

Introduction to Machine Learning

Foundations and Applications

Paul J. Atzberger
University of California
Santa Barbara

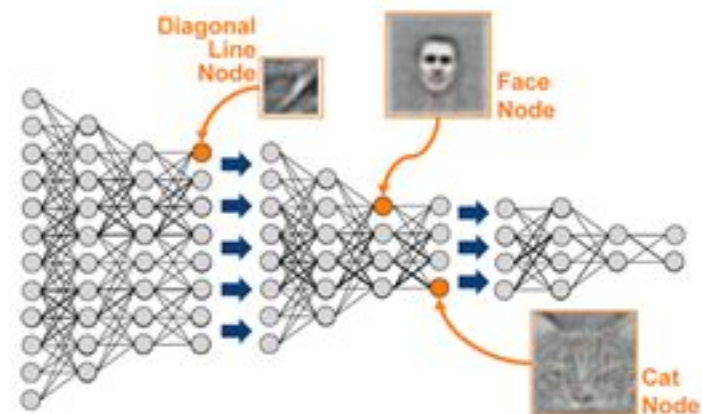




Deep Learning and Neural Networks

Deep Learning Overview

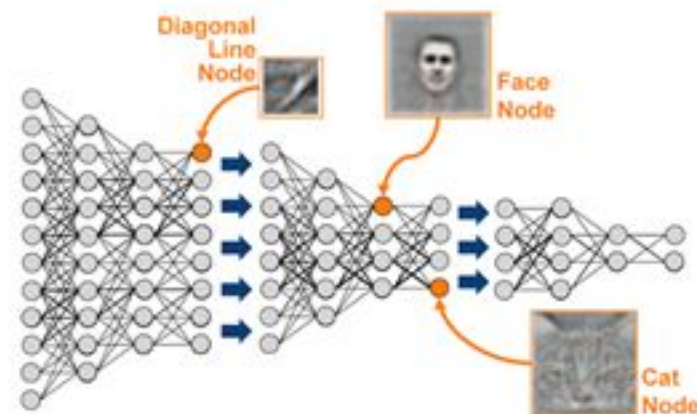
Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.



Deep Learning Overview

Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.

One approach: For regression use linear hypotheses $h(\mathbf{x}; W, b) = \mathbf{x}^T W + b$ or for classification $h(\mathbf{x}; W, b) = \text{sign}(\mathbf{x}^T W + b)$, [convex optimization problems].

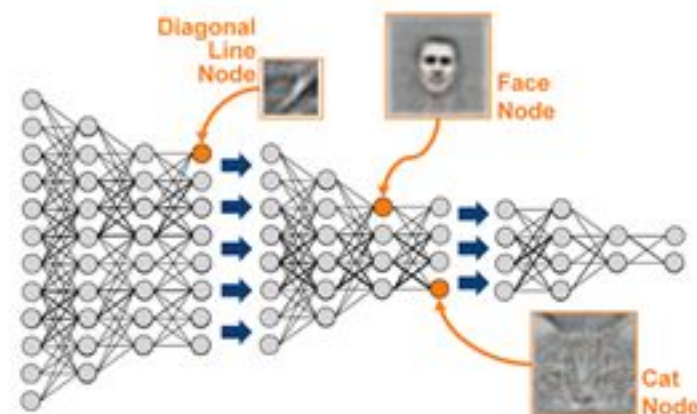


Deep Learning Overview

Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.

One approach: For regression use linear hypotheses $h(\mathbf{x}; W, b) = \mathbf{x}^T W + b$ or for classification $h(\mathbf{x}; W, b) = \text{sign}(\mathbf{x}^T W + b)$, [convex optimization problems].

Non-linear case: We can use **kernel trick** with feature map $\mathbf{z} = \phi(\mathbf{x})$ and use linear methods in feature space $h(\mathbf{x}; W, b) = \phi(\mathbf{x})^T W + b$.



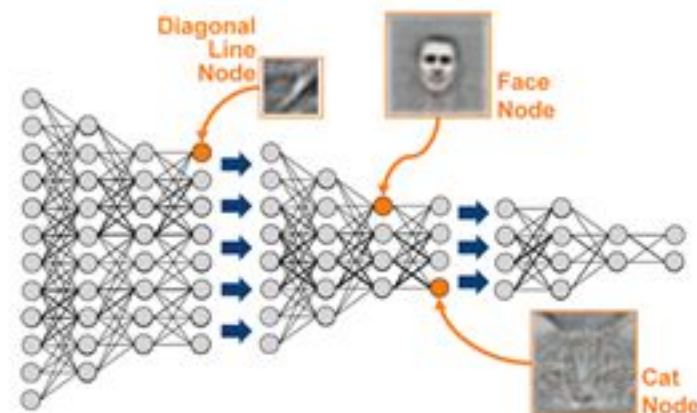
Deep Learning Overview

Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.

One approach: For regression use linear hypotheses $h(\mathbf{x}; W, b) = \mathbf{x}^T W + b$ or for classification $h(\mathbf{x}; W, b) = \text{sign}(\mathbf{x}^T W + b)$, [convex optimization problems].

Non-linear case: We can use **kernel trick** with feature map $\mathbf{z} = \phi(\mathbf{x})$ and use linear methods in feature space $h(\mathbf{x}; W, b) = \phi(\mathbf{x})^T W + b$.

Design of kernel $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ incorporates prior knowledge and traditionally was done manually by experts (computer vision, speech recognition, natural language processing).



Deep Learning Overview

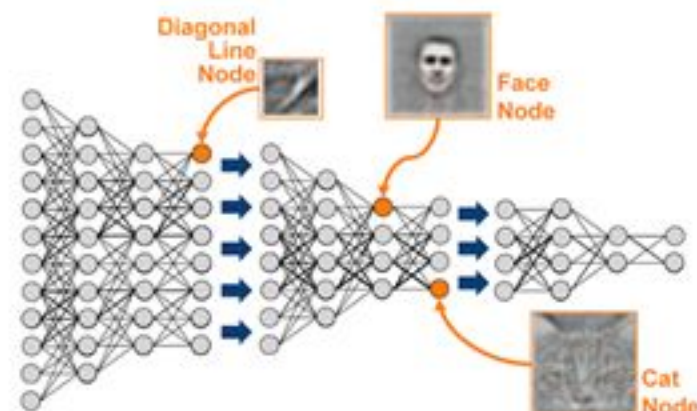
Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.

One approach: For regression use linear hypotheses $h(\mathbf{x}; W, b) = \mathbf{x}^T W + b$ or for classification $h(\mathbf{x}; W, b) = \text{sign}(\mathbf{x}^T W + b)$, [convex optimization problems].

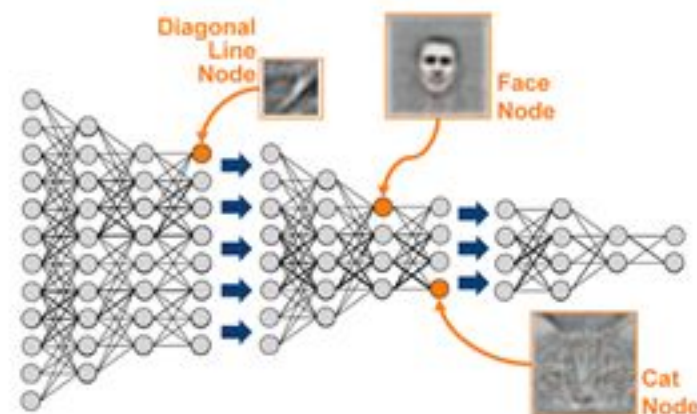
Non-linear case: We can use **kernel trick** with feature map $\mathbf{z} = \phi(\mathbf{x})$ and use linear methods in feature space $h(\mathbf{x}; W, b) = \phi(\mathbf{x})^T W + b$.

Design of kernel $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ incorporates prior knowledge and traditionally was done manually by experts (computer vision, speech recognition, natural language processing).

Data-driven approach: **Try to learn the feature map ϕ from the data!** Find best parameters θ over some class of feature maps $\phi(\mathbf{x}; \theta)$. Choice of feature map class can incorporate prior knowledge.



Deep Learning Overview



Machine Learning: Typical task is to try to learn a function $h(\mathbf{x}; \theta)$ from data $S = \{(x_i, y_i)\}$ that approximates $y = f(x)$.

One approach: For regression use linear hypotheses $h(\mathbf{x}; W, b) = \mathbf{x}^T W + b$ or for classification $h(\mathbf{x}; W, b) = \text{sign}(\mathbf{x}^T W + b)$, [convex optimization problems].

Non-linear case: We can use **kernel trick** with feature map $\mathbf{z} = \phi(\mathbf{x})$ and use linear methods in feature space $h(\mathbf{x}; W, b) = \phi(\mathbf{x})^T W + b$.

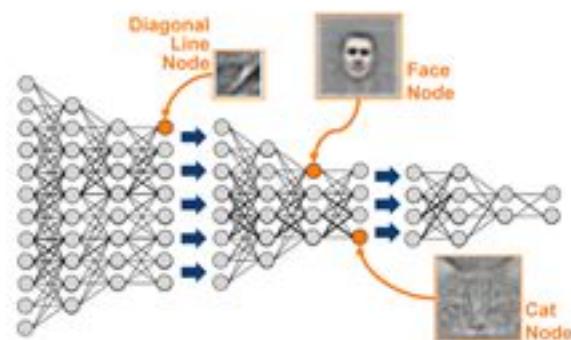
Design of kernel $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ incorporates prior knowledge and traditionally was done manually by experts (computer vision, speech recognition, natural language processing).

Data-driven approach: **Try to learn the feature map ϕ from the data!** Find best parameters θ over some class of feature maps $\phi(\mathbf{x}; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions we try to learn in applications often can be approximated well by a composition of functions: $y = f(\mathbf{x}) = f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$, where $f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}$.

Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

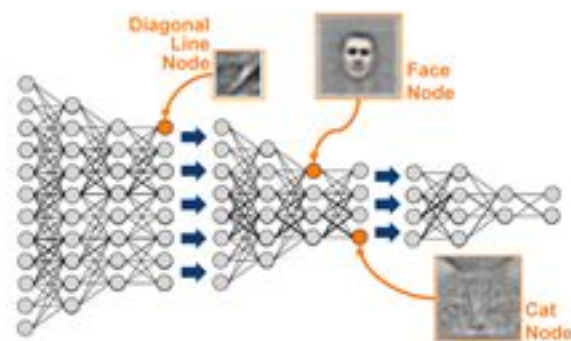


Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions in applications can be approximated well by a composition of functions:

$$y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right), \text{ where } f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}.$$



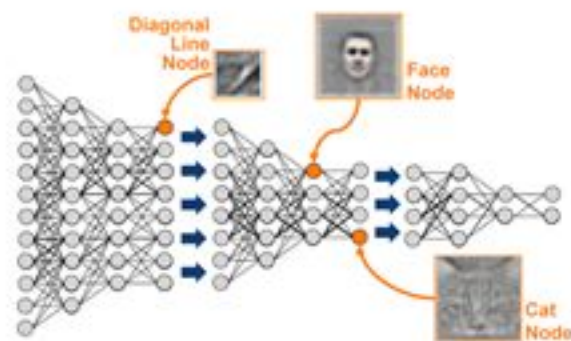
Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions in applications can be approximated well by a composition of functions:

$$y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right), \text{ where } f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}.$$

Ex: Regression $y = |\cos(x_1 x_2) + 1|$ would have $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ with $f^{(1)}(x_1, x_2) = x_1 x_2$ and $f^{(2)}(h^{(1)}) = \cos(h^{(1)}) + 1$ and $f^{(3)}(h^{(2)}) = \text{abs}(h^{(2)})$.



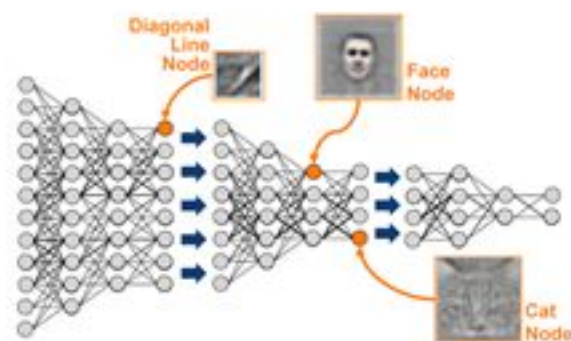
Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions in applications can be approximated well by a composition of functions:

$$y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right), \text{ where } f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}.$$

Ex: Regression $y = |\cos(x_1 x_2) + 1|$ would have $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ with $f^{(1)}(x_1, x_2) = x_1 x_2$ and $f^{(2)}(h^{(1)}) = \cos(h^{(1)}) + 1$ and $f^{(3)}(h^{(2)}) = \text{abs}(h^{(2)})$.



Generally, we can think about $y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right)$ in terms of components as

$$y = f(x_1, x_2, \dots, x_N) = f^{(L)} \left(f_1^{(L-1)}(z_{i_1}, z_{i_2}, \dots, z_{i_{N_1}}), f_2^{(L-1)}(z_{j_1}, z_{j_2}, \dots, z_{j_{N_2}}), \dots, f_{n_2}^{(L-1)}(z_{k_1}, z_{k_2}, \dots, z_{k_{N_{n_2}}}) \right), z = f^{(L-2)}$$

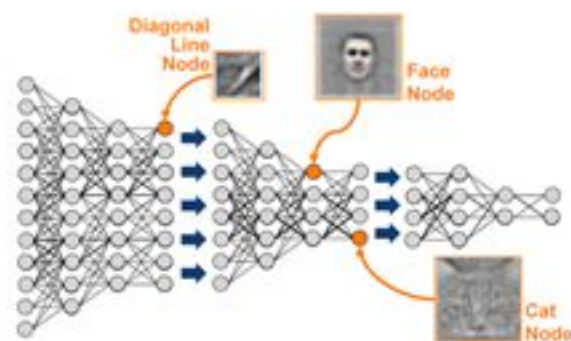
Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions in applications can be approximated well by a composition of functions:

$$y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right), \text{ where } f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}.$$

Ex: Regression $y = |\cos(x_1 x_2) + 1|$ would have $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ with $f^{(1)}(x_1, x_2) = x_1 x_2$ and $f^{(2)}(h^{(1)}) = \cos(h^{(1)}) + 1$ and $f^{(3)}(h^{(2)}) = \text{abs}(h^{(2)})$.



