

Intro to Neural Networks

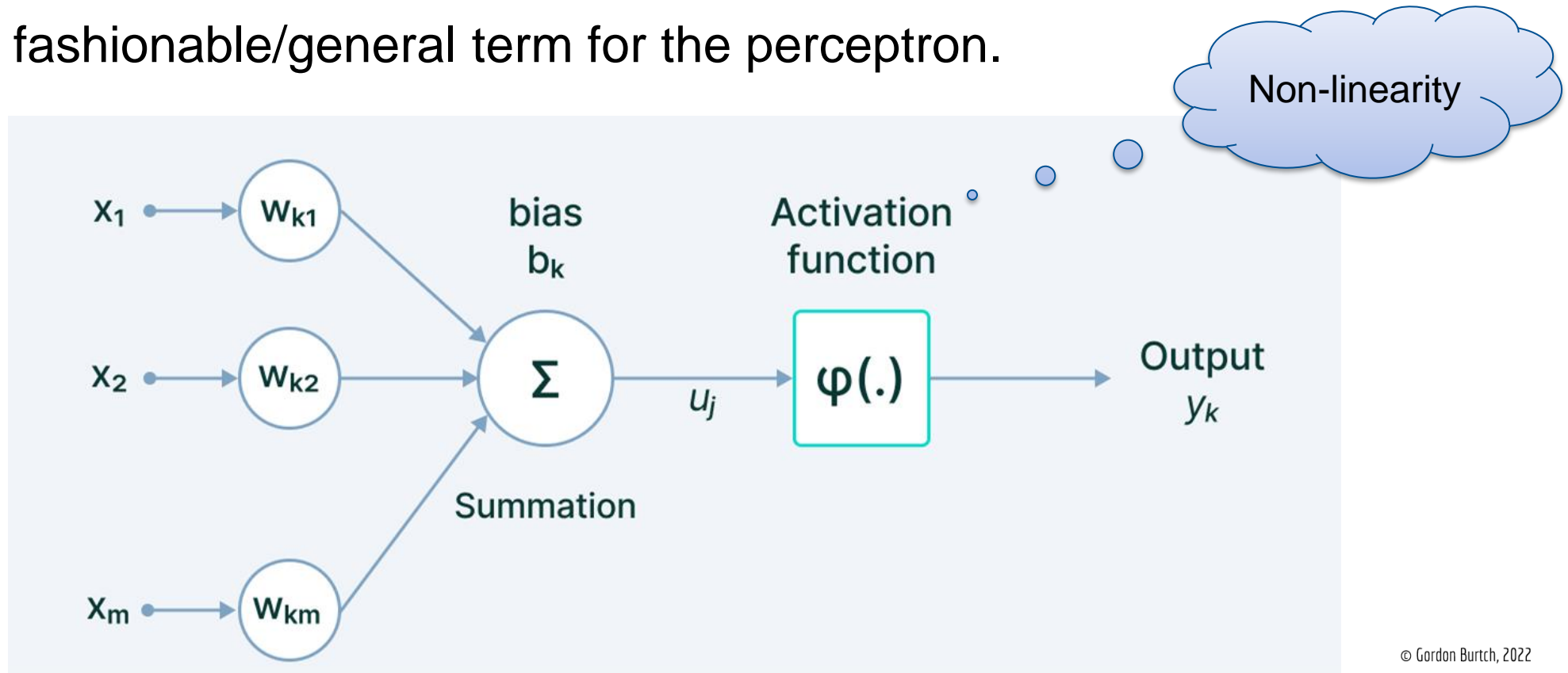
BA865 – Mohannad Elhamod

MLPs

The Multi-Layer Perceptron

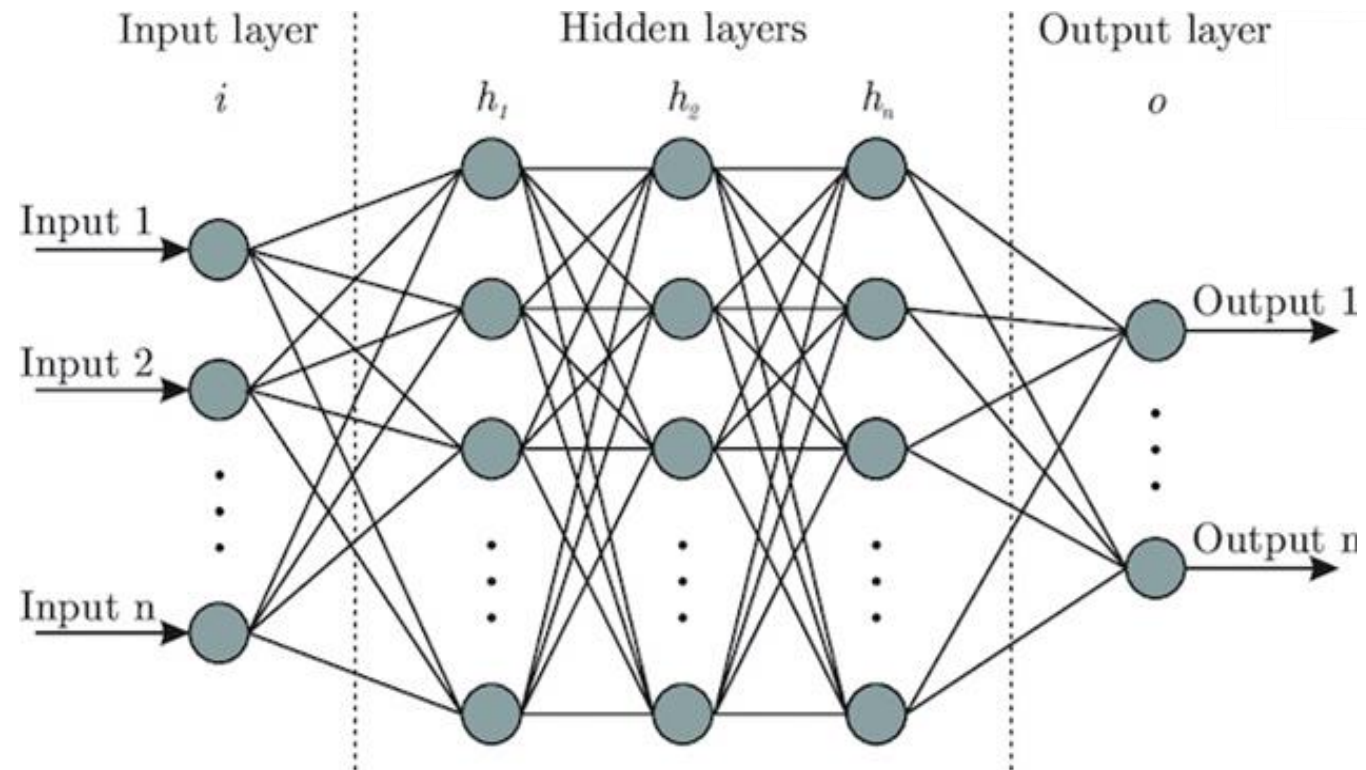
A Neuron

- A more fashionable/general term for the perceptron.



