# Johns Hopkins Engineering

## Applied Machine Learning for Mechanical Engineers

Machine Learning Fundamentals, Part 1, D

# Learning Algorithms
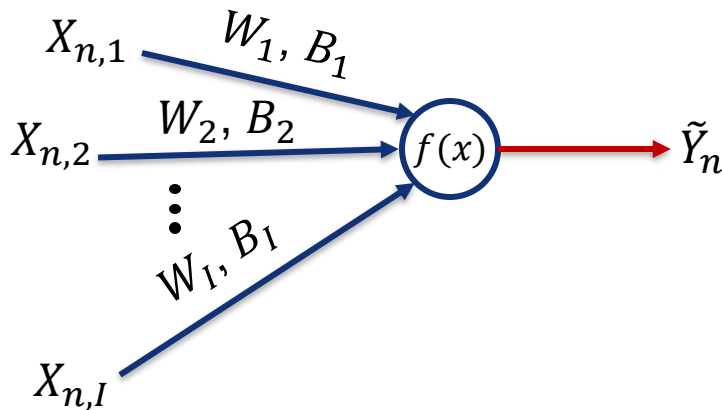
- By the end of this lecture you will be able to:

  o Describe gradient descent

  o Describe backpropagation

  o Describe universal approximation theorem

# Learning Algorithms

- Loss functions

  - Mean squared error

  - Mean absolute error

  - Cross entropy

# Learning Algorithms

- **Loss functions**

    o Loss function in multi-layer neural network; extremely complicated!

$$X_{n,1} \xrightarrow{W_1, B_1}$$
$$X_{n,2} \xrightarrow{W_2, B_2} f(x) \longrightarrow \tilde{Y}_n$$
$$\vdots$$
$$X_{n,I} \xrightarrow{W_I, B_I}$$