Generally, we can think about $y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right)$ in terms of components as

$$y = f(x_1, x_2, \dots, x_N) = f^{(L)} \left(f_1^{(L-1)}(z_{i_1}, z_{i_2}, \dots, z_{i_{N_1}}), f_2^{(L-1)}(z_{j_1}, z_{j_2}, \dots, z_{j_{N_2}}), \dots, f_{n_2}^{(L-1)}(z_{k_1}, z_{k_2}, \dots, z_{k_{N_{n_2}}}) \right), z = f^{(L-2)}$$

Concepts represented by hierarchy of distributed features. For instance: position of robotic actuator from angles in arm, identity of person from parts of the face, meaning of a sentence from phrases/words.

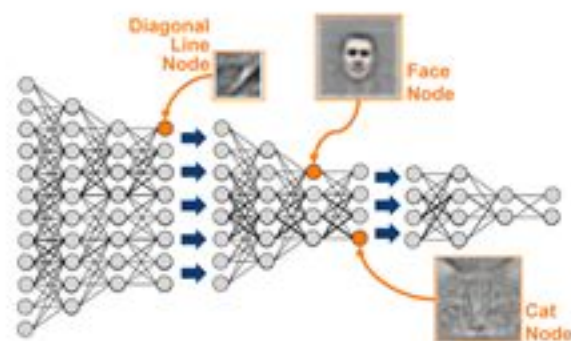
Deep Learning Overview

Data-driven approaches: Try to learn the feature map ϕ from the data! Find best parameters θ over some class of feature maps $\phi(x; \theta)$. Choice of feature map class can incorporate prior knowledge.

Deep Learning: Many functions in applications can be approximated well by a composition of functions:

$$y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right), \text{ where } f^{(k)}: \mathbb{R}^{N_k} \rightarrow \mathbb{R}^{N_{k+1}}.$$

Ex: Regression $y = |\cos(x_1 x_2) + 1|$ would have $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ with $f^{(1)}(x_1, x_2) = x_1 x_2$ and $f^{(2)}(h^{(1)}) = \cos(h^{(1)}) + 1$ and $f^{(3)}(h^{(2)}) = \text{abs}(h^{(2)})$.



Generally, we can think about $y = f(x) = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)}(x) \dots \right) \right)$ in terms of components as

$$y = f(x_1, x_2, \dots, x_N) = f^{(L)} \left(f_1^{(L-1)}(z_{i_1}, z_{i_2}, \dots, z_{i_{N_1}}), f_2^{(L-1)}(z_{j_1}, z_{j_2}, \dots, z_{j_{N_2}}), \dots, f_{n_2}^{(L-1)}(z_{k_1}, z_{k_2}, \dots, z_{k_{N_{n_2}}}) \right), z = f^{(L-2)}$$

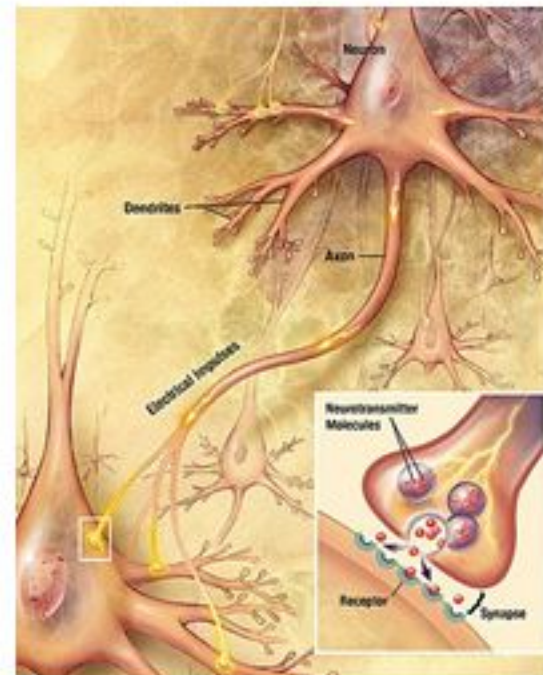
Concepts represented by hierarchy of distributed features. For instance: position of robotic actuator from angles in arm, identity of person from parts of the face, meaning of a sentence from phrases/words.

View as building up function from "hidden units" that detect particular features of input \mathbf{z} as $h^{(k)} = f^{(k)}(\mathbf{z})$. Popular way to do this is to use Artificial Neural Networks (ANNs), $h^{(k)} = g(\mathbf{z}^T W + b)$.

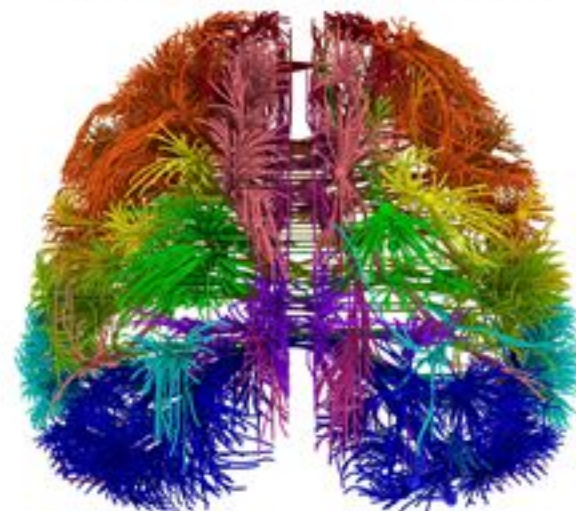
What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors:
action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Neuron: Axons and Dendrites



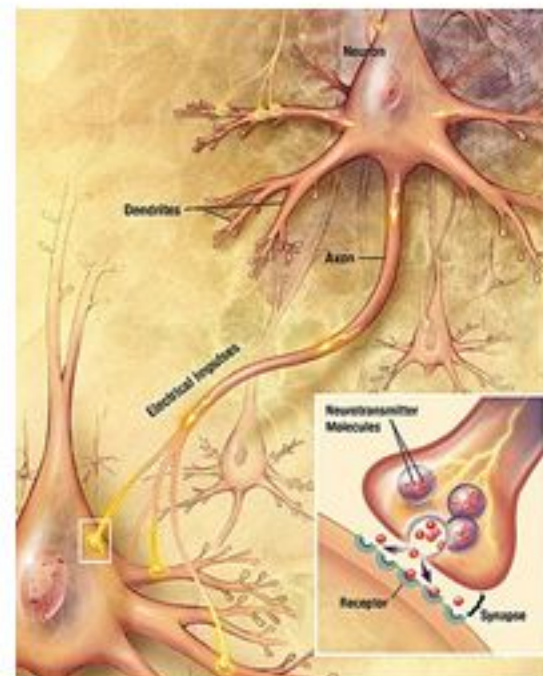
Human Connectome Project



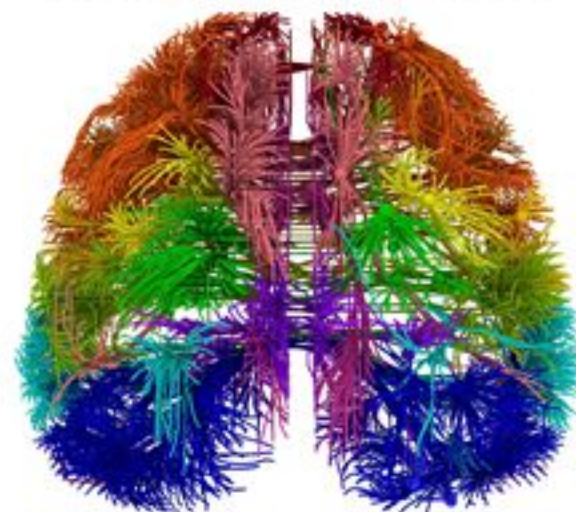
What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Neuron: Axons and Dendrites



Human Connectome Project

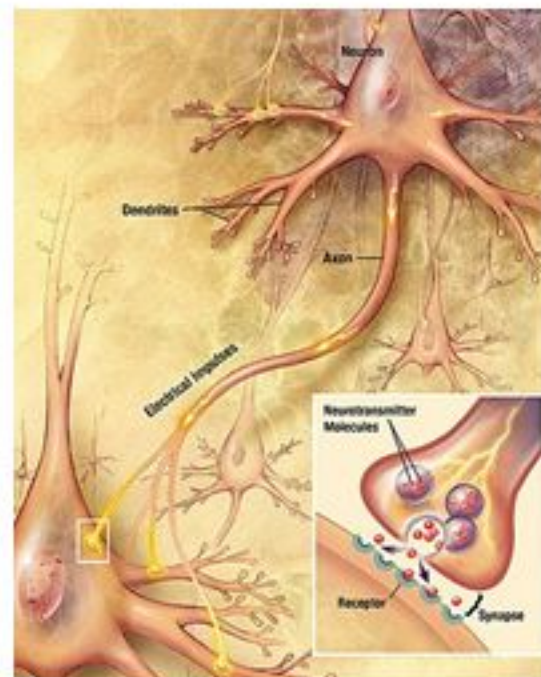


NN's are engineering tools inspired by nature. **Not a model of real neurons!**

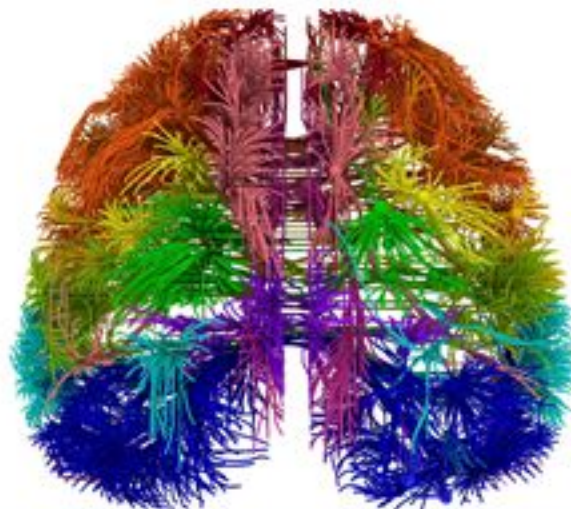
What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Neuron: Axons and Dendrites



Human Connectome Project



NN's are engineering tools inspired by nature. **Not a model of real neurons!** Biology is more complex / temporal dynamics, refraction / other relevant factors.

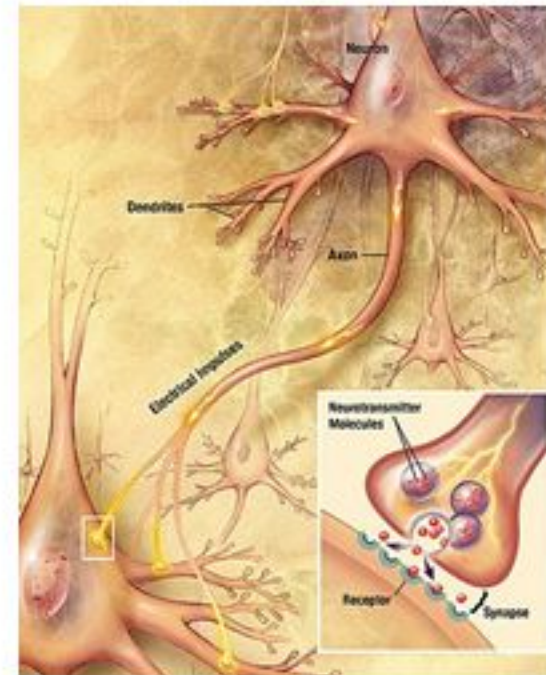
What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

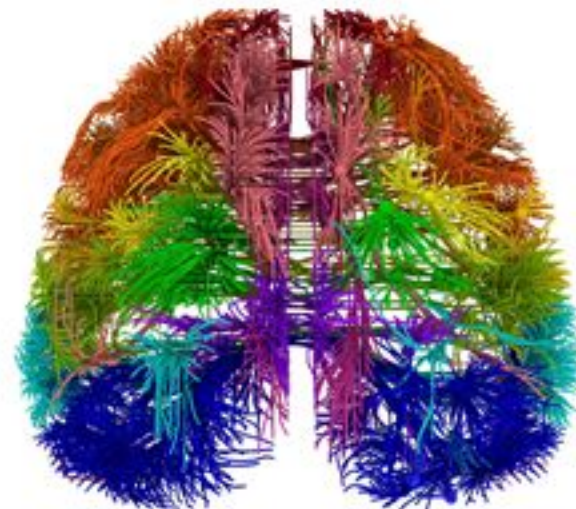
Winged Flight



Neuron: Axons and Dendrites



Human Connectome Project



NN's are engineering tools inspired by nature. **Not a model of real neurons!**
Biology is more complex / temporal dynamics, refraction / other relevant factors.

What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials \rightarrow voltage-gated ion channels / neurotransmitters \rightarrow collective neural activity.

