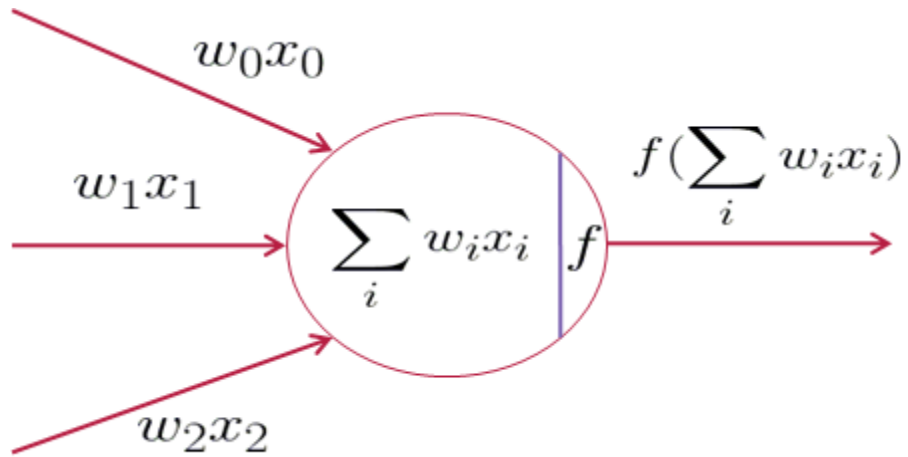


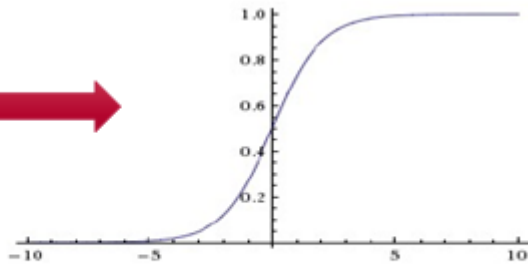
Backpropagation and Training

C. V. Jawahar
IIIT Hyderabad

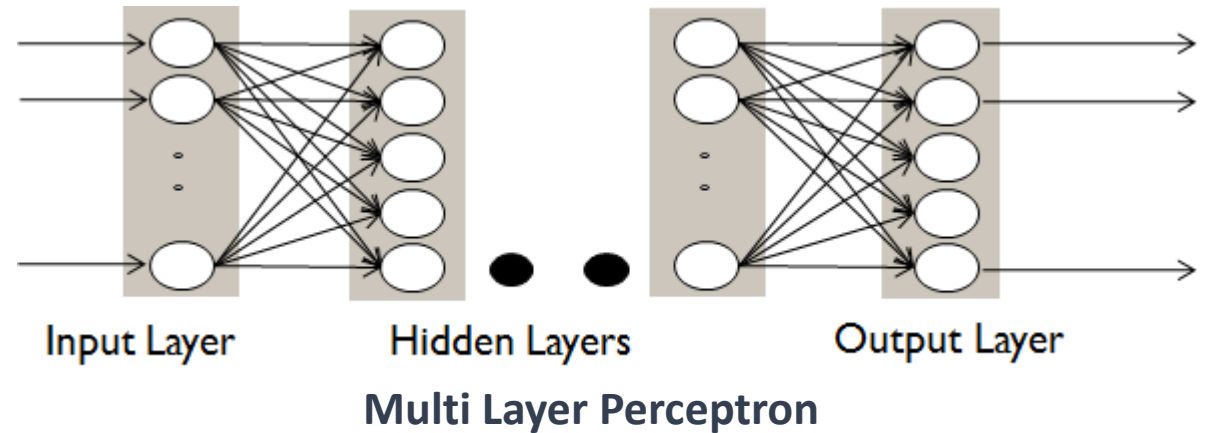
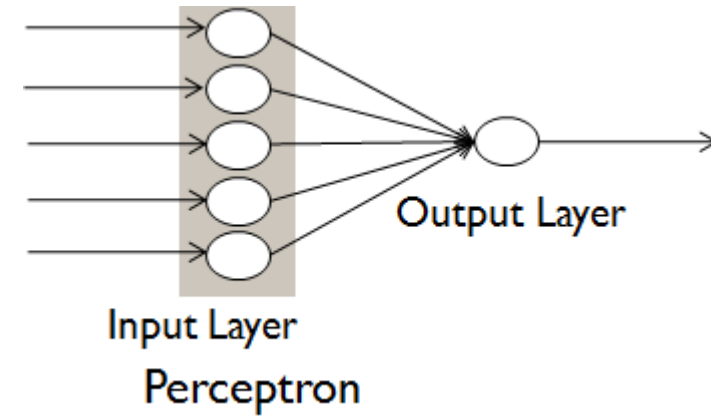
Neuron, Perceptron and MLP