# Learning Algorithms

- Loss functions

  ○ Loss function in perceptron

$$l(p) = \frac{1}{2N} \sum_{n=1}^{N} \left\| Y_n - \tilde{Y}_n \right\|^2$$

Estimated

Actual/real

where

$$\tilde{Y}_n = f(x) = f\left( \sum_{i=1}^{I} X_{n,i} W_i + B_i \right)$$

$X_{n,1}$  $W_1, B_1$

$X_{n,2}$  $W_2, B_2$

$f(x)$  $\tilde{Y}_n$
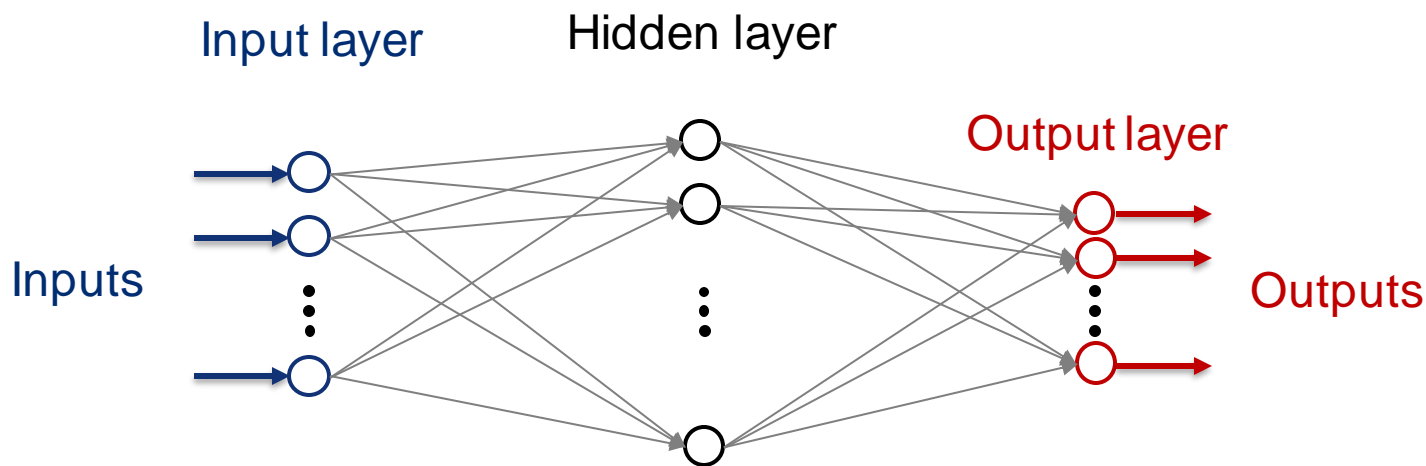
$W_I, B_I$

$X_{n,I}$

$p$ includes $I$ weights and $I$ biases to be optimized

# Learning Algorithms

- Loss functions

    o 10 inputs, 100 hidden neurons, and 5 outputs; 3000 weights and biases!



Input layer          Hidden layer

Output layer

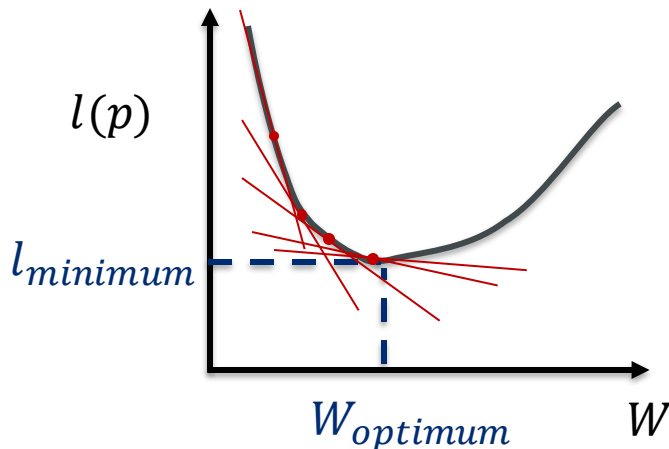Inputs

Outputs

# Learning Algorithms

- Optimizers or learning algorithms

  o Gradient descent

  o Stochastic gradient descent

  o Conjugate gradient descent

  o Adaptive moment estimation (ADAM)

# Learning Algorithms

- Gradient descent and backpropagation

  ○ Gradient descent is simply gradient (slope in 2d) calculation over loss function in a proper step, called learning rate (similar to penalty constant in mathematical optimization).



Small learning rate
  ✓ Time consuming
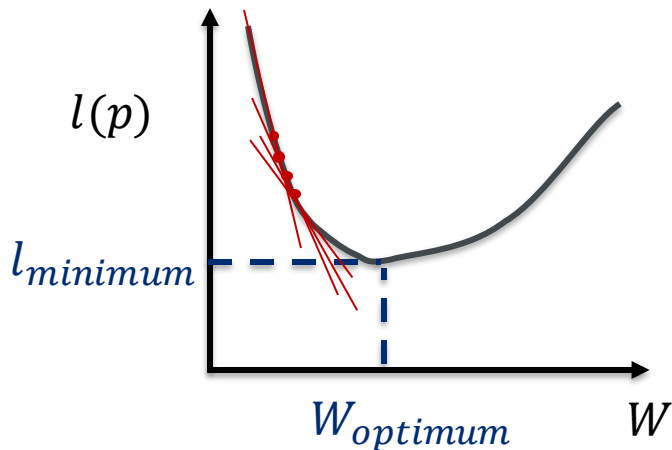Large learning rate
  ✓ Higher error, less accurate
Adoptive learning rate
  ✓ Start with large magnitude and make it smaller as slope decreases

# Learning Algorithms

- Gradient descent and backpropagation

  o Gradient descent is simply gradient (slope in 2d) calculation over loss function in a proper step, called learning rate (similar to penalty constant in mathematical optimization).