Winged Flight

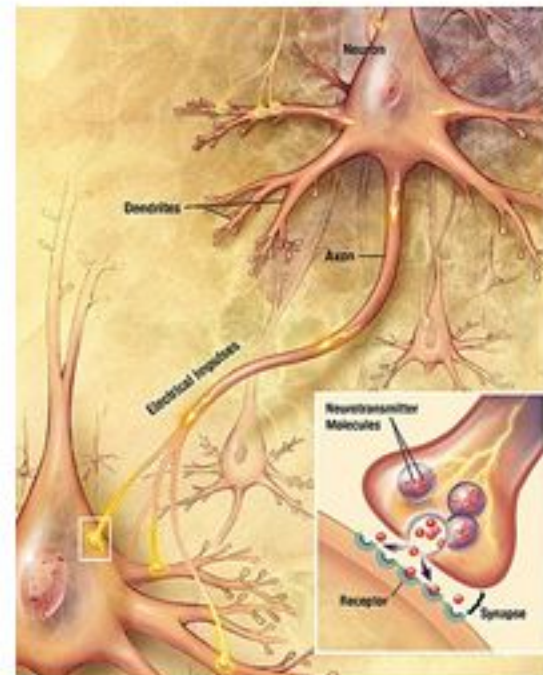


inspiration \rightarrow

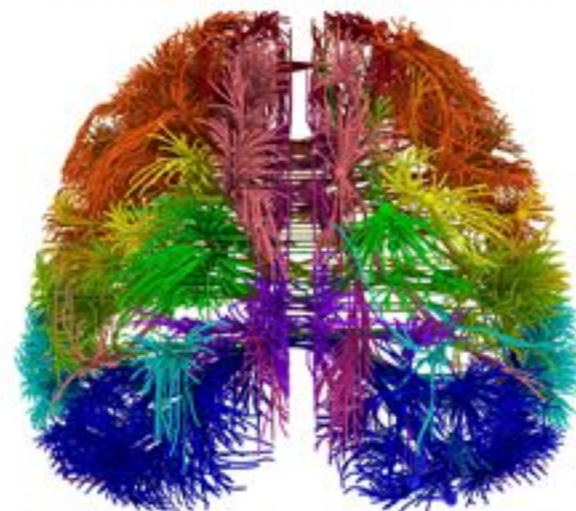
Wright Brothers



Neuron: Axons and Dendrites



Human Connectome Project



NN's are engineering tools inspired by nature. **Not a model of real neurons!**
Biology is more complex / temporal dynamics, refraction / other relevant factors.

What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Winged Flight



inspiration →

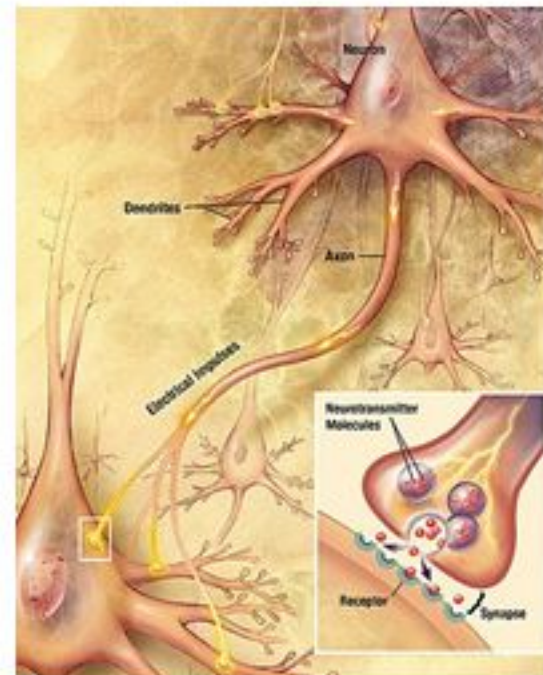
Wright Brothers



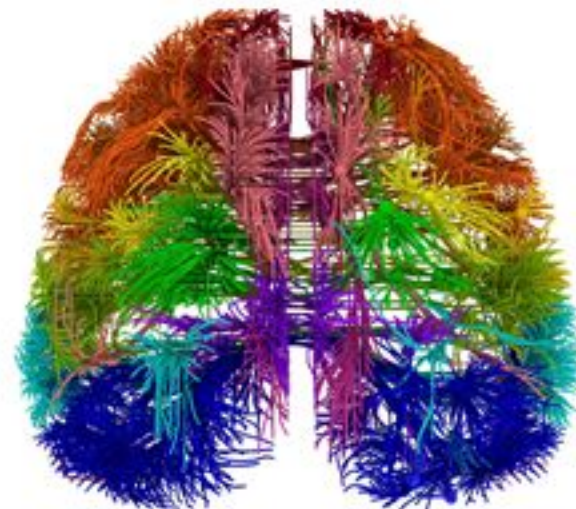
Modern Airplanes



Neuron: Axons and Dendrites



Human Connectome Project



NN's are engineering tools inspired by nature. **Not a model of real neurons!**
Biology is more complex / temporal dynamics, refraction / other relevant factors.

What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Winged Flight

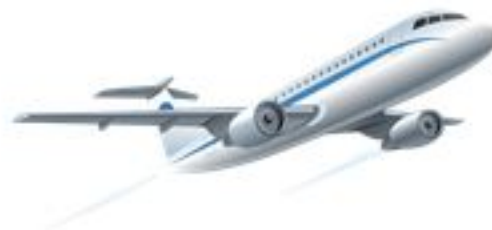


inspiration

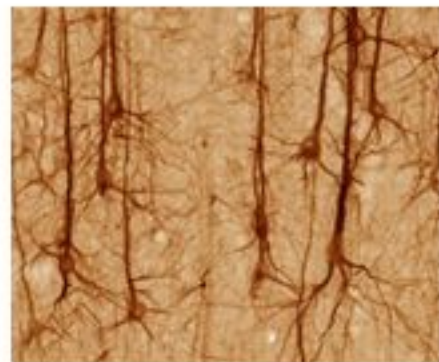
Wright Brothers



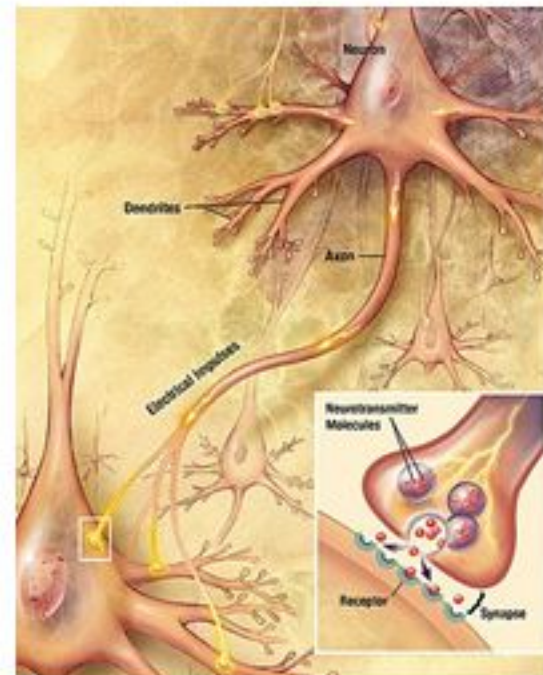
Modern Airplanes



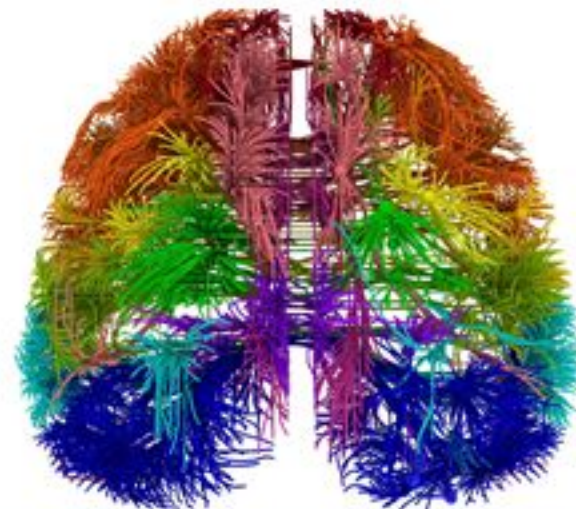
Neural Information Processing



Neuron: Axons and Dendrites



Human Connectome Project



NN's are engineering tools inspired by nature. **Not a model of real neurons!**
Biology is more complex / temporal dynamics, refraction / other relevant factors.

What are Artificial Neural Networks?

Biological neurons process information by firing to excite or inhibit neighbors: action potentials → voltage-gated ion channels / neurotransmitters → collective neural activity.

Winged Flight



inspiration →

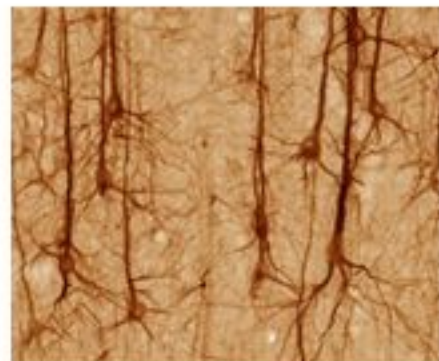
Wright Brothers



Modern Airplanes

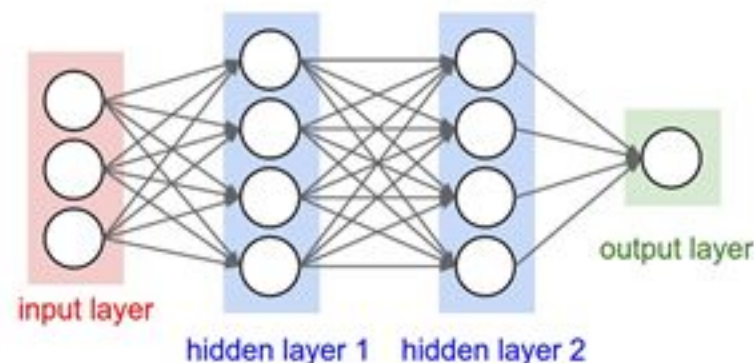


Neural Information Processing

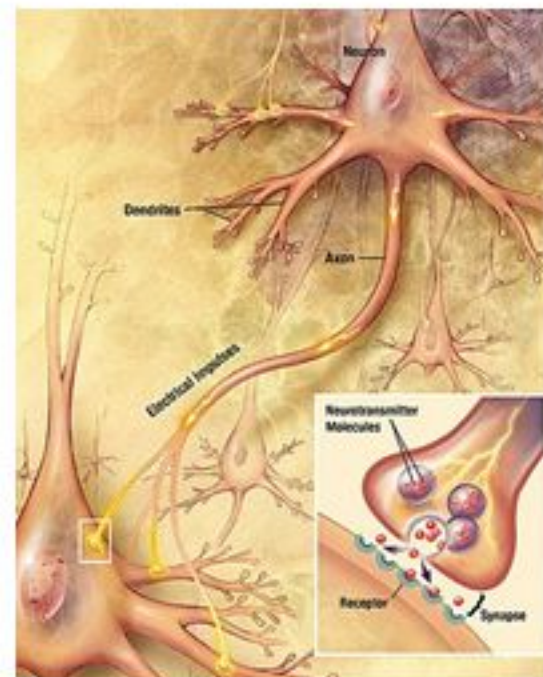


inspiration →

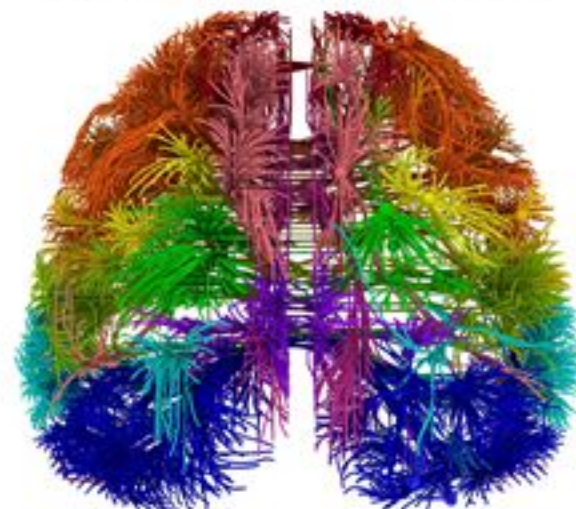
Artificial Neural Network (ANN)



Neuron: Axons and Dendrites

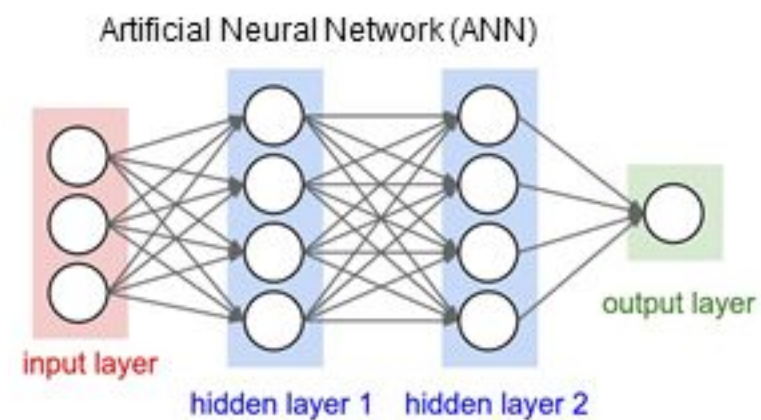


Human Connectome Project



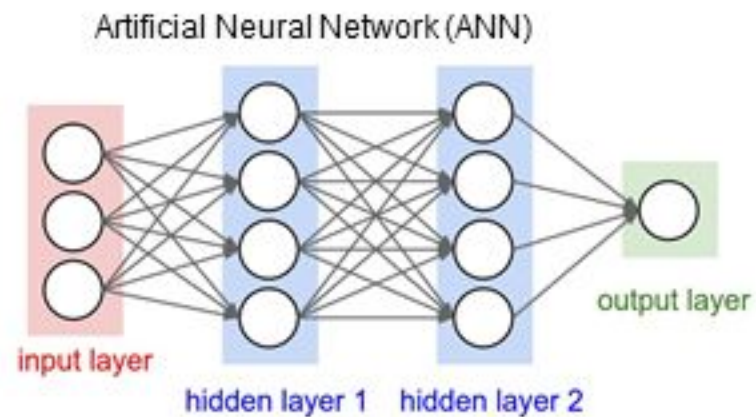
NN's are engineering tools inspired by nature. **Not a model of real neurons!**
Biology is more complex / temporal dynamics, refraction / other relevant factors.

What are Neural Networks (NNs)?