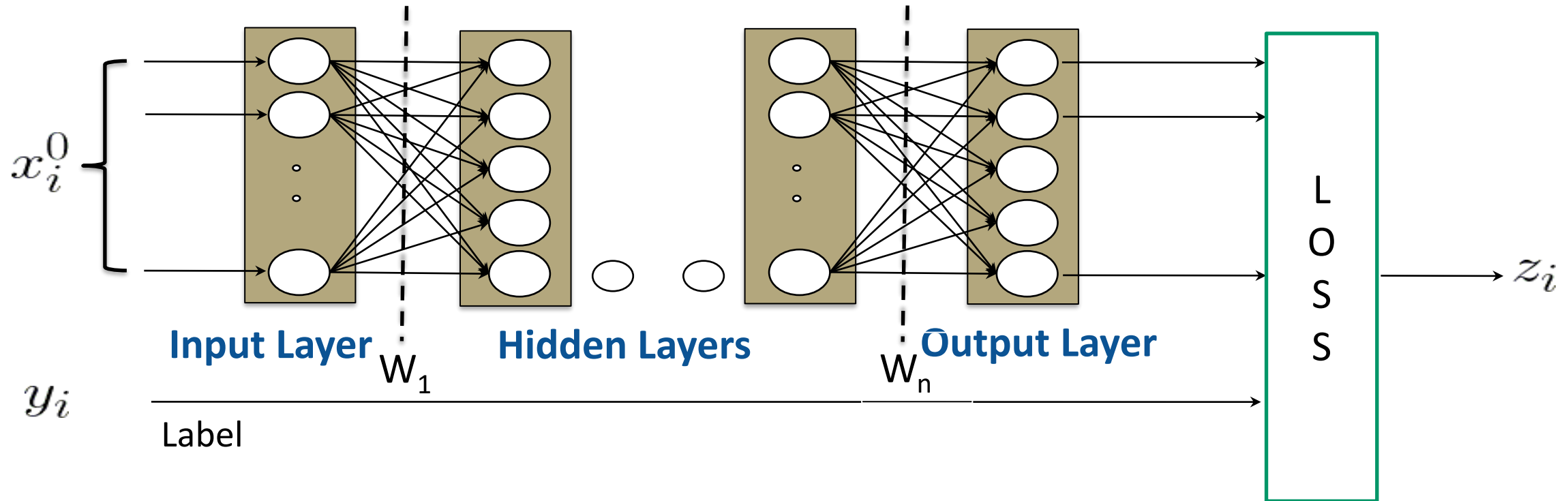
f



E.g. Sigmoid Activation Function



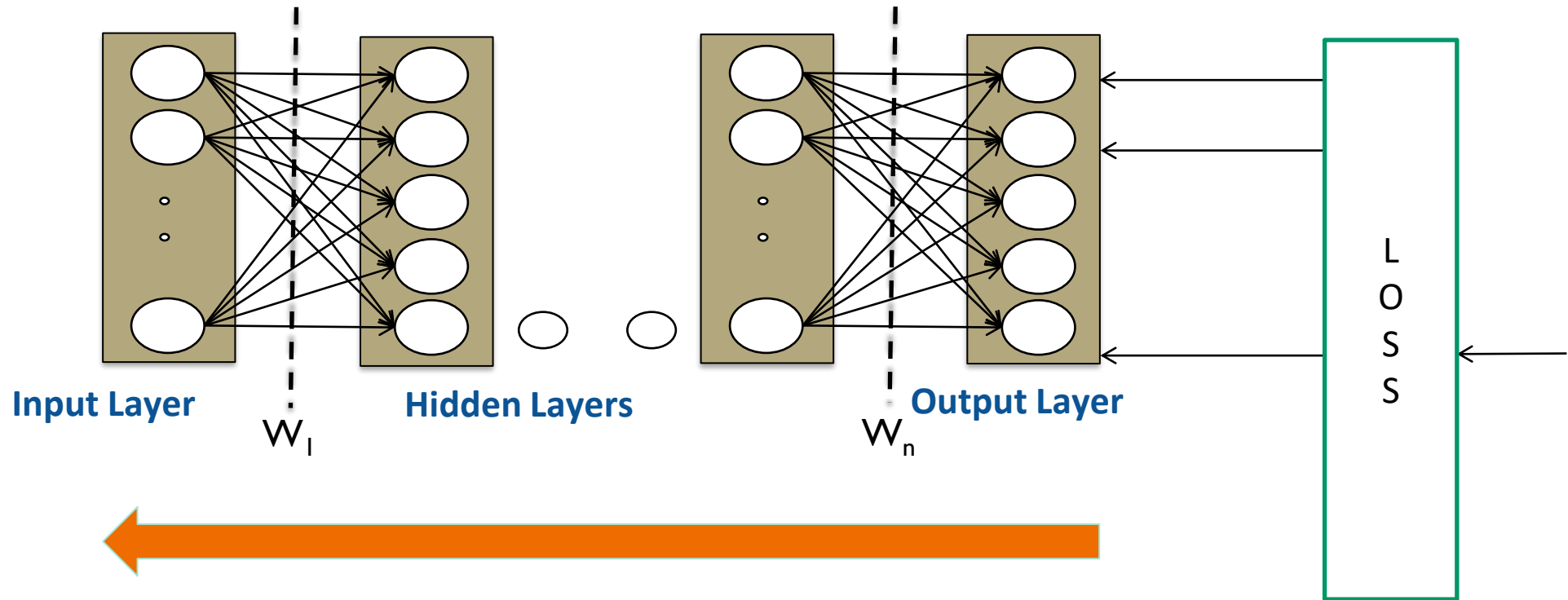
Loss or Objective



Objective: Find out the best parameters which will minimize the loss.

$$W^* = \arg \min_W \sum_{i=1}^N L(x_i^n, y_i; W) \longrightarrow \text{Weight vector}$$
$$z_i = \frac{1}{2} \| x_i^n - y_i \|_2^2 \quad \text{E.g. Squared Loss}$$