© Gordon Burtch, 2022

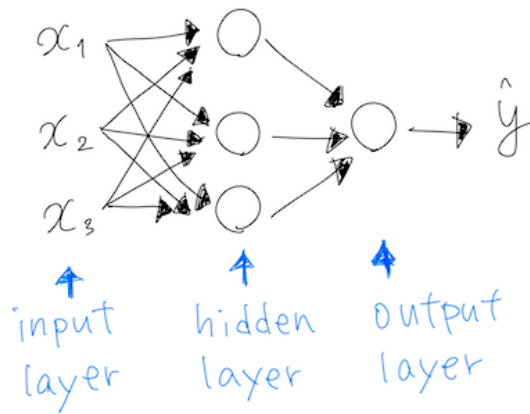
Neural Networks



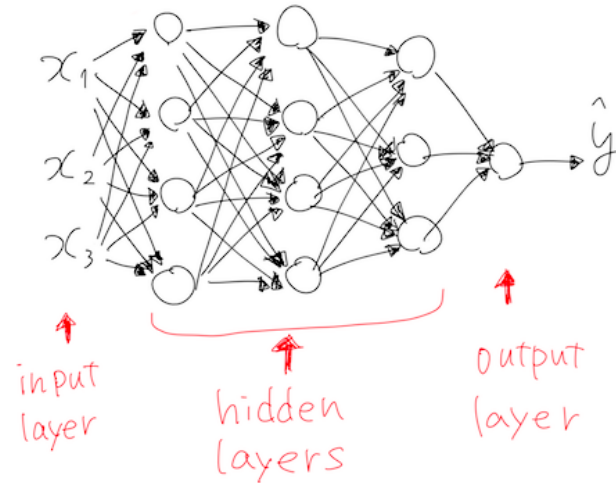
Deep Networks

- Deep = More and more layers...
- leading to more complexity and better capacity for capturing complex phenomena.