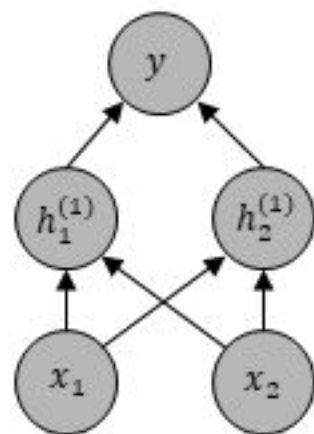
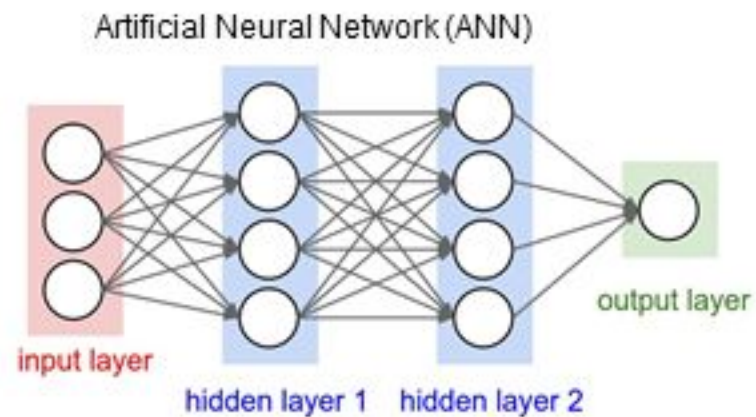
What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}\mathbf{W}^k + b^k)$.



What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}\mathbf{W}^k + b^k)$.

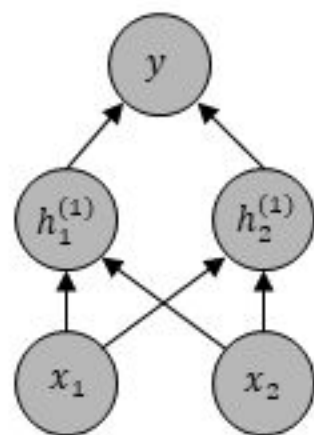
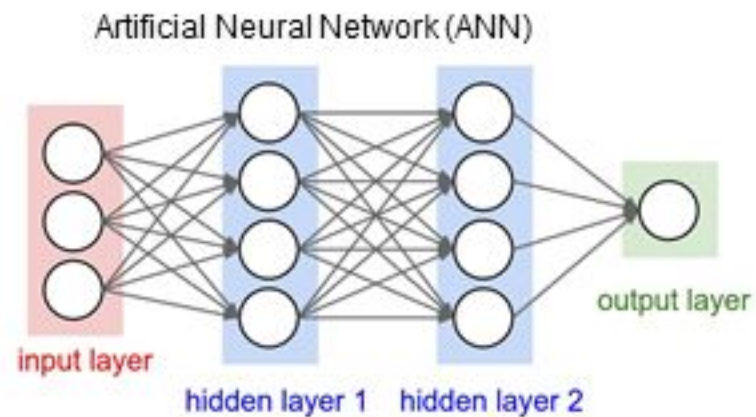


What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}\mathbf{W}^k + b^k)$.

The $g(z)$ is called the **activation function**. Common choices include

- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

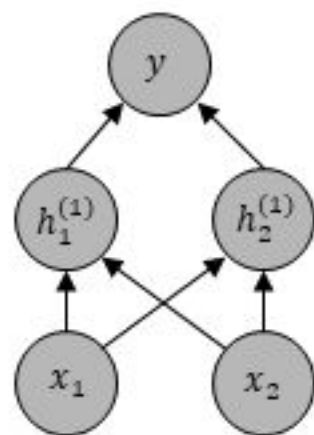
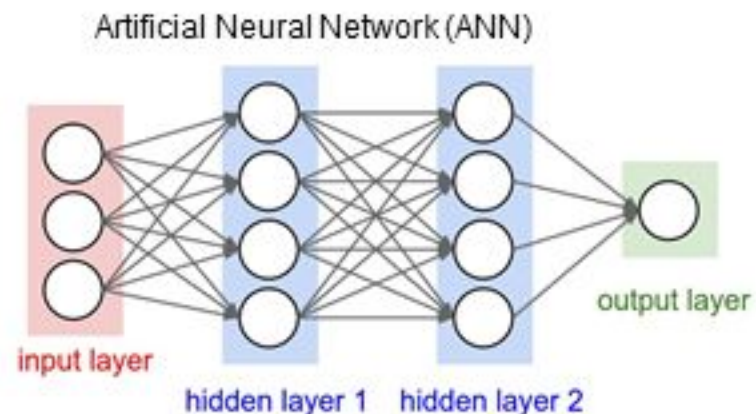


What are Neural Networks (NNs)?

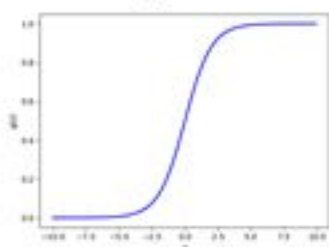
NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}\mathbf{W}^k + b^k)$.

The $g(z)$ is called the **activation function**. Common choices include

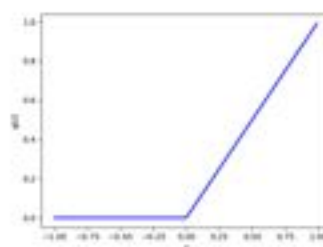
- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.



Sigmoid



ReLU



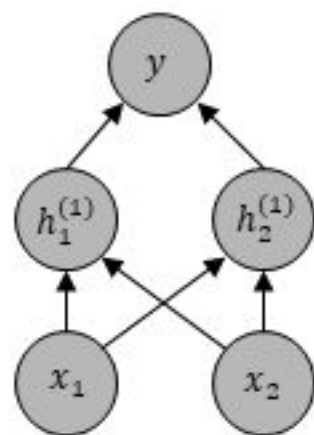
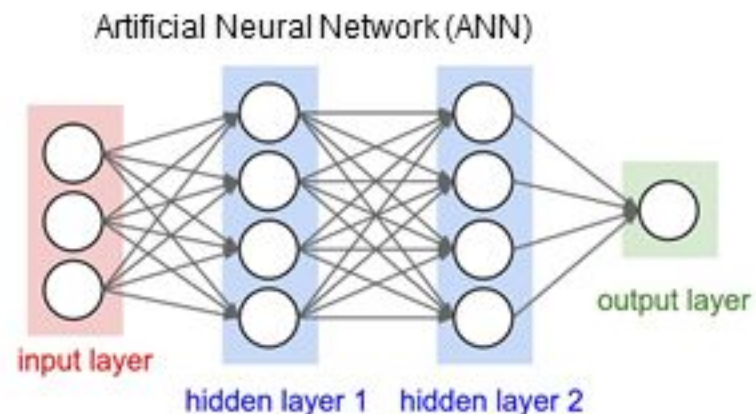
What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

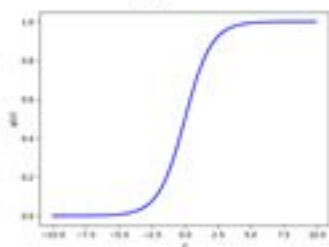
The $g(z)$ is called the **activation function**. Common choices include

- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

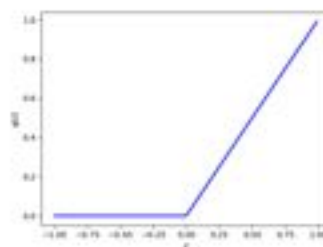
Find best **weights** W and **bias** b for each layer to minimize some **loss** $\ell(\dots)$.



Sigmoid



ReLU



What are Neural Networks (NNs)?

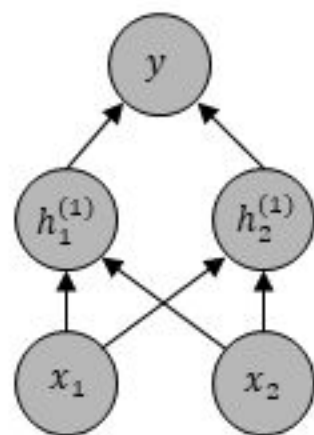
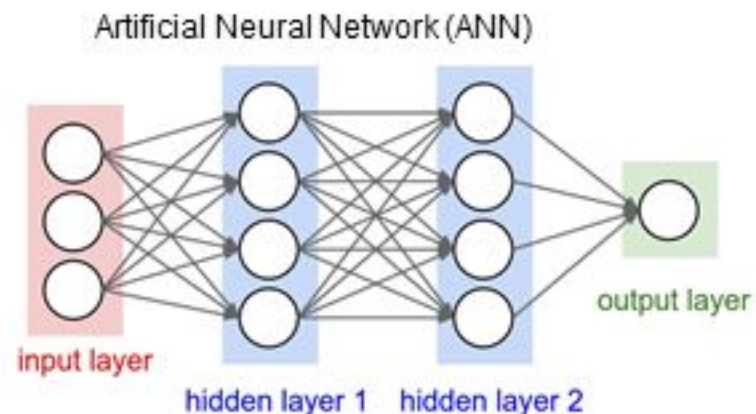
NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

The $g(z)$ is called the **activation function**. Common choices include

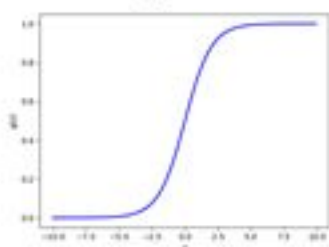
- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

Find best **weights** W and **bias** b for each layer to minimize some **loss** $\ell(\dots)$.

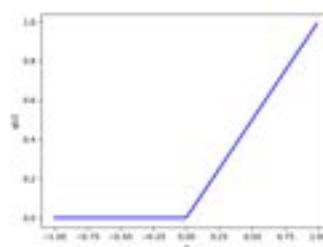
Optimization non-convex so use local gradient learning methods.



Sigmoid



ReLU



What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

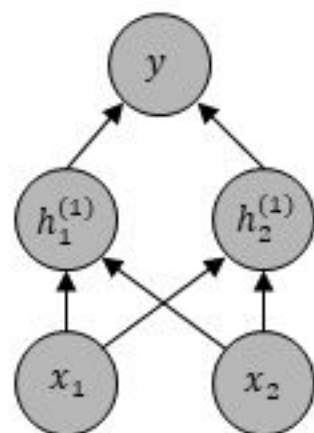
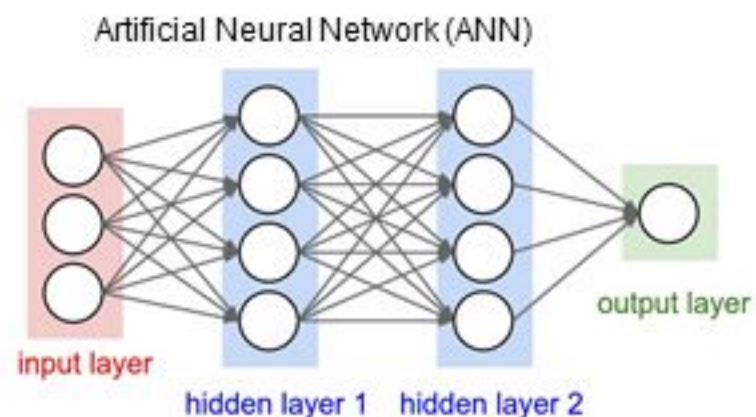
The $g(z)$ is called the **activation function**. Common choices include

- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

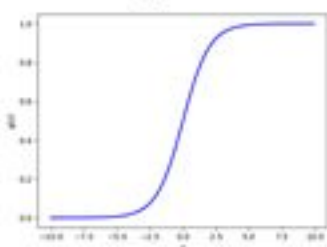
Find best **weights** W and **bias** b for each layer to minimize some **loss** $\ell(\dots)$.

Optimization non-convex so use local gradient learning methods.

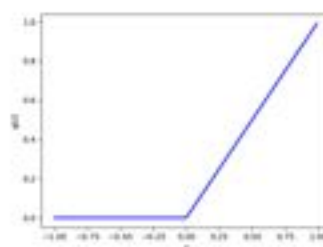
Universal approximation of any smooth function $y = f(x)$ just **two layers sufficient** for broad class of $g(z)$. **However**, may need many hidden units in layer.



Sigmoid



ReLU



What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

The $g(z)$ is called the **activation function**. Common choices include

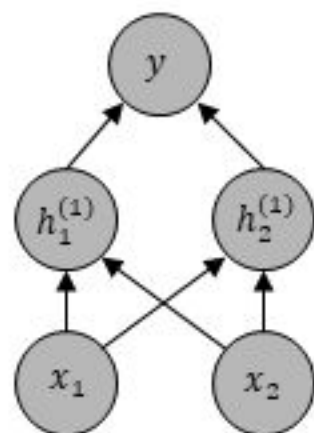
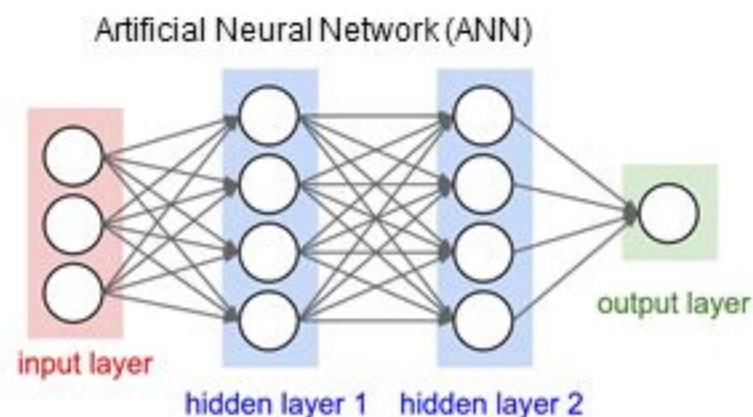
- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

Find best **weights** W and **bias** b for each layer to minimize some **loss** $\ell(\dots)$.