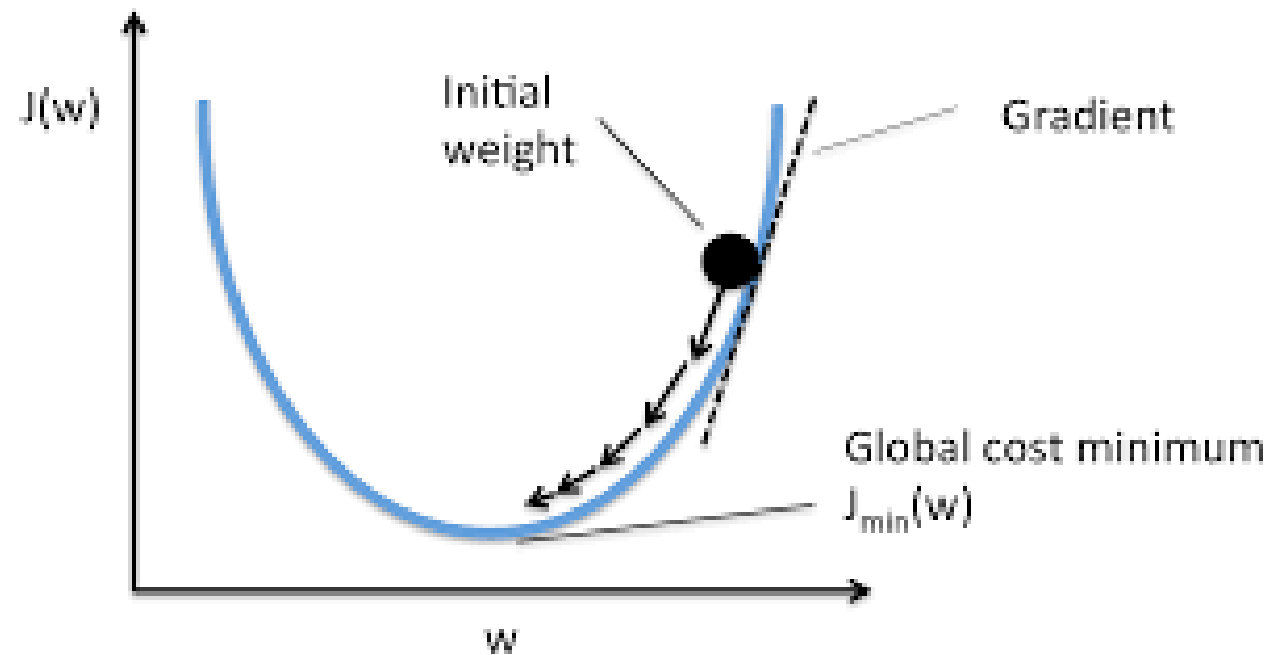
Back Propagation



Solution: Iteratively update W along the direction where loss decreases.

Each layer's weights are updated based on the derivative of its output w.r.t. input and weights

Gradient Descent



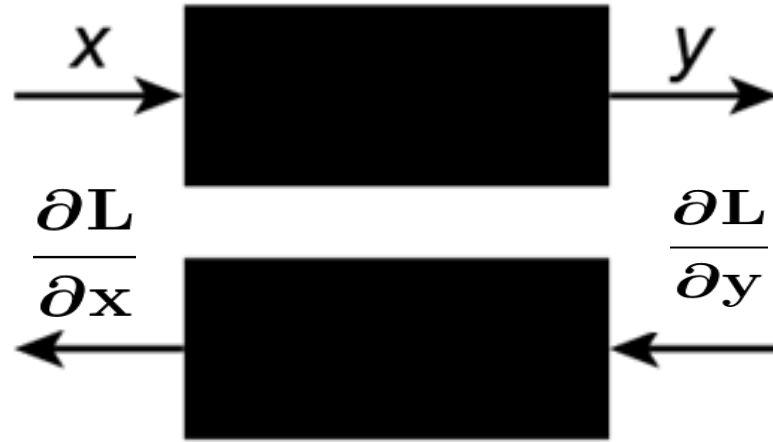
Neural Network Training

- **Step 1:** Compute loss for samples (on mini-batch) [F-Pass]
- **Step 2:** Compute gradients w.r.t parameters and update [B-Pass]



$$W \leftarrow W - \eta \frac{\partial L}{\partial W}$$

Chain Rule



Given $y(x)$ and $\frac{\partial L}{\partial y}$
What is $\frac{\partial L}{\partial x}$?

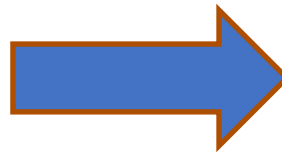


$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

Chain Rule



Given $\frac{\partial y}{\partial x}$ and $\frac{\partial L}{\partial y}$



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

What is $\frac{\partial L}{\partial x}$?

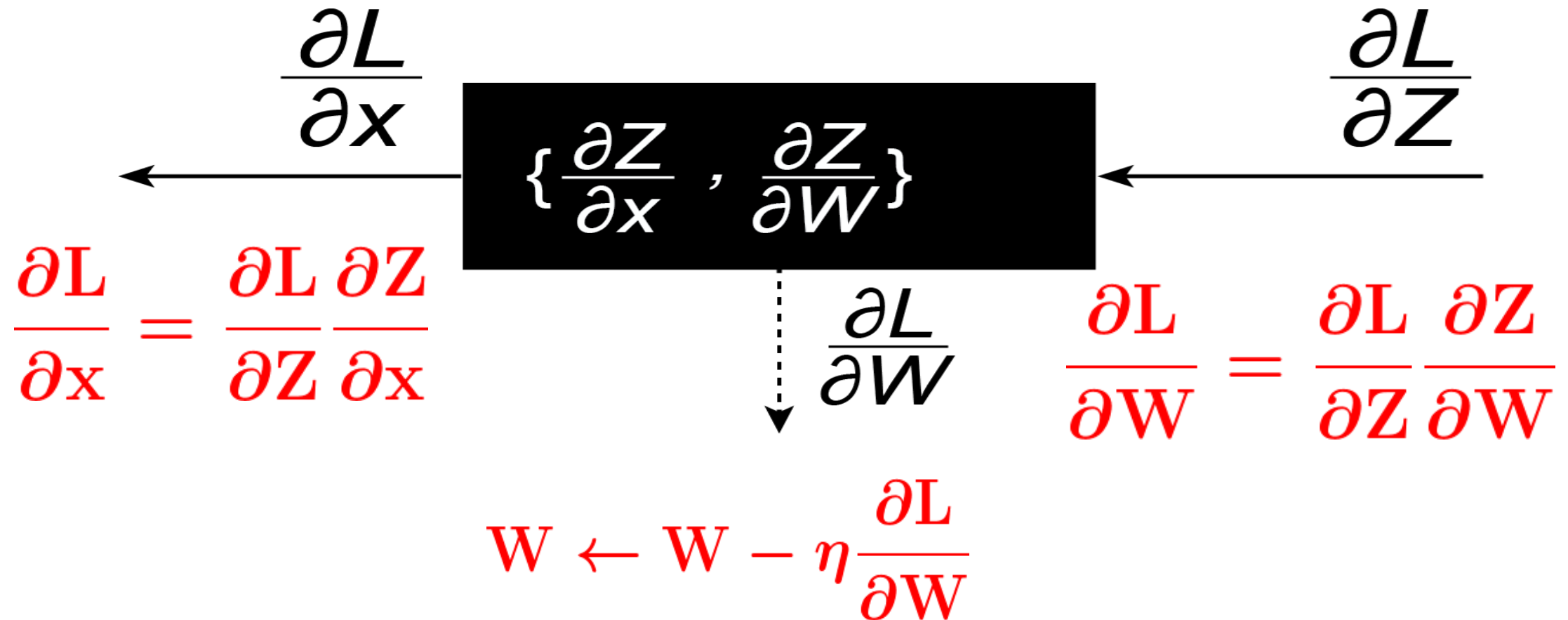
For each block/parameters, we only need to find $\frac{\partial y}{\partial x}$

Key Computation: Forward-Propagation



W is the parameter, say the weights within the box.