Shallow Neural Network

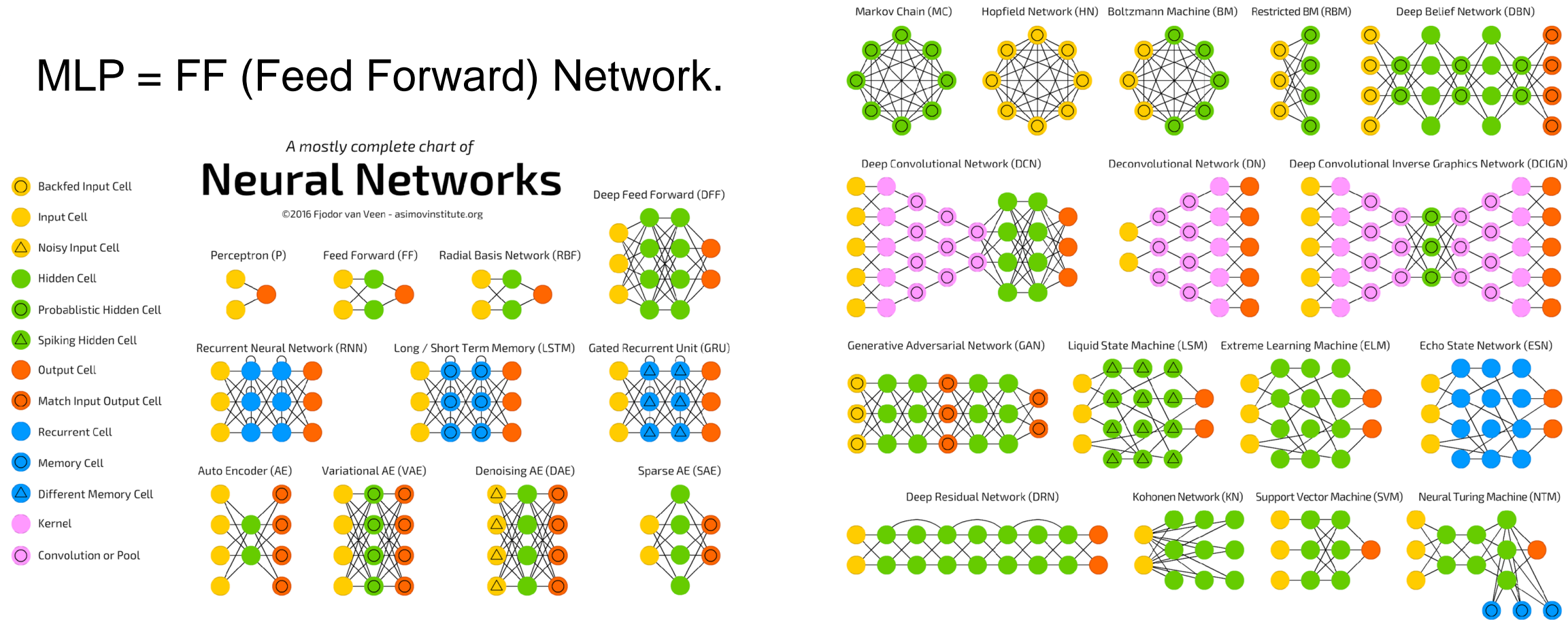


Deep Neural Network



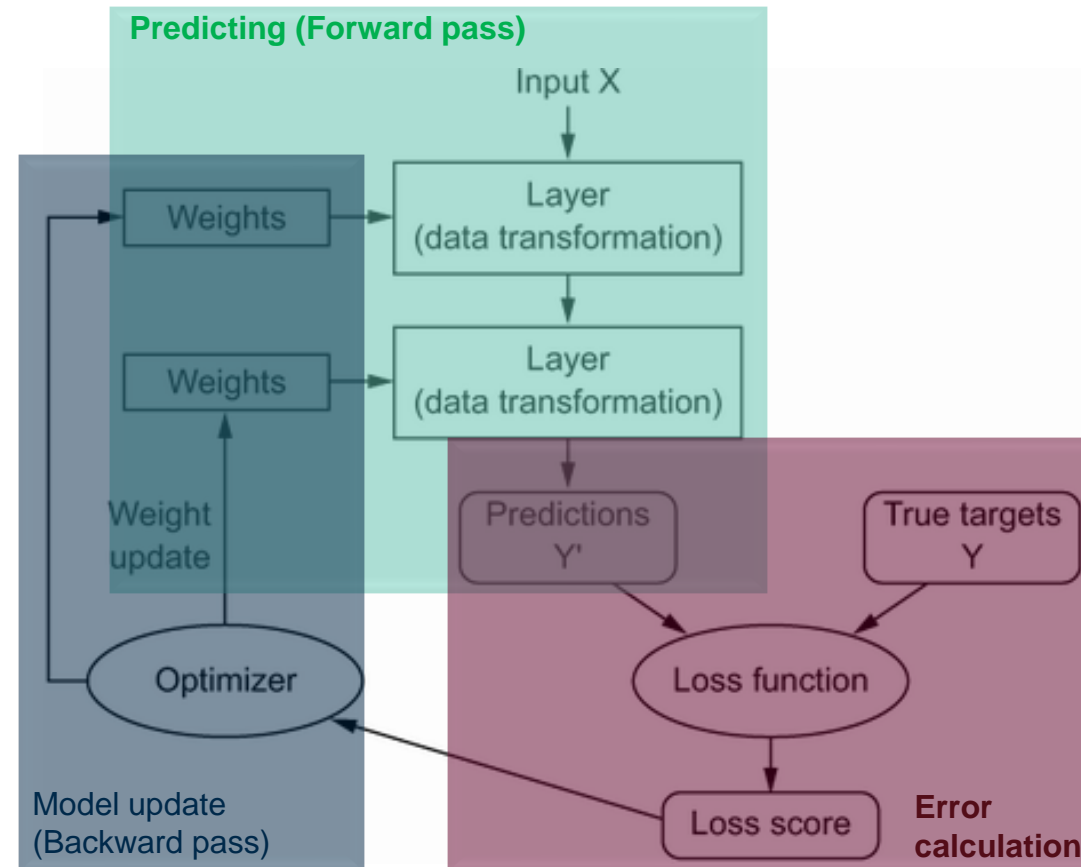
MLPs, One of Many Types...