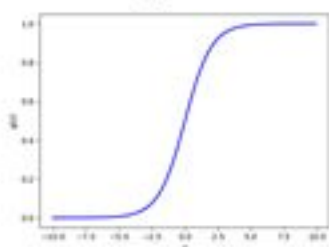
Optimization non-convex so use local gradient learning methods.

Universal approximation of any smooth function $y = f(x)$ just **two layers sufficient** for broad class of $g(z)$. **However**, may need many hidden units in layer.

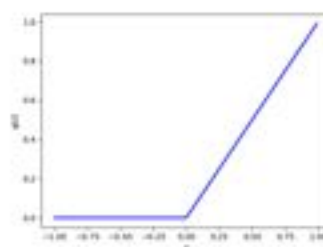
Deep architectures often less hidden units needed (symmetries). However, training can be more challenging for deep architectures.



Sigmoid



ReLU



What are Neural Networks (NNs)?

NN Transform: $y = f(\mathbf{x}; \theta)$ where f is obtained by compositions $f^{(L)}(f^{(L-1)}(\dots f^{(1)}(\mathbf{x}) \dots))$ with $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

The $g(z)$ is called the **activation function**. Common choices include

- **Sigmoid σ :** $g(z) = 1/(1 + e^{-z})$.
- **Rectified Linear Unit (ReLU):** $g(z) = \max(0, z)$.

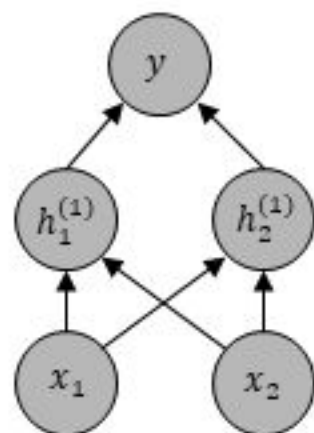
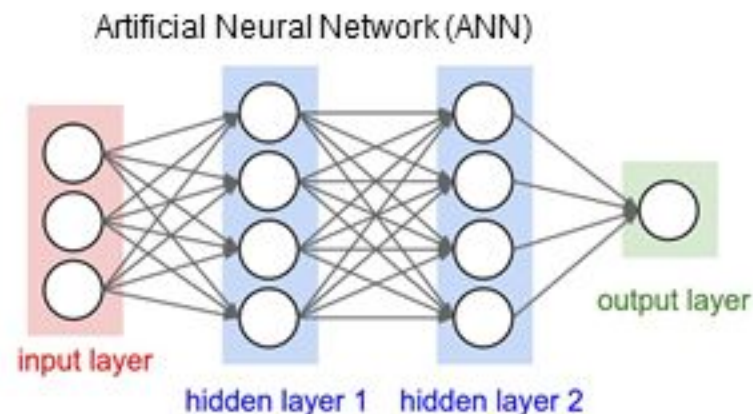
Find best **weights** W and **bias** b for each layer to minimize some **loss** $\ell(\dots)$.

Optimization non-convex so use local gradient learning methods.

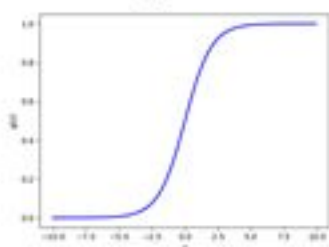
Universal approximation of any smooth function $y = f(x)$ just **two layers sufficient** for broad class of $g(z)$. **However**, may need many hidden units in layer.

Deep architectures often less hidden units needed (symmetries). However, training can be more challenging for deep architectures.

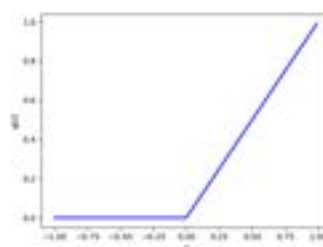
Many possible choices for network architectures, depth, activation functions.



Sigmoid



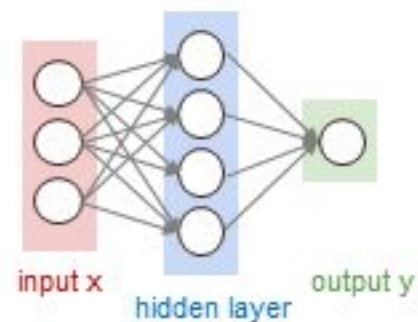
ReLU



Universal Approximation Theorem

For many activation functions $g(z)$ just two layers is sufficient for universal approximation of any continuous function $y = f(x)$ on a compact set.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

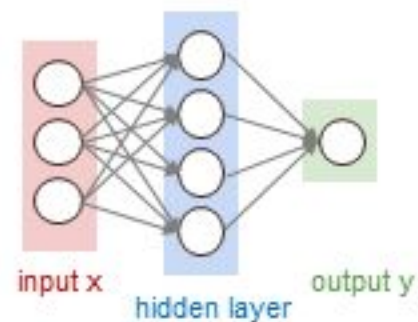
Universal Approximation Theorem

For many activation functions $g(z)$ just two layers is sufficient for universal approximation of any continuous function $y = f(x)$ on a compact set.

Definition: We say a subspace V **zeros-out** a measure μ . If for all $v \in V$, $\int v(x) d\mu(x) = 0$ holds then the measure must be zero $\mu \equiv 0$.

Definition: We say a function $g(z)$ is **discriminatory** if for all \mathbf{w}, b , we have $\int g(\mathbf{w}^T \mathbf{x} + b) d\mu(x) = 0$ then the measure must be zero $\mu \equiv 0$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

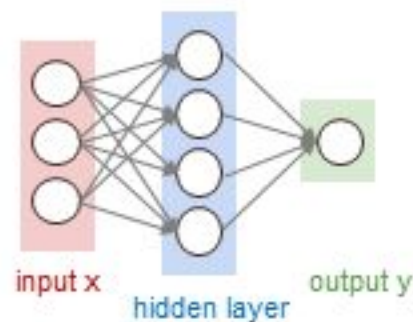
For many activation functions $g(z)$ just two layers is sufficient for universal approximation of any continuous function $y = f(x)$ on a compact set.

Definition: We say a subspace V **zeros-out** a measure μ . If for all $v \in V$, $\int v(x) d\mu(x) = 0$ holds then the measure must be zero $\mu \equiv 0$.

Definition: We say a function $g(z)$ is **discriminatory** if for all \mathbf{w}, b , we have $\int g(\mathbf{w}^T \mathbf{x} + b) d\mu(x) = 0$ then the measure must be zero $\mu \equiv 0$.

Lemma: The $g(z)$ is discriminatory for Borel measures μ iff the subspace $V = \{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ zeros-out the measures μ .

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

For many activation functions $g(z)$ just two layers is sufficient for universal approximation of any continuous function $y = f(x)$ on a compact set.

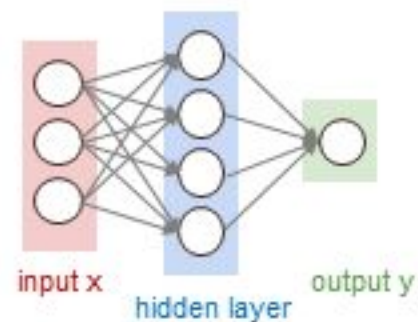
Definition: We say a subspace V **zeros-out** a measure μ . If for all $v \in V$, $\int v(x) d\mu(x) = 0$ holds then the measure must be zero $\mu \equiv 0$.

Definition: We say a function $g(z)$ is **discriminatory** if for all \mathbf{w}, b , we have $\int g(\mathbf{w}^T \mathbf{x} + b) d\mu(x) = 0$ then the measure must be zero $\mu \equiv 0$.

Lemma: The $g(z)$ is discriminatory for Borel measures μ iff the subspace $V = \{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ zeros-out the measures μ .

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V = \{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense in the space of continuous functions $C(I_n)$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

For many activation functions $g(z)$ just two layers is sufficient for universal approximation of any continuous function $y = f(x)$ on a compact set.

Definition: We say a subspace V **zeros-out** a measure μ . If for all $v \in V$, $\int v(x) d\mu(x) = 0$ holds then the measure must be zero $\mu \equiv 0$.

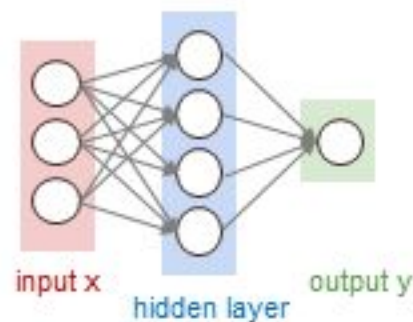
Definition: We say a function $g(z)$ is **discriminatory** if for all \mathbf{w}, b , we have $\int g(\mathbf{w}^T \mathbf{x} + b) d\mu(x) = 0$ then the measure must be zero $\mu \equiv 0$.

Lemma: The $g(z)$ is discriminatory for Borel measures μ iff the subspace $V = \{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ zeros-out the measures μ .

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V = \{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense in the space of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Neural Network: 1-Hidden Layer

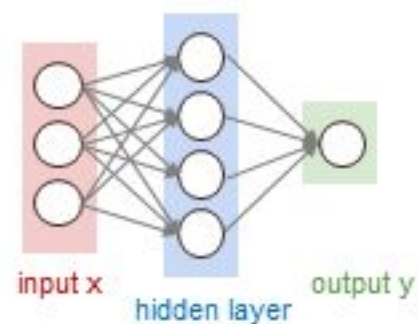


I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Neural Network: 1-Hidden Layer



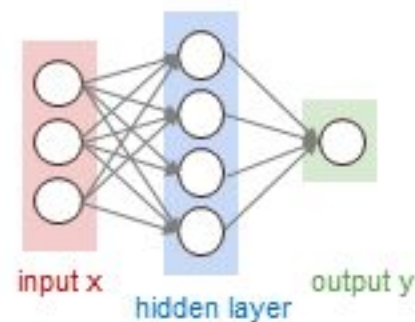
I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

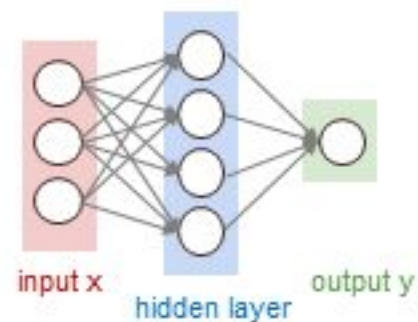
Universal Approximation Theorem

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

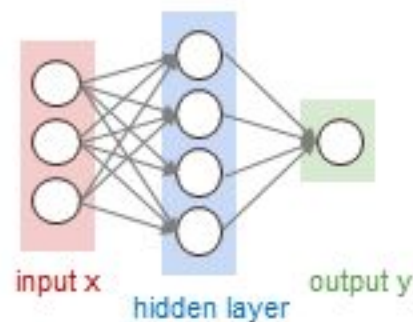
Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Follows from: $v(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j) = \sum_{j=1}^N \alpha_j (c_1 (\mathbf{w}_j^T \mathbf{x} + b_j) + c_2) = \mathbf{w}^T \mathbf{x} + b$, so
 $0 = \int v(x) d\mu(x) = \int (\mathbf{w}^T \mathbf{x} + b) d\mu(x) = \mathbf{w}^T \int \mathbf{x} d\mu(x) + b \int d\mu(x)$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

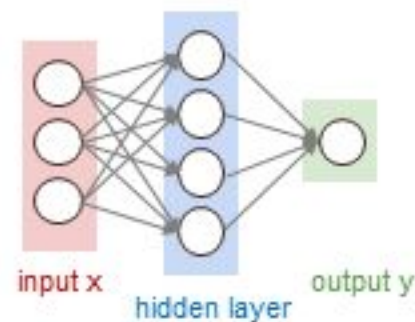
Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Follows from: $v(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j) = \sum_{j=1}^N \alpha_j (c_1 (\mathbf{w}_j^T \mathbf{x} + b_j) + c_2) = \mathbf{w}^T \mathbf{x} + b$, so $0 = \int v(x) d\mu(x) = \int (\mathbf{w}^T \mathbf{x} + b) d\mu(x) = \mathbf{w}^T \int \mathbf{x} d\mu(x) + b \int d\mu(x)$. This can be made to hold if we can find a Borel measure μ so that both $\int d\mu(x) = 0$ and $\int \mathbf{x} d\mu(x) = 0$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

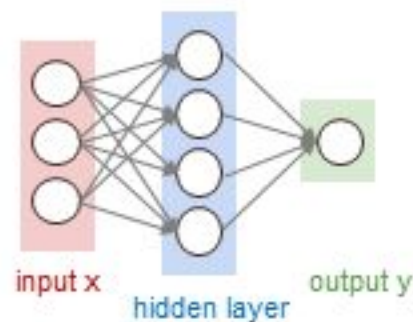
Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Follows from: $v(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j) = \sum_{j=1}^N \alpha_j (c_1 (\mathbf{w}_j^T \mathbf{x} + b_j) + c_2) = \mathbf{w}^T \mathbf{x} + b$, so $0 = \int v(x) d\mu(x) = \int (\mathbf{w}^T \mathbf{x} + b) d\mu(x) = \mathbf{w}^T \int \mathbf{x} d\mu(x) + b \int d\mu(x)$. This can be made to hold if we can find a Borel measure μ so that both $\int d\mu(x) = 0$ and $\int \mathbf{x} d\mu(x) = 0$. Let $\mu(x) = a_1 \delta(x_1 - r_1) + a_2 \delta(x_1 - r_2) + a_3 \delta(x_1 - r_2) = -\delta(x_1) + 2\delta(x_1 - \frac{1}{2}) - \delta(x_1 - 1)$.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

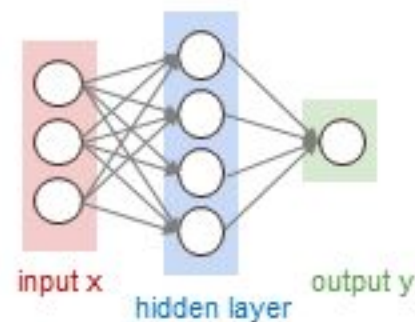
Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Follows from: $v(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j) = \sum_{j=1}^N \alpha_j (c_1 (\mathbf{w}_j^T \mathbf{x} + b_j) + c_2) = \mathbf{w}^T \mathbf{x} + b$, so $0 = \int v(x) d\mu(x) = \int (\mathbf{w}^T \mathbf{x} + b) d\mu(x) = \mathbf{w}^T \int \mathbf{x} d\mu(x) + b \int d\mu(x)$. This can be made to hold if we can find a Borel measure μ so that both $\int d\mu(x) = 0$ and $\int \mathbf{x} d\mu(x) = 0$. Let $\mu(x) = a_1 \delta(x_1 - r_1) + a_2 \delta(x_1 - r_2) + a_3 \delta(x_1 - r_2) = -\delta(x_1) + 2\delta(x_1 - \frac{1}{2}) - \delta(x_1 - 1)$.