Key Computation: Backward-Propagation



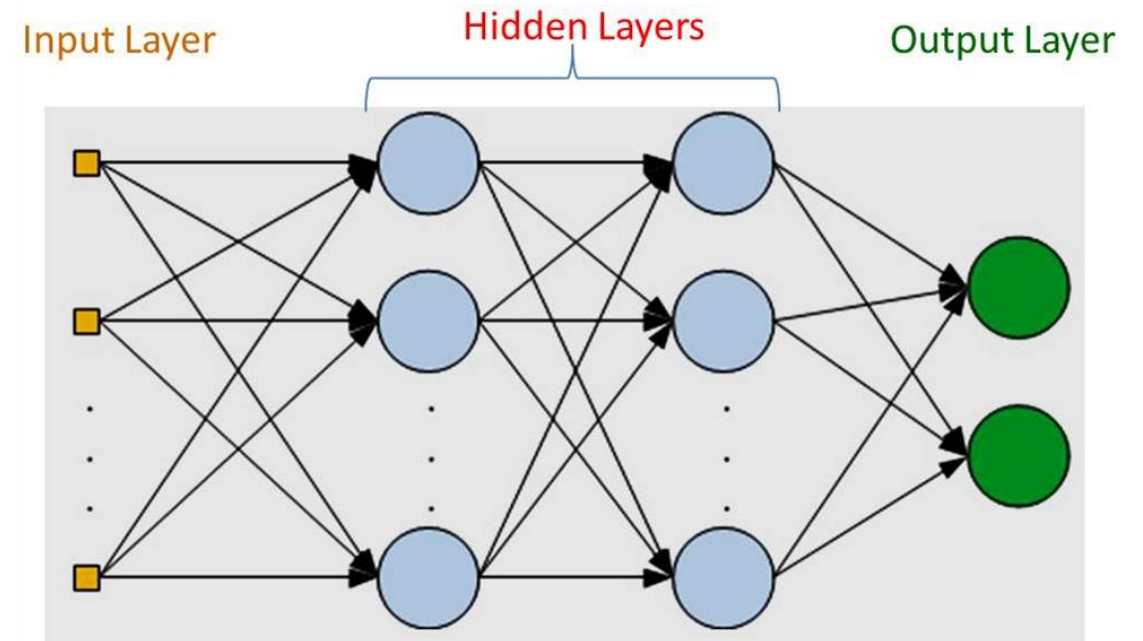
Back Propagation for MLP

- Two computational blocks/steps

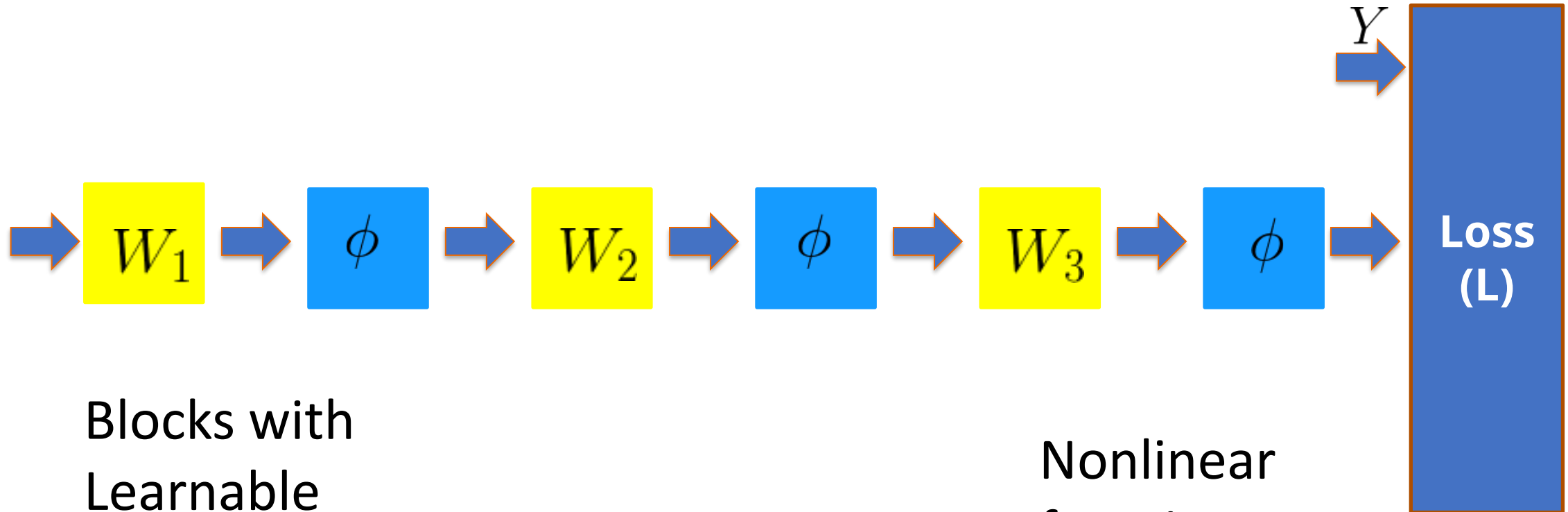
$$y = Wx$$

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

- In either case we can compute easily. $\frac{\partial y}{\partial x}$



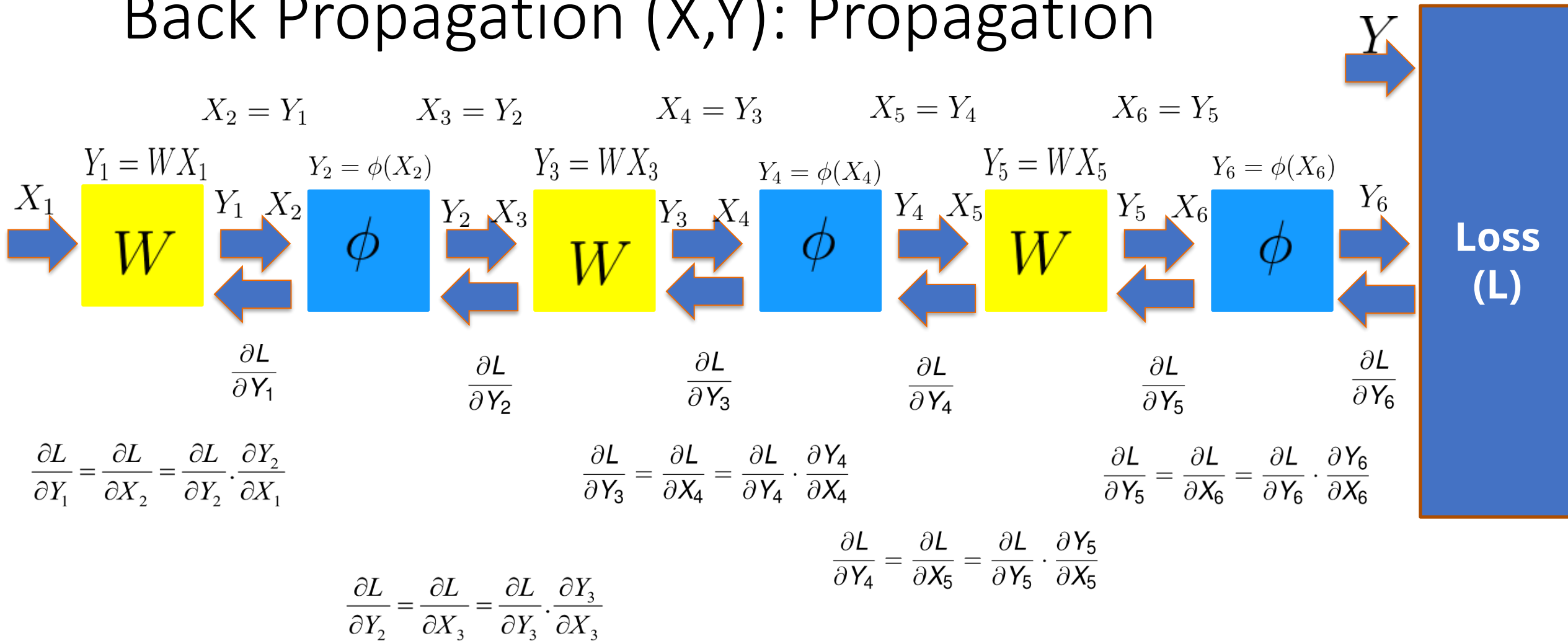
A simpler view point



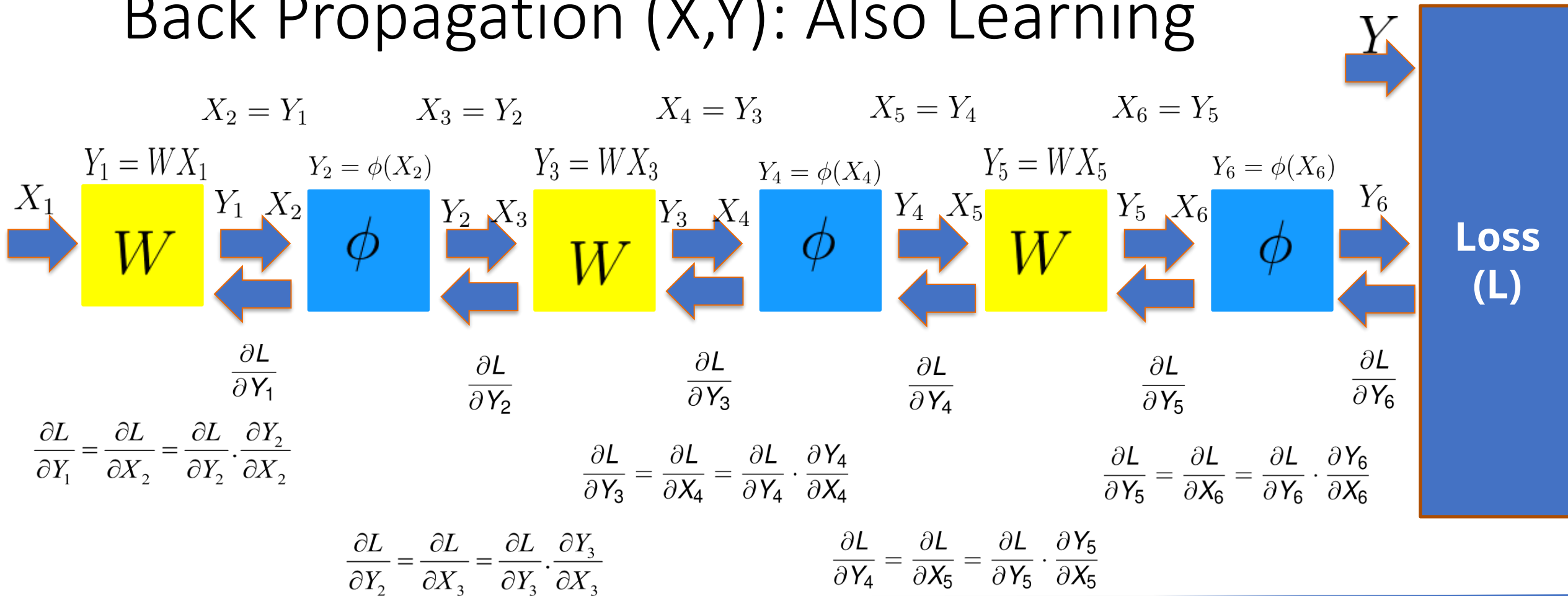
Blocks with
Learnable
parameters
Matrix
Multiplication

Nonlinear
functions
(often non learnable)

Back Propagation (X,Y): Propagation



Back Propagation (X,Y): Also Learning



$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y_1} \cdot \frac{\partial Y_1}{\partial W} \qquad \frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial W} \qquad \frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y_5} \cdot \frac{\partial Y_5}{\partial W} \qquad W^{n+1} = W^n - \eta \frac{dL}{dW}$$

Summary

- **Step 0:**
 - Initialize the Network, weights
- **Step 1:**
 - Do forward pass for a batch of randomly selected samples.
 - Predict outputs with the existing weights.
- **Step 2:**
 - Compute Loss for the set of samples.