- MLP = FF (Feed Forward) Network.

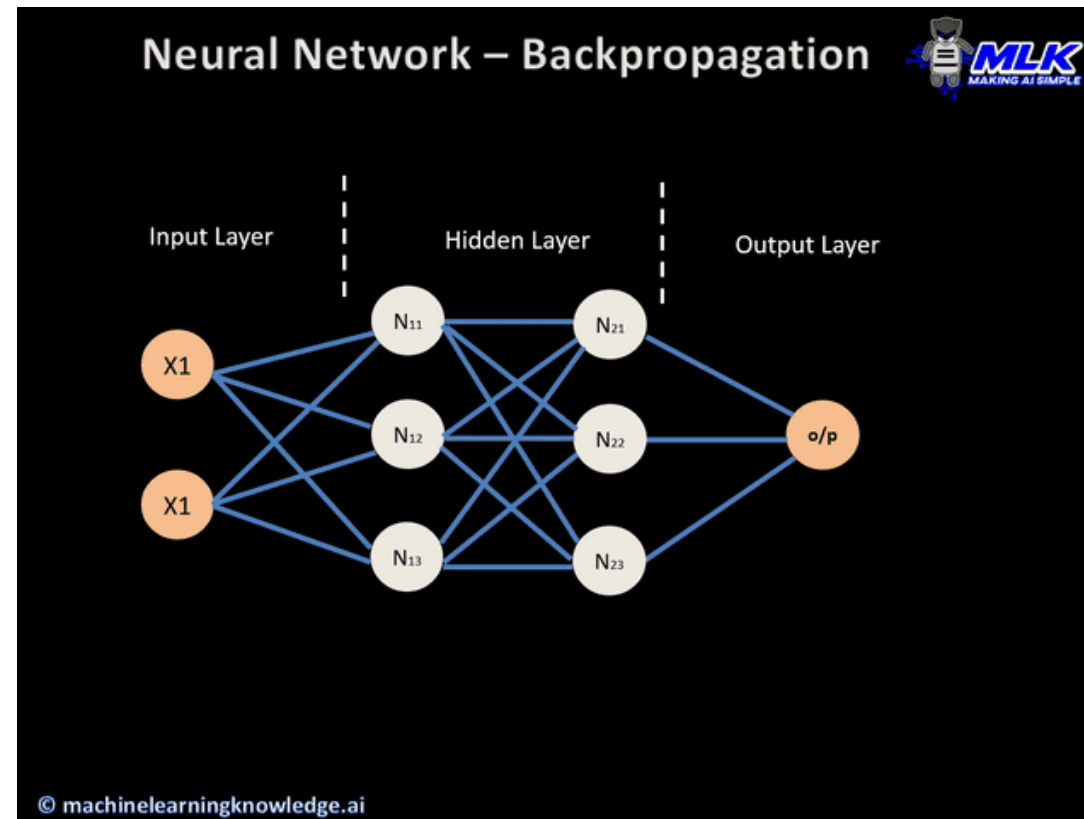


Disecting The Neural Network

The Framework

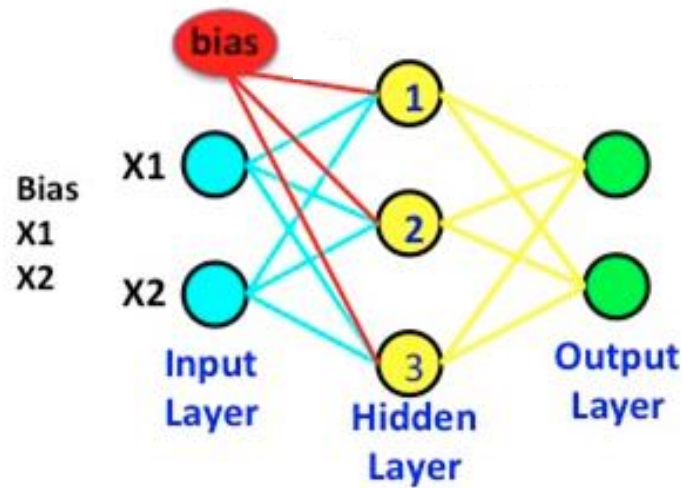


The Framework



Predicting

- What is a layer actually doing?
- Each layer is a matrix multiplication followed by a non-linearity!
 - Why bother with the non-linearity?!