Example: ReLU activations generate subspace that zeros-out Borel measures, so ReLU-NN's have the universal approximation property.

Neural Network: 1-Hidden Layer



I_n is the unit cube in \mathbb{R}^n

Universal Approximation Theorem

Theorem (Cybenko 1989): Let $g(z)$ be a continuous activation function that generates a subspace $V =$ in the space $\{q \mid q(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j), N \in \mathbb{N}\}$ that zeros-out all Borel measures on I_n , then V is dense of continuous functions $C(I_n)$.

Remark: For any function $f \in C(I_n)$ and $\epsilon > 0$ there exists $q \in V$ with N and weights $\alpha_j, \mathbf{w}_j^T, b_j$ such that $|f(x) - q(x)| < \epsilon$ for all $x \in I_n$.

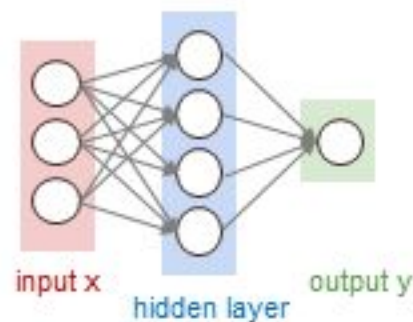
Example: Linear/affine activation functions are **not discriminatory**, $g(z) = c_1 z + c_2$.

Follows from: $v(x) = \sum_{j=1}^N \alpha_j g(\mathbf{w}_j^T \mathbf{x} + b_j) = \sum_{j=1}^N \alpha_j (c_1 (\mathbf{w}_j^T \mathbf{x} + b_j) + c_2) = \mathbf{w}^T \mathbf{x} + b$, so $0 = \int v(x) d\mu(x) = \int (\mathbf{w}^T \mathbf{x} + b) d\mu(x) = \mathbf{w}^T \int \mathbf{x} d\mu(x) + b \int d\mu(x)$. This can be made to hold if we can find a Borel measure μ so that both $\int d\mu(x) = 0$ and $\int \mathbf{x} d\mu(x) = 0$. Let $\mu(x) = a_1 \delta(x_1 - r_1) + a_2 \delta(x_1 - r_2) + a_3 \delta(x_1 - r_2) = -\delta(x_1) + 2\delta(x_1 - \frac{1}{2}) - \delta(x_1 - 1)$.

Example: ReLU activations generate subspace that zeros-out Borel measures, so ReLU-NN's have the universal approximation property.

ReLU-NN's include all pieces-wise linear approximations, realizable with enough hidden nodes, also provides another way to prove the universal approximation.

Neural Network: 1-Hidden Layer

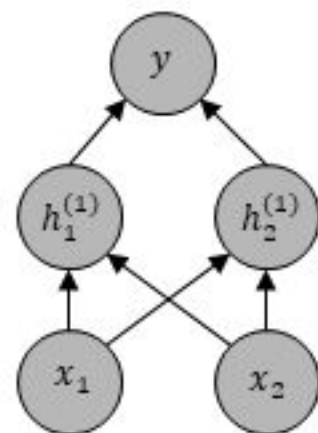
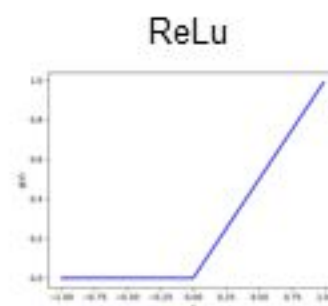
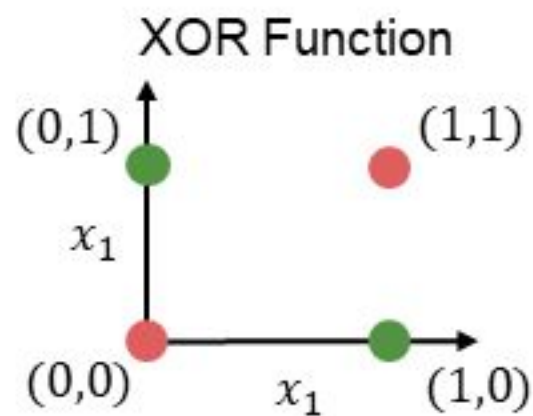


I_n is the unit cube in \mathbb{R}^n

Neural Networks: Example

NN Transform: $y = f(\mathbf{x}; \theta)$ from compositions $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}\mathbf{W}^k + b^k)$.

Example: Compute the XOR function $y = f(\mathbf{x}) = x_1 \oplus x_2$.



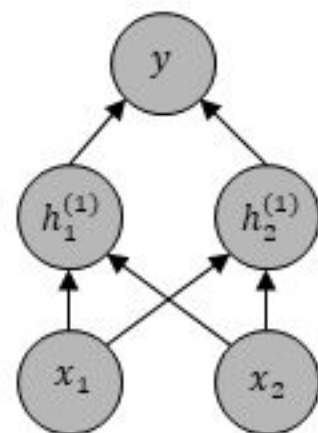
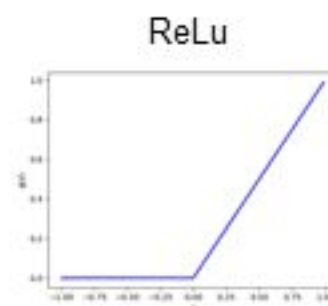
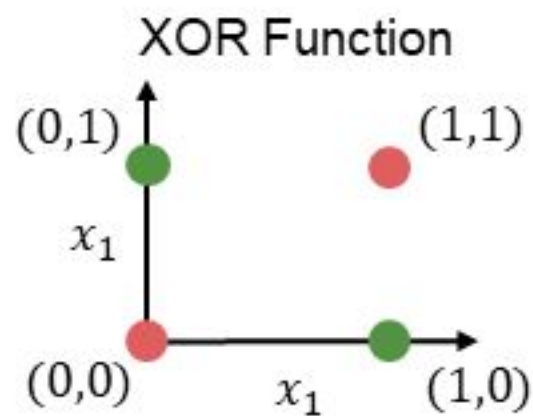
Neural Networks: Example

NN Transform: $y = f(\mathbf{x}; \theta)$ from compositions $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

Example: Compute the XOR function $y = f(\mathbf{x}) = x_1 \oplus x_2$.

Linear models are insufficient $h(\mathbf{x}) = \text{sign}(\mathbf{x}^T W + b)$, no choice W, b works.

Non-linearity important to approximate functions such as XOR.



Neural Networks: Example

NN Transform: $y = f(\mathbf{x}; \theta)$ from compositions $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

Example: Compute the XOR function $y = f(\mathbf{x}) = x_1 \oplus x_2$.

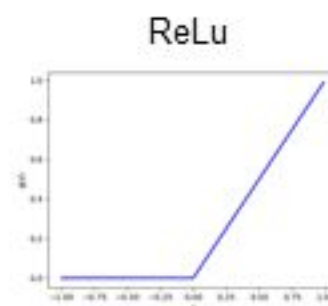
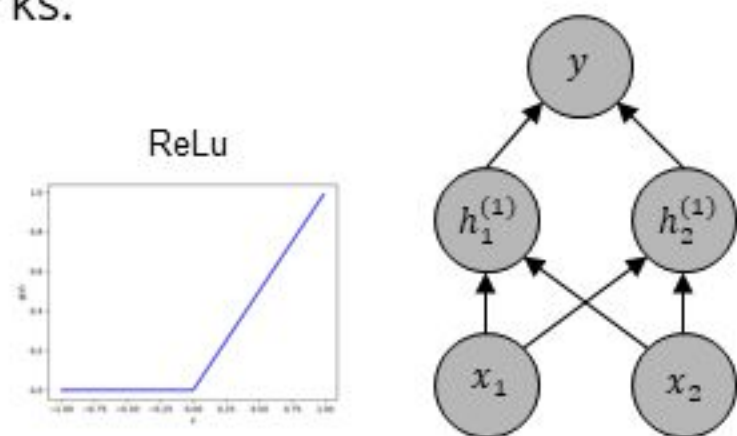
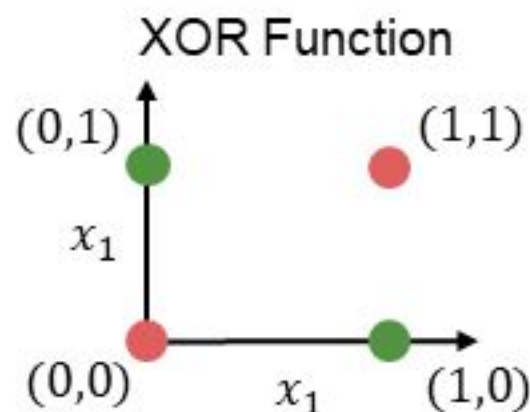
Linear models are insufficient $h(\mathbf{x}) = \text{sign}(\mathbf{x}^T W + b)$, no choice W, b works.

Non-linearity important to approximate functions such as XOR.

NN with two layers using non-linear ReLU activation g yields

$$\tilde{y} = f(\mathbf{x}; \theta) = (\max(0, \mathbf{x}^T W^{(1)} + b^{(1)}))^T W^{(2)} + b^{(2)}$$

Find parameters $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$ to try to obtain correct classification $y = \text{sign}(\tilde{y})$.



Neural Networks: Example

NN Transform: $y = f(\mathbf{x}; \theta)$ from compositions $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

Example: Compute the XOR function $y = f(\mathbf{x}) = x_1 \oplus x_2$.

Linear models are insufficient $h(\mathbf{x}) = \text{sign}(\mathbf{x}^T W + b)$, no choice W, b works.

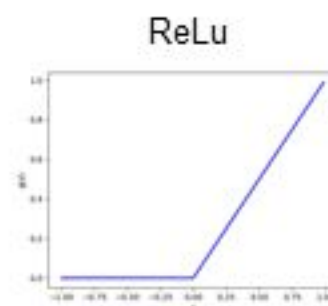
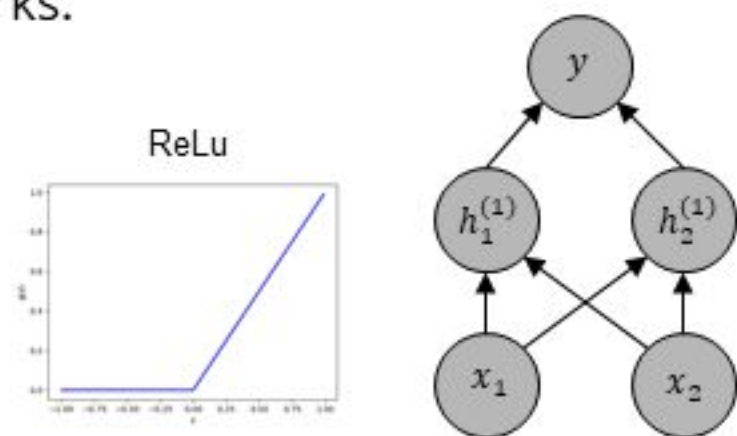
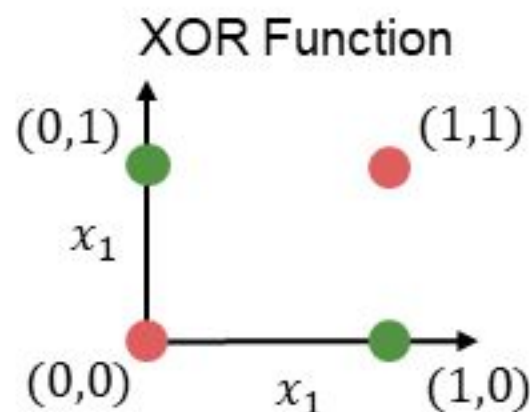
Non-linearity important to approximate functions such as XOR.

NN with two layers using non-linear ReLU activation g yields

$$\tilde{y} = f(\mathbf{x}; \theta) = (\max(0, \mathbf{x}^T W^{(1)} + b^{(1)}))^T W^{(2)} + b^{(2)}$$

Find parameters $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$ to try to obtain correct classification $y = \text{sign}(\tilde{y})$.

$$W^{(1)} = \begin{bmatrix} +1 & +1 \\ +1 & +1 \end{bmatrix}, b^{(1)} = [0, -1], W^{(2)} = \begin{bmatrix} 2 \\ -4 \end{bmatrix}, b^{(2)} = [-1].$$



Neural Networks: Example

NN Transform: $y = f(\mathbf{x}; \theta)$ from compositions $\mathbf{h}^k = f^{(k)}(\mathbf{h}^{k-1}) = g(\mathbf{h}^{k-1}W^k + b^k)$.

Example: Compute the XOR function $y = f(\mathbf{x}) = x_1 \oplus x_2$.

Linear models are insufficient $h(\mathbf{x}) = \text{sign}(\mathbf{x}^T W + b)$, no choice W, b works.

Non-linearity important to approximate functions such as XOR.