Summary

- Step 3:
 - Update all the weights using gradient descent.

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial \mathbf{L}}{\partial \mathbf{W}}$$

- Step 4:
 - Repeat all steps till the Loss is less than a threshold.

Loss Functions

- **For Classification**

- Binary Cross Entropy
- Hinge Loss

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

$$L = \max(0, 1 - y * f(x))$$

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

- **For Regression**

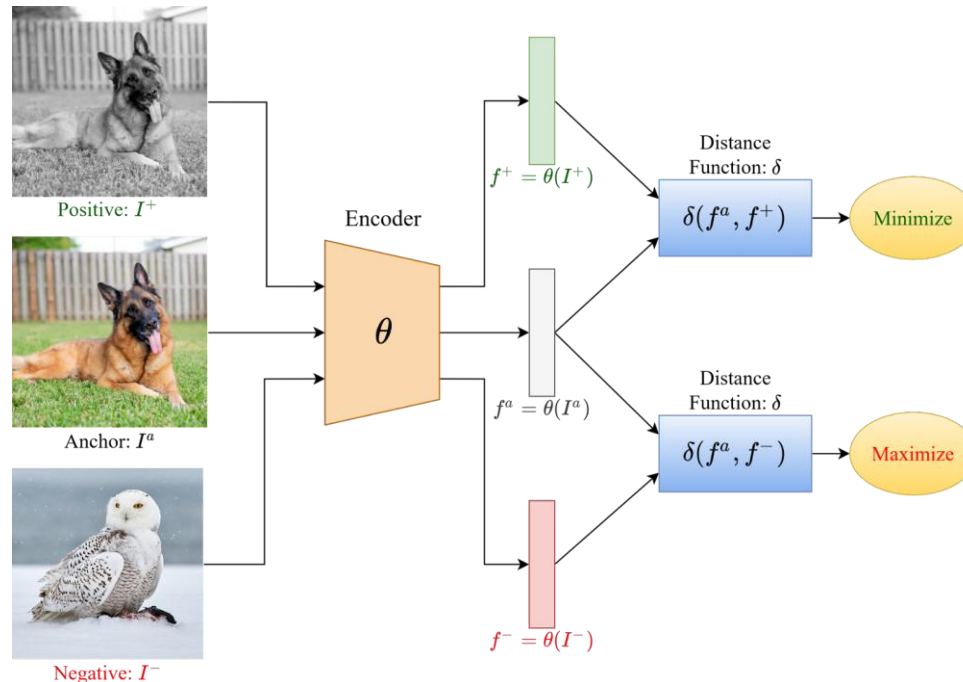
- Mean Square Error / Quadratic Loss / L2 Loss
- Mean Absolute Error / L1 Loss

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

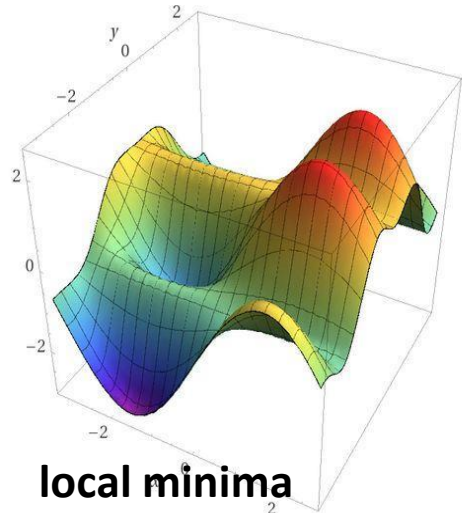
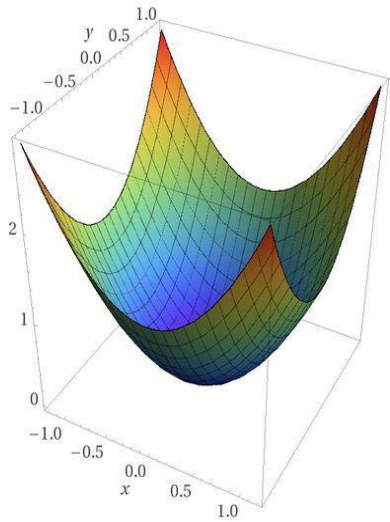
$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

Loss Functions

- For Comparing Structured Data
 - Strings
 - Graphs
 - Trees
- For Comparing Multiple Samples
 - Contrastive Learning
 - Siamese Learning
 - Triplet Loss

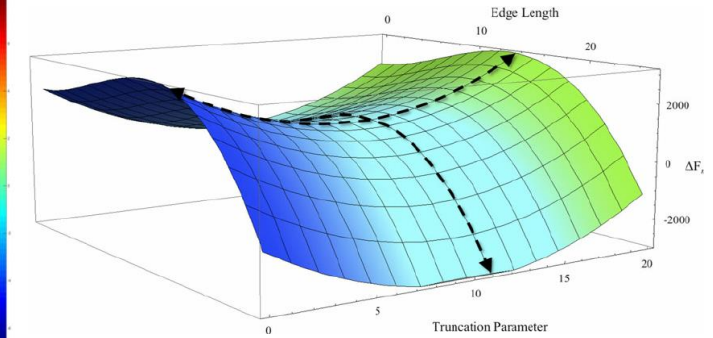
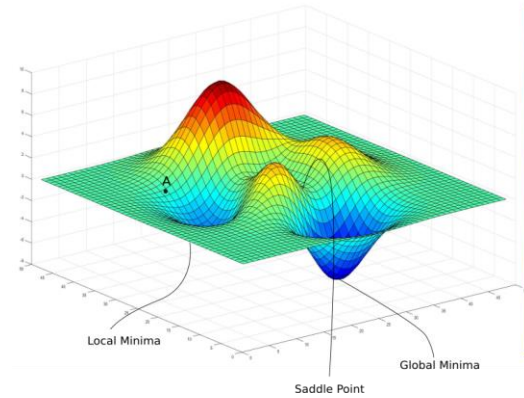


Why Challenging?