Small learning rate
  ✓ Time consuming
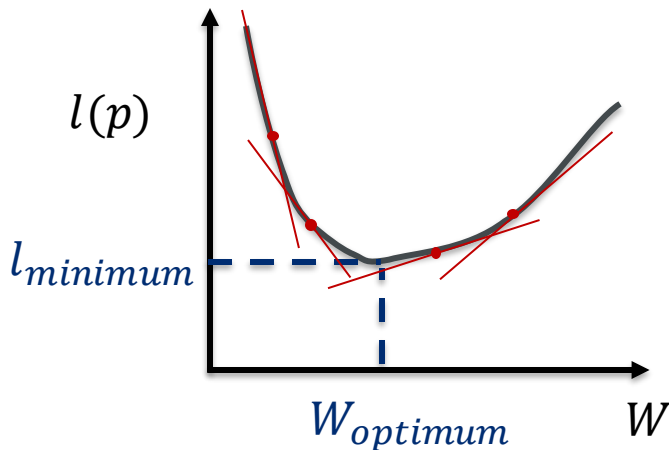Large learning rate
  ✓ Higher error, less accurate
Adoptive learning rate
  ✓ Start with large magnitude and make it smaller as slope decreases

# Learning Algorithms

- Gradient descent and backpropagation

  o Gradient descent is simply gradient (slope in 2d) calculation over loss function in a proper step, called learning rate (similar to penalty constant in mathematical optimization).



Small learning rate
  ✓ Time consuming
Large learning rate
  ✓ Higher error, less accurate
Adoptive learning rate
  ✓ Start with large magnitude and make it smaller as slope decreases

# Learning Algorithms

- **Gradient descent and backpropagation**

  - Gradient descent is simply gradient (slope in 2d) calculation over loss function in a proper step, called learning rate (similar to penalty constant in mathematical optimization).



Small learning rate
  ✓ Time consuming
Large learning rate
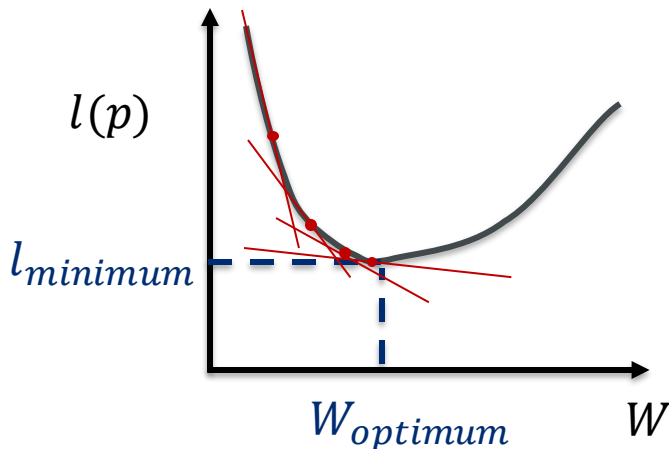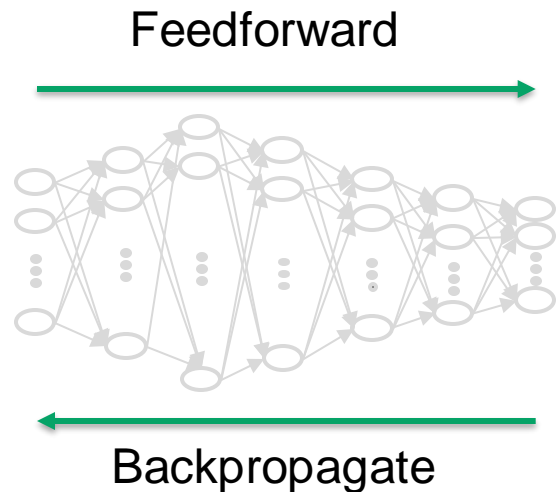  ✓ Higher error, less accurate
Adoptive learning rate
  ✓ Start with large magnitude and make it smaller as slope decreases