Input Layer

bias	X1	X2
1	0	1
1	0	0
1	0	0
1	1	0

4 x 3

Weights w^T (transposed)

.5	.5	.5
.5	.5	.5
.5	.5	.5

3 x 3

Go to Hidden Nodes

1 2 3

Hidden Layer

1	1	1
.5	.5	.5
.5	.5	.5
1	1	1

4 x 3

Sigmoid Function

Weights

.2	.1
.4	.1
.4	.1

3 x 2

Output Layer

1	.3
.5	.15
.5	.15
1	.3

4 x 2

Formula: $\frac{1}{1 + e^{-(wx+b)}}$

Figure (modified) courtesy of Baber Zaimen

Predicting

- [Demo](#)
- The non-linearities allow the neural net to “warp” a non-linear problem into a linear one!

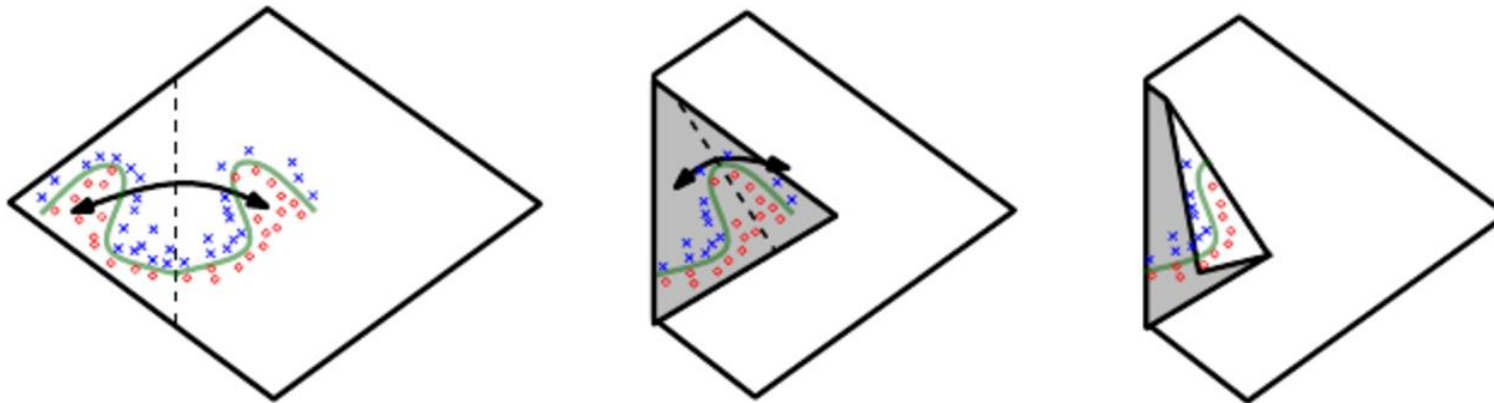


Figure courtesy of [Deep Learning Book](#)

Optimization

- Using Gradient Descent (or some other optimizer) to “update the network.”
 - What do we exactly mean by “updating the network”?

Optimization

- Gradient descent is performed with respect to the weights/biases.

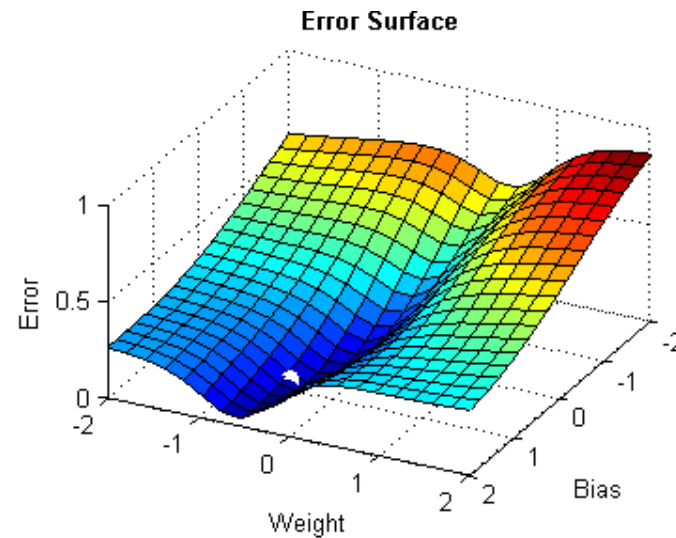
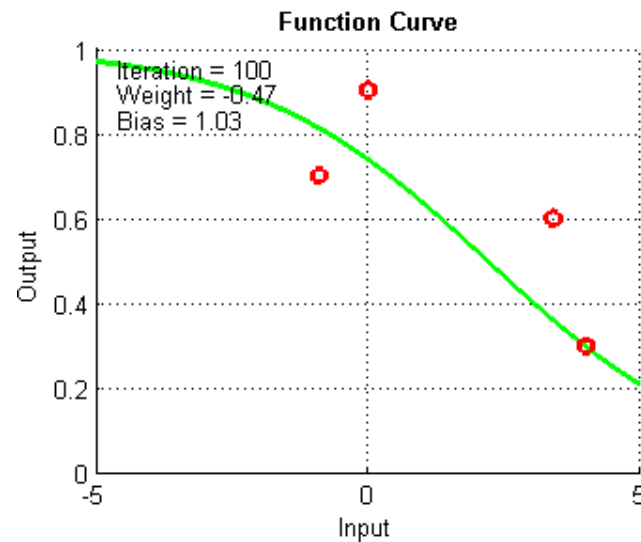
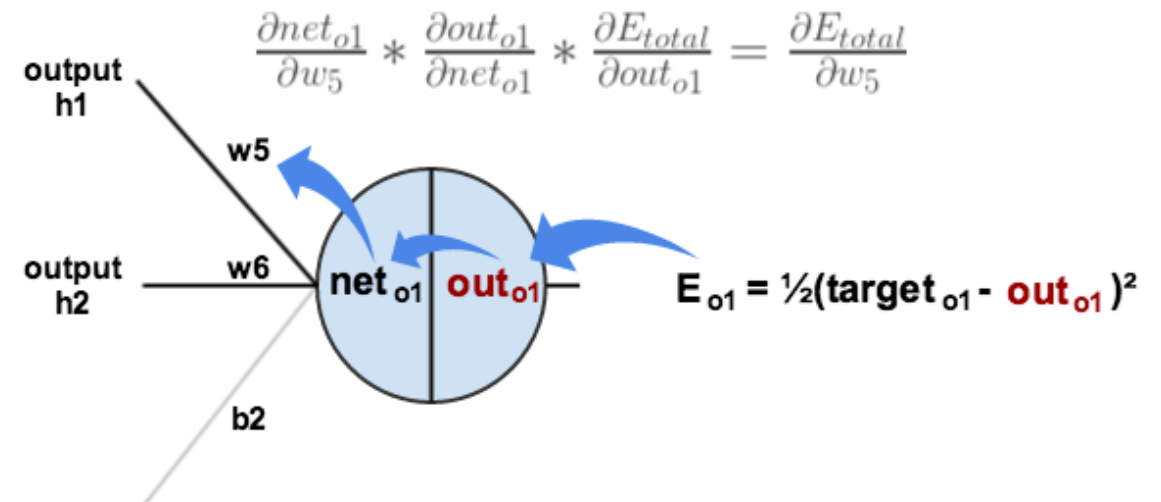