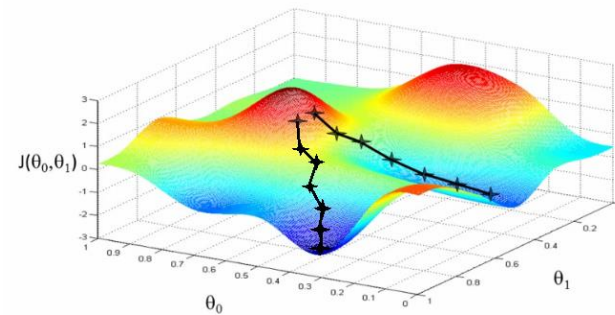
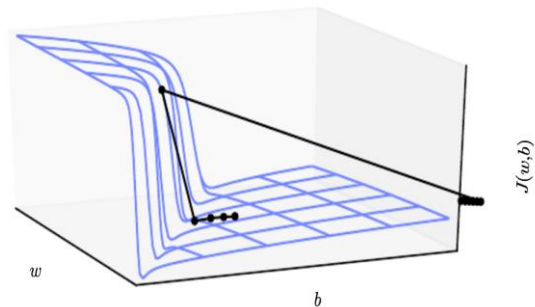
local minima

$$\nabla J = 0$$



Saddle Point

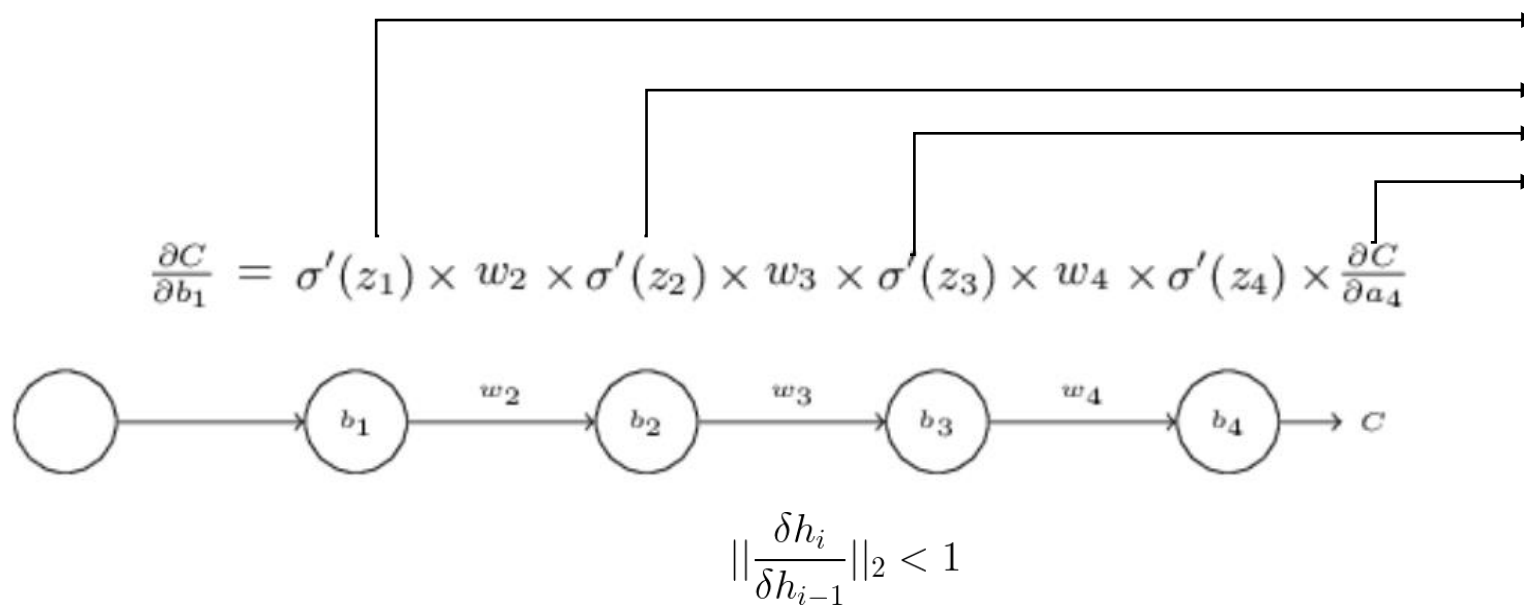
Cliffs



Saddle Point

Vanishing and Exploding Gradients

The Vanishing Gradient Problem:



Solution: Gradient clipping

When gradient is too small, the repetitive multiplication results in vanishing gradient

When gradient is too high, the repetitive multiplication results in exploding gradient

Solution: Better Estimate of the Gradient

Ideal optimizer:

- Finds minimum fast and reliably well
- Doesn't get stuck in local minima, saddle points, or plateau region

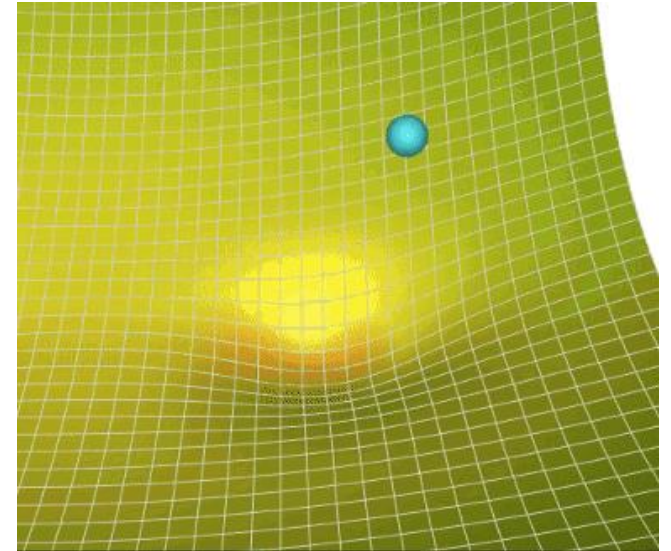
$$w_{t+1} = w_t - \alpha * \frac{\delta L}{\delta w_t}$$

Vanilla Gradient Descent: One step for the entire dataset

Stochastic Gradient Descent: One step for each stochastically chosen sample

Mini-batch Gradient Descent: One step for each mini-batch of samples chosen stochastically

Best of both worlds!