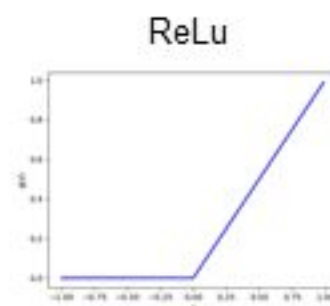
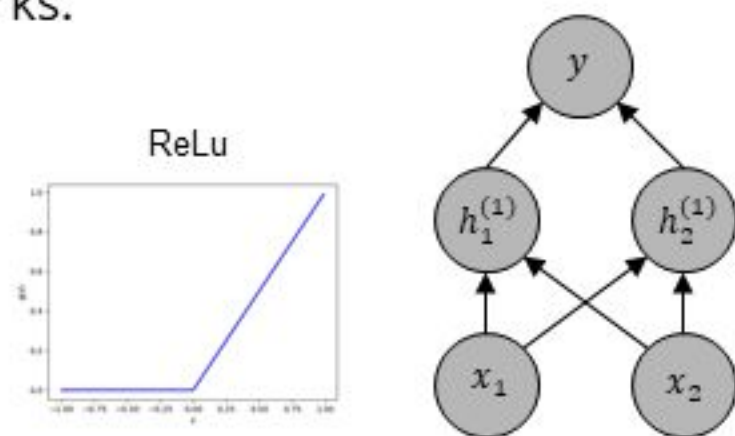
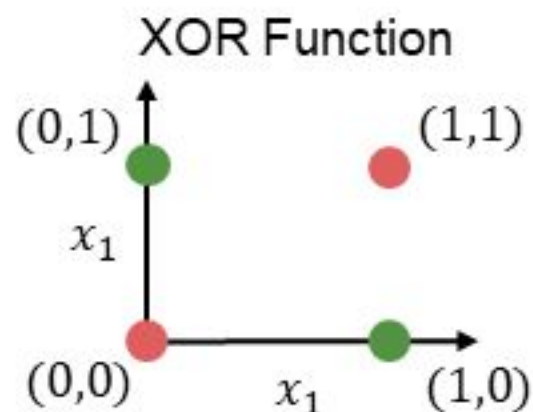
NN with two layers using non-linear ReLU activation g yields

$$\tilde{y} = f(\mathbf{x}; \theta) = (\max(0, \mathbf{x}^T W^{(1)} + b^{(1)}))^T W^{(2)} + b^{(2)}$$

Find parameters $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$ to try to obtain correct classification $y = \text{sign}(\tilde{y})$.

$$W^{(1)} = \begin{bmatrix} +1 & +1 \\ +1 & +1 \end{bmatrix}, b^{(1)} = [0, -1], W^{(2)} = \begin{bmatrix} 2 \\ -4 \end{bmatrix}, b^{(2)} = [-1].$$

In general, we need methods to learn from data such weights to minimize a loss function.





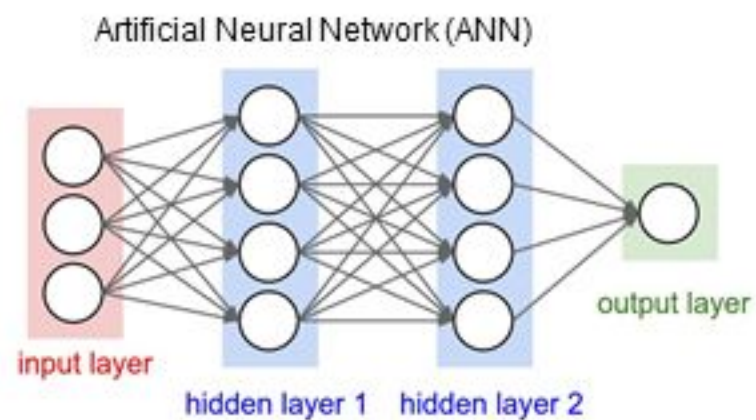
Learning with Neural Networks



How do we find weights of Neural Networks?

Optimization Problem:

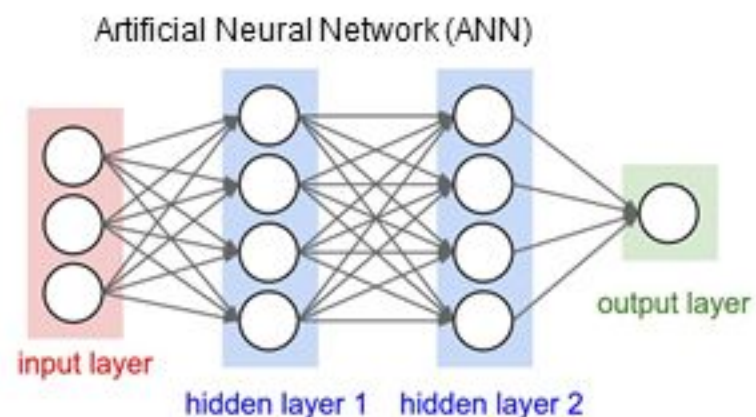
$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$



How do we find weights of Neural Networks?

Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

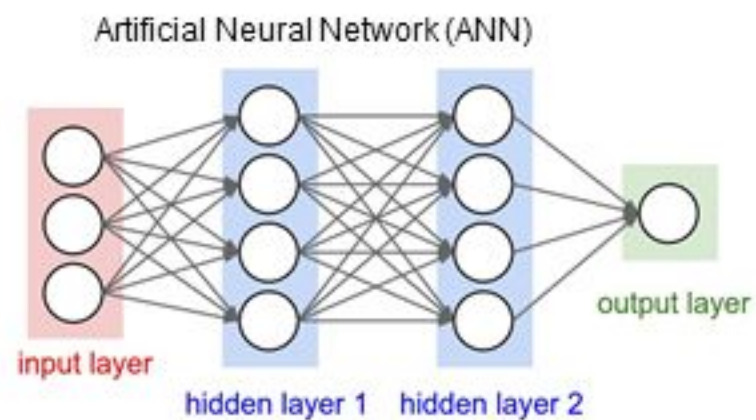


Non-convex problems typically have non-unique solutions and many local minima.

How do we find weights of Neural Networks?

Optimization Problem:

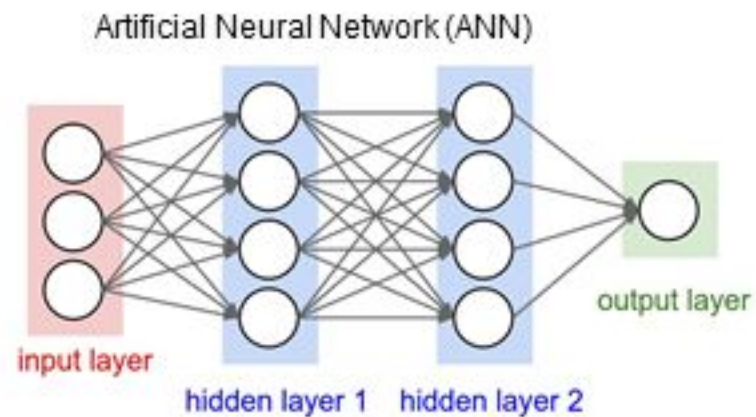
$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$



Non-convex problems typically have non-unique solutions and many local minima.

Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

How do we find weights of Neural Networks?



Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

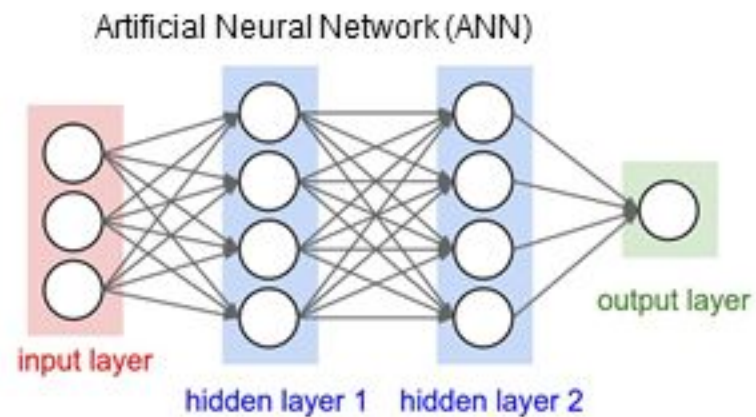
Non-convex problems typically have non-unique solutions and many local minima.

Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

Stochastic Gradient Descent:

$$\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} Q^n(\theta^n), \text{ with } Q^n(\theta^n) = Q^n(\mathbf{X}; \theta^n) = \frac{1}{m_b} \sum_{k=1}^{m_b} \ell(y_{i_k}, f(\mathbf{x}_{i_k}; \theta^n)).$$

How do we find weights of Neural Networks?



Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

Non-convex problems typically have non-unique solutions and many local minima.

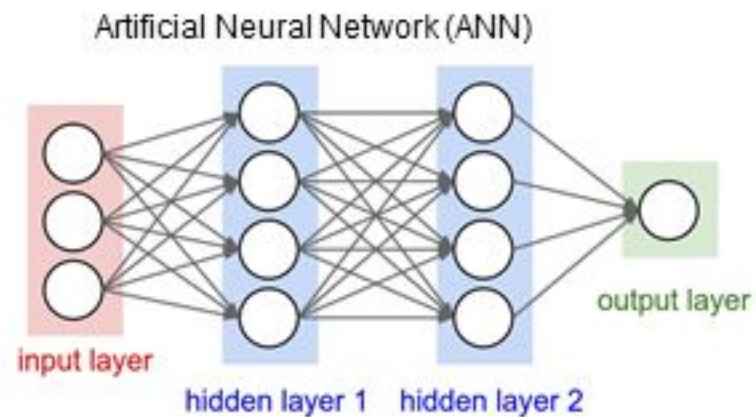
Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

Stochastic Gradient Descent:

$$\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} Q^n(\theta^n), \text{ with } Q^n(\theta^n) = Q^n(\mathbf{X}; \theta^n) = \frac{1}{m_b} \sum_{k=1}^{m_b} \ell(y_{i_k}, f(\mathbf{x}_{i_k}; \theta^n)).$$

A subset (batch) of the available data is used of size m_b to estimate expected loss, (provides regularization).

How do we find weights of Neural Networks?



Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

Non-convex problems typically have non-unique solutions and many local minima.

Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

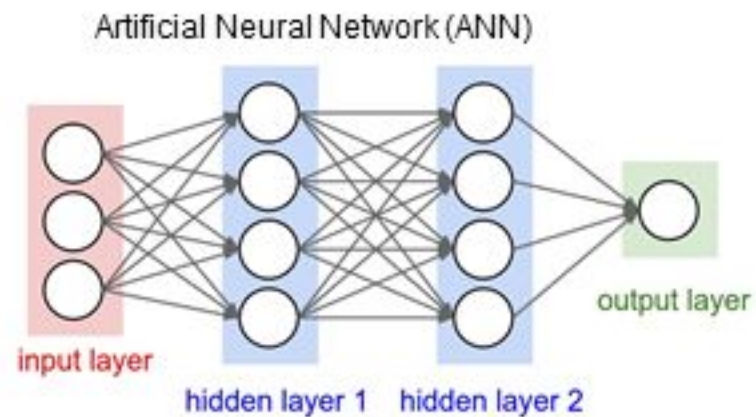
Stochastic Gradient Descent:

$$\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} Q^n(\theta^n), \text{ with } Q^n(\theta^n) = Q^n(\mathbf{X}; \theta^n) = \frac{1}{m_b} \sum_{k=1}^{m_b} \ell(y_{i_k}, f(\mathbf{x}_{i_k}; \theta^n)).$$

A subset (batch) of the available data is used of size m_b to estimate expected loss, (provides regularization).

Still needs computation of gradients $\nabla_{\theta} f(\mathbf{x}; \theta)$.

How do we find weights of Neural Networks?



Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

Non-convex problems typically have non-unique solutions and many local minima.

Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

Stochastic Gradient Descent:

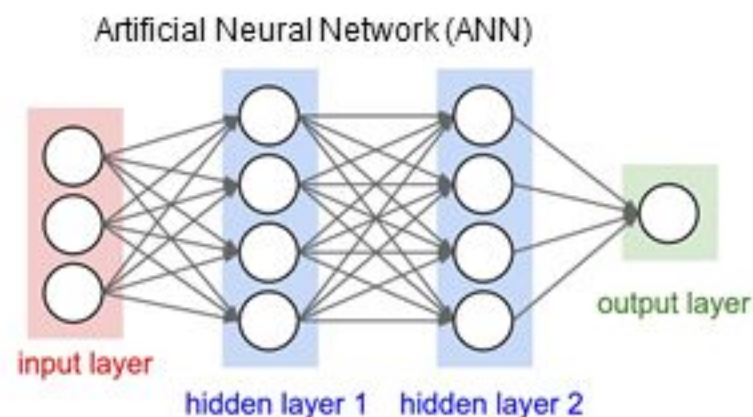
$$\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} Q^n(\theta^n), \text{ with } Q^n(\theta^n) = Q^n(\mathbf{X}; \theta^n) = \frac{1}{m_b} \sum_{k=1}^{m_b} \ell(y_{i_k}, f(\mathbf{x}_{i_k}; \theta^n)).$$

A subset (batch) of the available data is used of size m_b to estimate expected loss, (provides regularization).

Still needs computation of gradients $\nabla_{\theta} f(\mathbf{x}; \theta)$.

Analytically straight-forward to compute by chain-rule, but naïve evaluation is **computationally expensive**.

How do we find weights of Neural Networks?



Optimization Problem:

$$\min_{\theta} L(\theta; \{\mathbf{x}_i, y_i\}) = E_{\mathbf{x}, y \sim \tilde{D}_{data}} [\ell(y, f(\mathbf{x}; \theta))]$$

Non-convex problems typically have non-unique solutions and many local minima.

Goal is to find sets of parameters with small loss. Gradient-based methods can be used.