Figure courtesy of [Devin Soni](#)

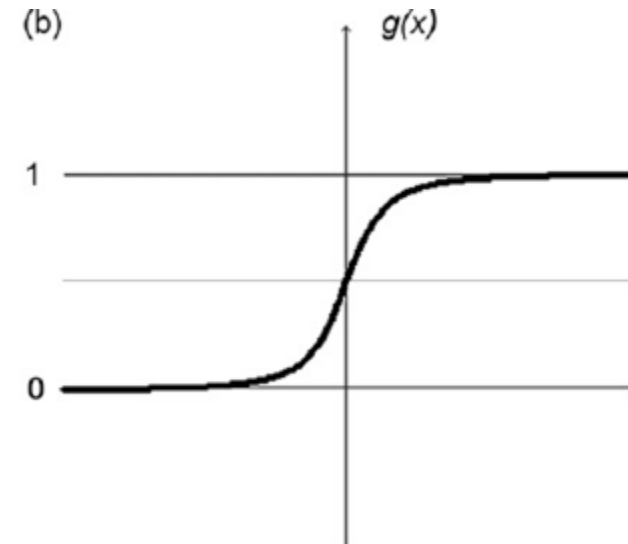
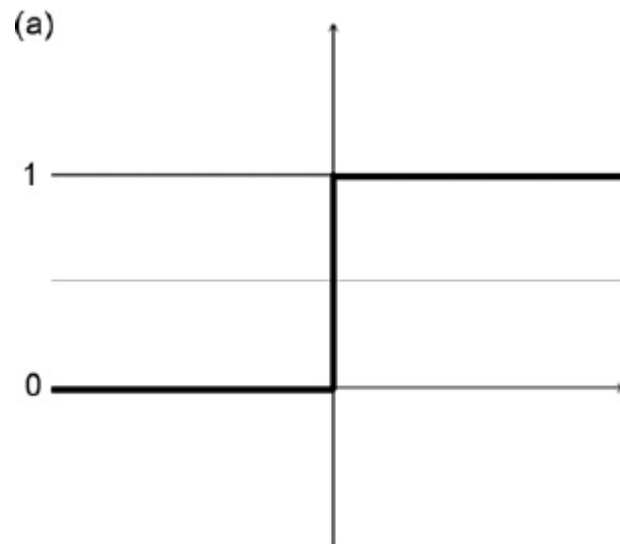
Optimization

- But how does the update get carried all the way back?
 - Chain rule!
 - This is called back-propagation.
- Luckily, these calculations can be automated with automatic differentiation.



Optimization

- We need to make sure the gradient is non-zero...
 - Otherwise, the gradient can't “flow”!
- Replace the step function with a continuous one!



Hyper-Parameters

Learning Rate

- Generally, the most important hyperparameter of them all!
 - Too low: Really slow convergence.
 - Too high: No convergence.