# Learning Algorithms

- Gradient descent and backpropagation

  Backpropagate the error and update the weight and biases

Feedforward

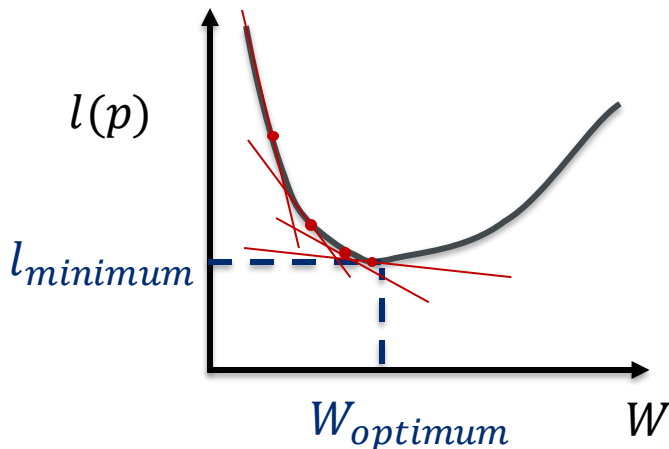Backpropagate

Learning rate
$(0 \leq \gamma \leq 1)$

$$p_{k,updated} = p_{k,old} - \gamma \frac{\partial l}{\partial p_k}$$

A network variable (e.g., the weight corresponding to the weight connection of neuron 5 in 3rd layer and neuron 23 in the 4th layer)

# Learning Algorithms

- Gradient descent and backpropagation

  Backpropagate the error and update the weight and biases



In each iteration:
- ✓ Compute $\frac{\partial l}{\partial p_k}$ using all training repository at once in each iteration
- ✓ Time-consuming and memory intensive for large networks and/or large repositories

# Learning Algorithms

- Stochastic gradient descent and backpropagation

  - Divide training repository to batches of training datapoints, each with a memory-friendly number of datapoints

    - Batch size
    - Batch number

  - Feed each batch to the network one at a time and update the variables

  - Once all batches pass the network once, it is counted as one iteration

  - Train the network for a number of iterations

# Learning Algorithms

- Stochastic gradient descent and backpropagation

  - Example: 32 batches each with 32 training datapoints for $N = 1000$ (notations are based on slide 3 of lecture B of current module):

| Batch 01 (32 datapoints) | Batch 02 (32 datapoints) | | Batch 31 (32 datapoints) | Batch 32 (8 datapoints) |
|---|---|---|---|---|
| $X_1, Y_1$ | $X_{33}, Y_{33}$ | | $X_{961}, Y_{961}$ | $X_{993}, Y_{993}$ |
| $X_2, Y_2$ | $X_{34}, Y_{34}$ | | $X_{962}, Y_{962}$ | $X_{994}, Y_{994}$ |
| . | . | . . . | . | . |
| . | . | | . | . |
| . | . | | . | . |
| $X_{32}, Y_{32}$ | $X_{64}, Y_{64}$ | | $X_{992}, Y_{992}$ | $X_{1000}, Y_{1000}$ |

# Learning Algorithms
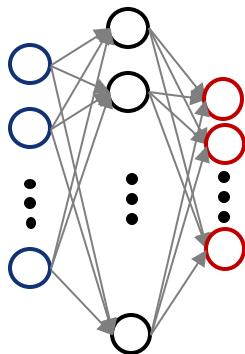
- Stochastic gradient descent versus gradient descent

Gradient Descent  Stochastic Gradient Descent

# Learning Algorithms

- ## Universal Approximation Theorem

  o With mild assumptions on the transfer function, a single hidden layer feed-forward neural network with a finite number of neurons can approximate any convex continues function.

Cybenko, G. (1989) "Approximations by superpositions of sigmoidal functions", Mathematics of Control, Signals, and Systems, 2(4), 303–314. doi:10.1007/BF02551274
Hanin, B. (2018). Approximating Continuous Functions by ReLU Nets of Minimal Width. arXiv preprint arXiv:1710.11278.

# Learning Algorithms

- In this lecture, you learned about:

  - Gradient descent

  - Stochastic gradient descent

  - Backpropagation

  - Universal approximation theorem

- In the next module, we will practice training and testing artificial neural networks over available programming packages