}_l$$

Assume \mathbf{W}_l has a shape of m_l by n_l , and \mathbf{x}_l has a size of $n_l \times 1$, then \mathbf{y}_l has a size of $m_l \times 1$.

For each element $y_{l,i}$ of \mathbf{y}_l , its variance $\text{var}(y_{l,i}) = \text{var}\left(\sum_{j=1}^{n_l} W_{l,i,j} x_{l,j}\right) = n_l \text{var}(W_{l,i,j} x_{l,j})$

Assume each pair of $W_{l,i,j}$ and $x_{l,j}$ are independent random variable, then we have:

$$\text{var}(W_{l,i,j} x_{l,j}) = E(W_{l,i,j}^2 x_{l,j}^2) - E^2(W_{l,i,j} x_{l,j}) = E(W_{l,i,j}^2) E(x_{l,j}^2) - E^2(W_{l,i,j}) E^2(x_{l,j})$$

Assume $W_{l,i,j}$ follows a normal distribution with mean of 0, that is $E(W_{l,i,j}) = 0$, then:

$$\text{var}(W_{l,i,j} x_{l,j}) = \text{var}(W_{l,i,j}) E(x_{l,j}^2)$$

$$\text{var}(y_{l,i}) = n_l \text{var}(W_{l,i,j} x_{l,j}) = n_l \text{var}(W_{l,i,j}) E(x_{l,j}^2)$$

Derivation

Let see how $E(x_{l,i}^2)$ is related to the variance of $y_{l-1,j}$, where $x_{l,i} = \text{ReLU}(y_{l-1,j})$

$$E(x_{l,i}^2) = E(\text{ReLU}^2(y_{l-1,j}))$$

Then we have: $E(\text{ReLU}(y_{l-1,j})^2)$

$$= E(\text{ReLU}(y_{l-1,j})^2 | y_{l-1,j} > 0)P(y_{l-1,j} > 0) + E(\text{ReLU}(y_{l-1,j})^2 | y_{l-1,j} < 0)P(y_{l-1,j} < 0)$$

$$= E(\text{ReLU}(y_{l-1,j})^2 | y_{l-1,j} > 0)P(y_{l-1,j} > 0) = 0.5E(y_{l-1,j}^2) = 0.5\text{var}(y_{l-1,j})$$

Therefore, we have: $E(x_{l,i}^2) = 0.5\text{var}(y_{l-1,j})$

Given this, we have:

$$\text{var}(y_{l,i}) = n_l \text{var}(W_{l,i,j}) E(x_{l,i}^2) = 0.5 n_l \text{var}(W_{l,i,j}) \text{var}(y_{l-1,j})$$

Derivation

$$\text{var}(y_{l,i}) = \left(\prod_{s=2}^l 0.5n_s \text{var}(W_{s,i,j}) \right) \text{var}(y_{1,j})$$

In order to ensure the variance of y does not change, we have to make sure:

$$\text{var}(W_{s,i,j}) = \frac{2}{n_s} \quad W_{s,i,j} \sim \mathcal{N}\left(0, \sqrt{\frac{2}{n_s}}\right)$$