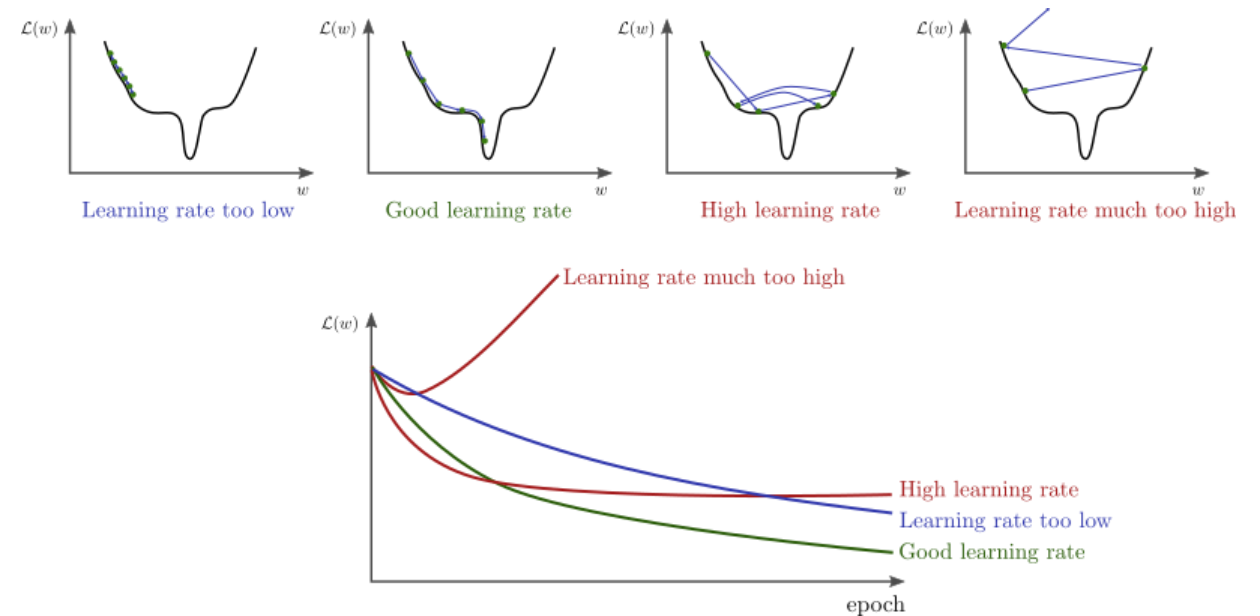


Figure courtesy of Stanford CS class CS231n

Learning Rate: Schedulers

- To get the best of both worlds, you could adjust the learning rate in phases.
 - This way, you still converge but faster.
- Using a scheduler is a common practice.

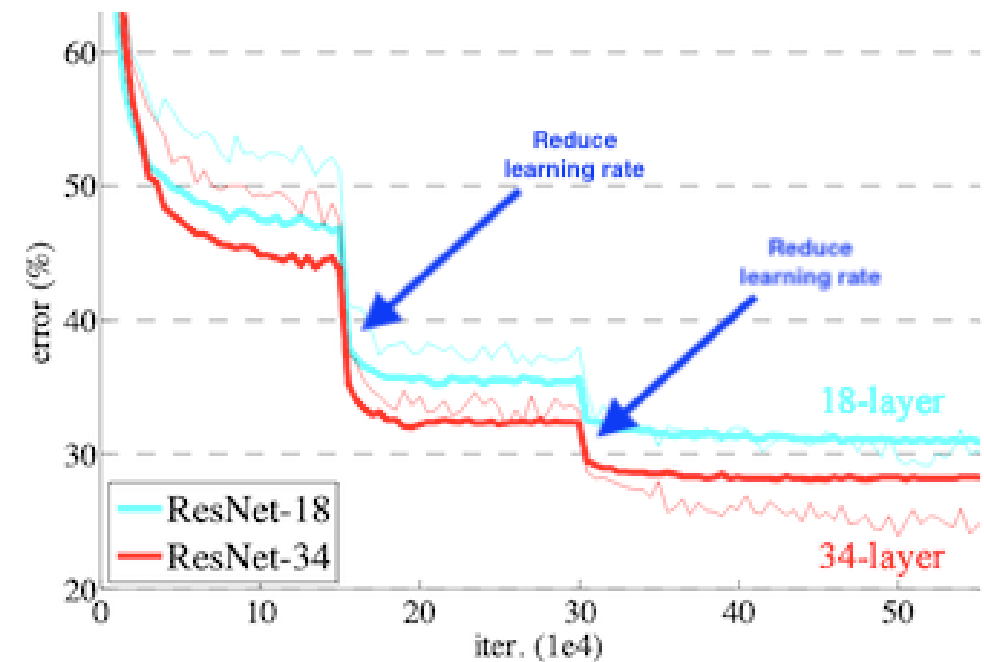


Figure courtesy of [B. D. Hammel](#)

Learning Rate: Early Stopping

- The number of epochs impacts the model's fitness.
- Training needs to stop at the “right” epoch.
- How do we achieve that?
 - Better to stop when validation error stops decreasing for a certain number (n) of epochs.
 - Setting n too small or too large will impact convergence.