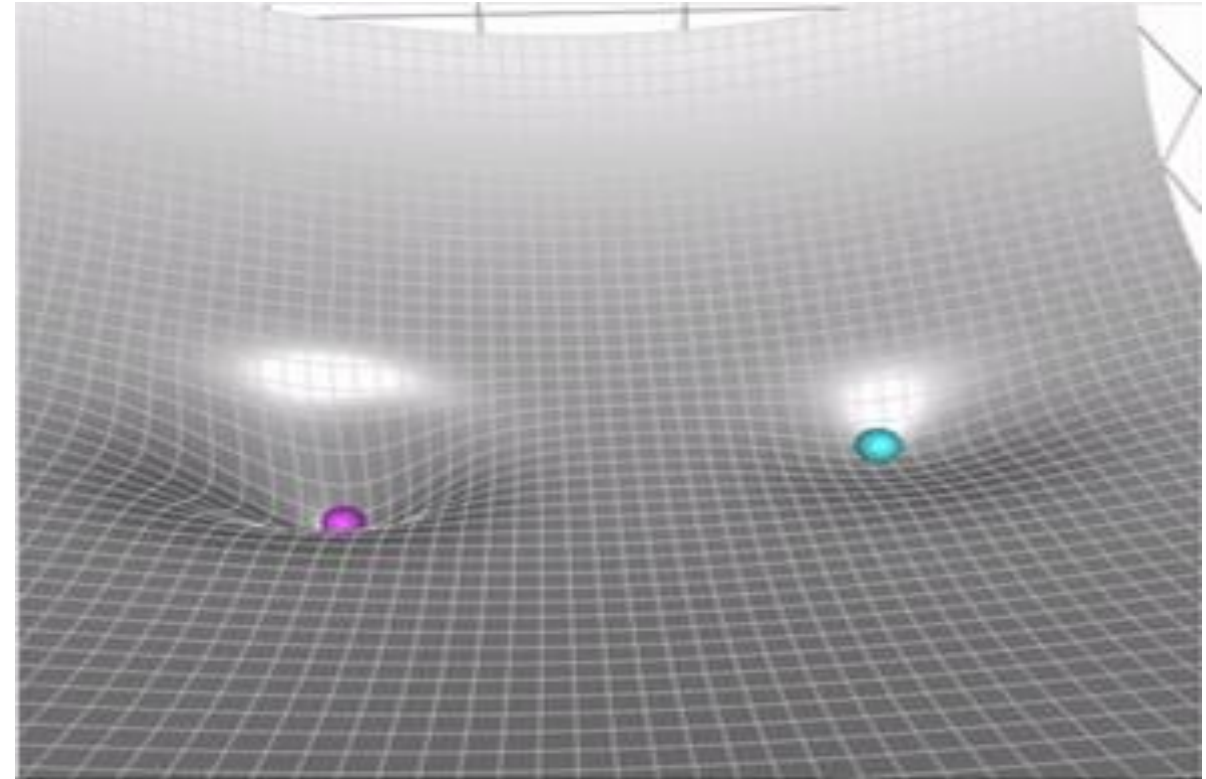


```
1 import torch.optim as optim
2
3 optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0)
```

Solution: Better Update Equation

Gradient Descent with Momentum

Intuitively, this helps us to come out of local minima



$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

$$\theta = \theta - v_t$$

Momentum

```
1 import torch.optim as optim
2
3 optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
```

Solution: Better Optimizers

Other Variants:

- RMSprop
- Adagrad
- **Adam**
- Adadelta
- etc.

```
1 import torch.optim as optim
2
3 optimizer = optim.RMSprop(model.parameters(), lr=0.01, alpha=0.99)
```

```
1 import torch.optim as optim
2
3 optimizer = optim.Adagrad(model.parameters(), lr=0.01)
```

```
1 import torch.optim as optim
2
3 optimizer = optim.Adam(model.parameters(), lr=0.01, betas=(0.9, 0.999))
```

```
1 import torch.optim as optim
2
3 optimizer = optim.Adadelta(model.parameters(), lr=0.01, rho=0.9)
```

Simplified View: Variable Learning Rates for Features/Dimensions

Solution: Better Weight Initialization

- All zero initialization
 - Initializing all the weights with zeros leads the neurons to learn the same features during training.
- Random initialization
 - Gaussian
 - Xavier
 - uniform
 - normal
 - Kaiming
 - uniform
 - normal

```
1  import torch
2
3  w = torch.empty(3, 5)
4  torch.nn.init.kaiming_uniform_(w, mode='fan_in', nonlinearity='relu')
```


Regularization of the Network

- Input Level
 - Data Augmentation
- Activation level
 - Dropout
 - Dropconnect
- Feature statistics level
 - Batch normalization
 - Layer normalization
 - Group normalization
- Decision level
 - Ensemble
- Constraining network weights
 - ℓ_1 norm, ℓ_2 norm
- Terminating early based of the performance on validation set
 - Early stopping

Summary

- Data Normalization
- Data Augmentation
- Weight Initialization
- Optimization Algorithms
- Regularizer
- Batch Norm
- Dropout
- Better Learning/Optimization
- Better Generalization/Regularization
- Better Loss Functions
- BP as a General Algorithm
- Abstraction of layers and ease in defining layers
- Good implementations

Thank you!

Questions?