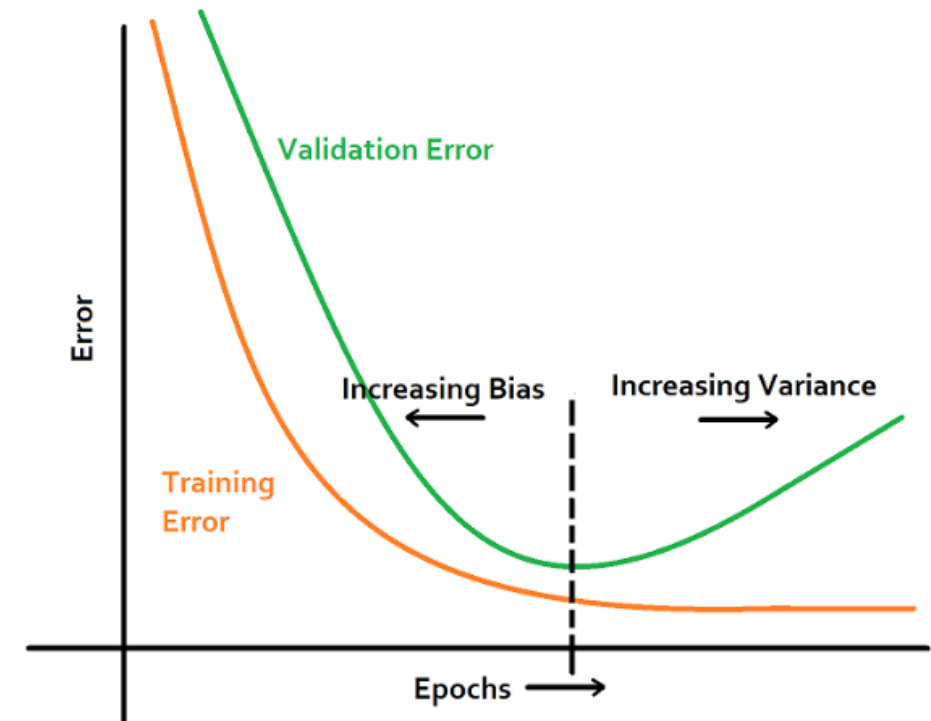


Figure courtesy of RAHUL JAIN

Optimizer: Momentum

- Adding a momentum term (i.e., gradients from previous epochs), helps the convergence process.

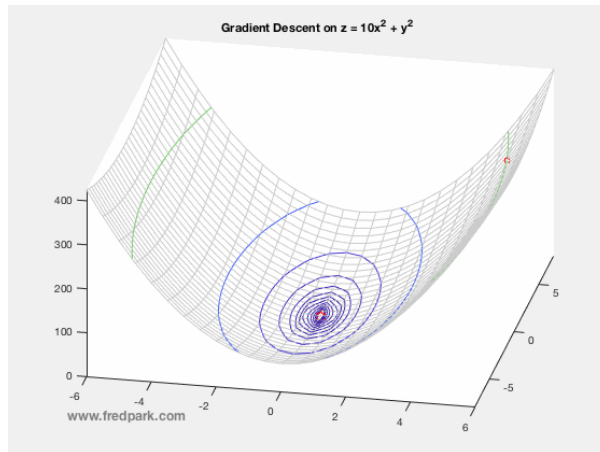


Figure courtesy of [Fred Park](http://www.fredpark.com)

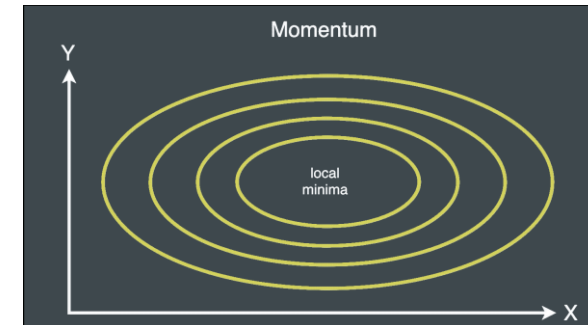
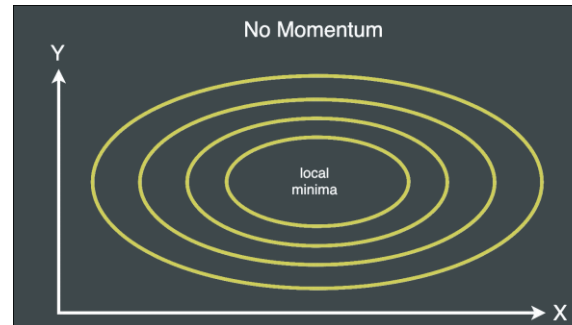


Figure courtesy of [Casper Hansen](#)

$$\begin{aligned} z^{k+1} &= \beta z^k + \nabla f(w^k) \\ w^{k+1} &= w^k - \alpha z^{k+1} \end{aligned}$$

Momentum coefficient