Stochastic Gradient Descent:

$$\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} Q^n(\theta^n), \text{ with } Q^n(\theta^n) = Q^n(\mathbf{X}; \theta^n) = \frac{1}{m_b} \sum_{k=1}^{m_b} \ell(y_{i_k}, f(\mathbf{x}_{i_k}; \theta^n)).$$

A subset (batch) of the available data is used of size m_b to estimate expected loss, (provides regularization).

Still needs computation of gradients $\nabla_{\theta} f(\mathbf{x}; \theta)$.

Analytically straight-forward to compute by chain-rule, but naïve evaluation is **computationally expensive**.

Automatic differentiation used in practice called **back-propagation**.



Computational Graphs and Back-Propagation

Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(\mathbf{X}; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.

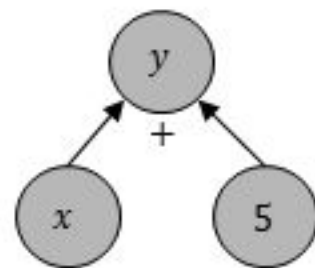
Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(\mathbf{X}; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.



Computational Graphs

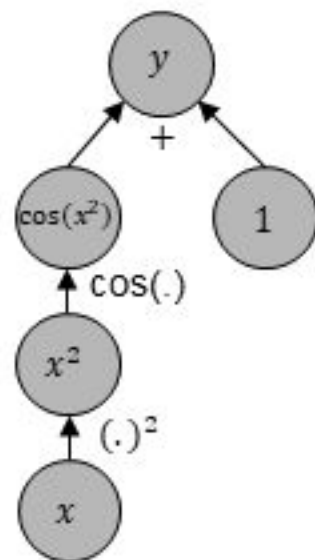
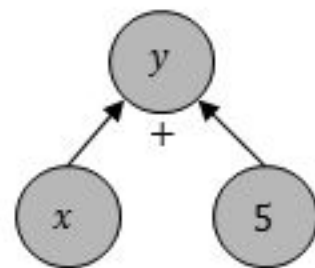
Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.

Example: $y = f(x) = \cos(x^2) + 1$.



Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

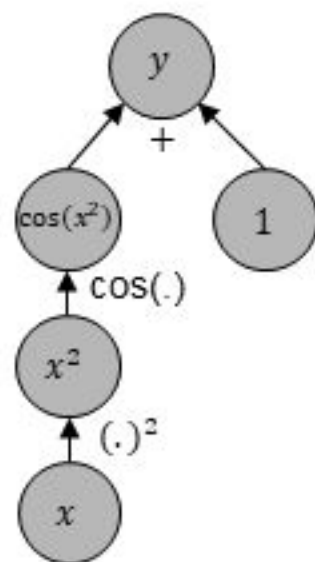
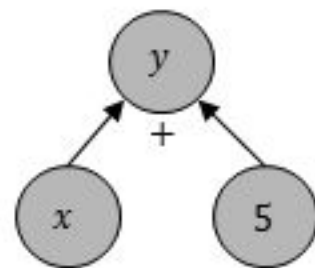
Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.

Example: $y = f(x) = \cos(x^2) + 1$.

Example: $y = f(x_1, x_2) = w_1^{(2)} g(w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2) + w_2^{(2)} g(w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2)$



Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

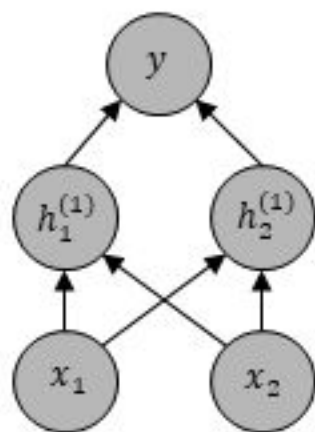
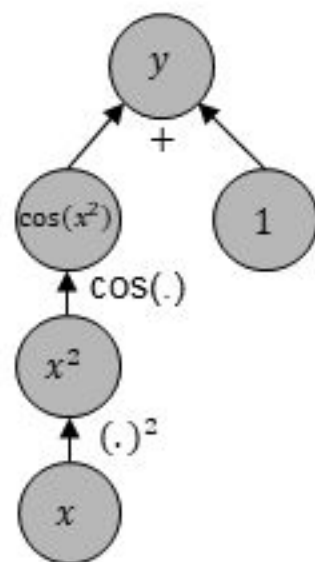
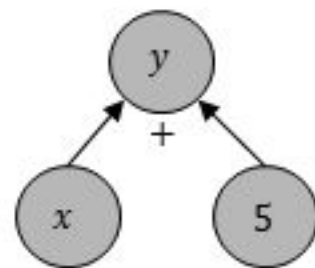
Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.

Example: $y = f(x) = \cos(x^2) + 1$.

Example: $y = f(x_1, x_2) = w_1^{(2)} g(w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2) + w_2^{(2)} g(w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2)$

$$y = w_1^{(2)} h_1^{(1)} + w_2^{(2)} h_2^{(1)}$$



Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

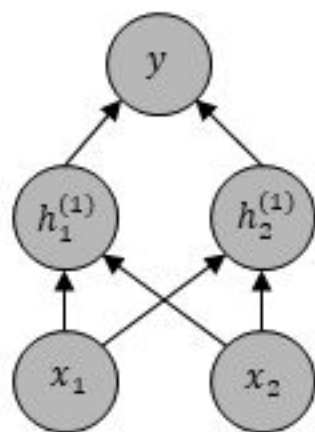
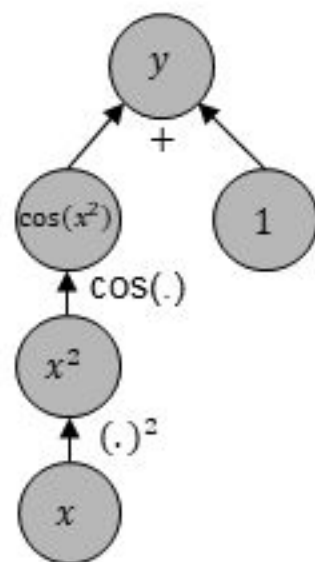
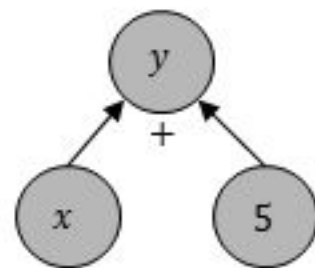
Example: $y = f(x) = x + 5$.

Example: $y = f(x) = \cos(x^2) + 1$.

Example: $y = f(x_1, x_2) = w_1^{(2)} g(w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2) + w_2^{(2)} g(w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2)$

$$y = w_1^{(2)} h_1^{(1)} + w_2^{(2)} h_2^{(1)}$$

Gradient can be computed provided we know how to differentiate result of each operation in terms of contributing terms.



Computational Graphs

Optimization methods often need gradients $\nabla_{\theta} f(X; \theta)$.

Symbolic representations of $f(\theta)$ useful for automatic differentiation to obtain $\nabla_{\theta} f$.

Computational graph represents function evaluation in terms of more basic operations.

Example: $y = f(x) = x + 5$.

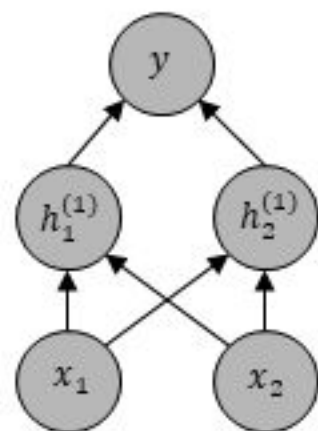
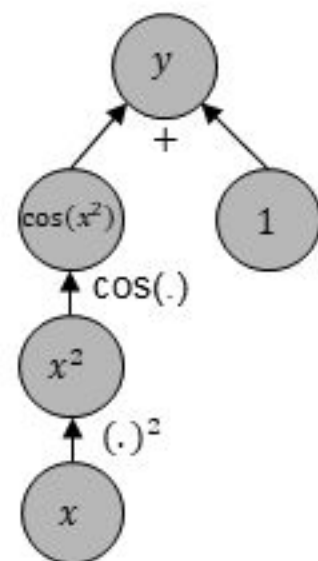
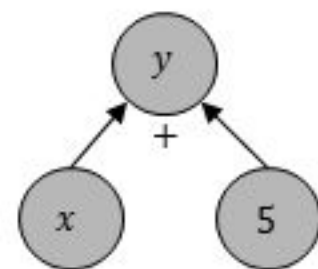
Example: $y = f(x) = \cos(x^2) + 1$.

Example: $y = f(x_1, x_2) = w_1^{(2)} g(w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2) + w_2^{(2)} g(w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2)$

$$y = w_1^{(2)} h_1^{(1)} + w_2^{(2)} h_2^{(1)}$$

Gradient can be computed provided we know how to differentiate result of each operation in terms of contributing terms.

Function derivatives can then be built up using the **chain-rule of calculus**.



Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$

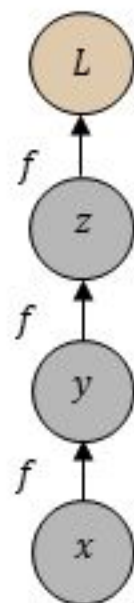
Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$



Chain-Rule

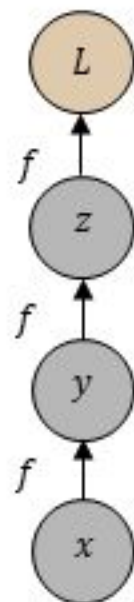
Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$, can also be expressed as (i) $\frac{\partial L}{\partial x} = f'(f(f(x)))f'(f(x))f'(x)$ or (ii) $\frac{\partial L}{\partial x} = f'(z)f'(y)f'(x)$.



Chain-Rule

Chain-Rule of Calculus:

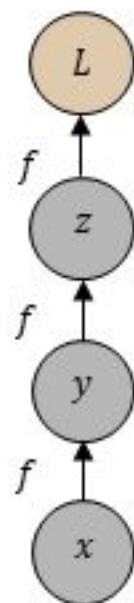
$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$, can also be expressed as (i) $\frac{\partial L}{\partial x} = f'(f(f(x)))f'(f(x))f'(x)$ or (ii) $\frac{\partial L}{\partial x} = f'(z)f'(y)f'(x)$.

First expression (i) "naive" chain-rule can require many function evaluations (exponential number some cases).



Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

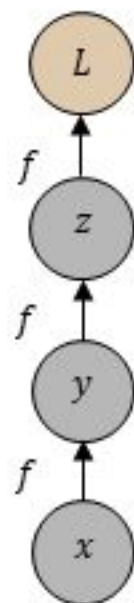
$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$, can also be expressed as (i) $\frac{\partial L}{\partial x} = f'(f(f(x)))f'(f(x))f'(x)$ or (ii) $\frac{\partial L}{\partial x} = f'(z)f'(y)f'(x)$.

First expression (i) "naive" chain-rule can require many function evaluations (exponential number some cases).

Second expression (ii) composite chain-rule reuses previous functional evaluations (cost in memory).



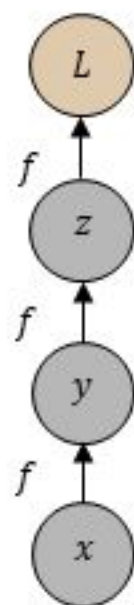
Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$



$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$, can also be expressed as (i) $\frac{\partial L}{\partial x} = f'(f(f(x)))f'(f(x))f'(x)$ or (ii) $\frac{\partial L}{\partial x} = f'(z)f'(y)f'(x)$.

First expression (i) "naive" chain-rule can require many function evaluations (exponential number some cases).

Second expression (ii) composite chain-rule reuses previous functional evaluations (cost in memory).

Computational can either assemble product to evaluate or store symbolic representation.

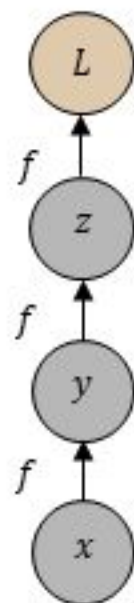
Chain-Rule

Chain-Rule of Calculus:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \quad \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Example: $L = f(z)$, $z = f(y)$, $y = f(x)$ we then have $L = f(f(f(x)))$



$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$, can also be expressed as (i) $\frac{\partial L}{\partial x} = f'(f(f(x)))f'(f(x))f'(x)$ or (ii) $\frac{\partial L}{\partial x} = f'(z)f'(y)f'(x)$.

First expression (i) "naive" chain-rule can require many function evaluations (exponential number some cases).

Second expression (ii) composite chain-rule reuses previous functional evaluations (cost in memory).

Computational can either assemble product to evaluate or store symbolic representation.

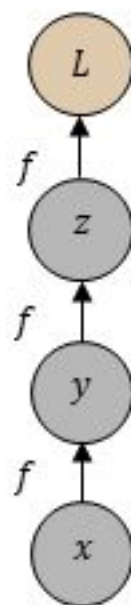
Advantages of (i) when memory storage issues, otherwise (ii) is usually preferred.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$



Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

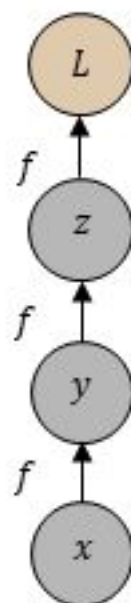
$$\mathbf{u}^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$\mathbf{u}^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $\mathbf{u}^{(n)}, \{u^{(i)}\}_{i=1}^n$.



Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

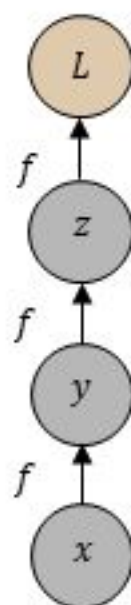
$$u^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$u^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $u^{(n)}, \{u^{(i)}\}_{i=1}^n$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

$$\mathbf{u}^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$\mathbf{u}^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $\mathbf{u}^{(n)}, \{u^{(i)}\}_{i=1}^n$.

Algorithm II (Backward-Pass):

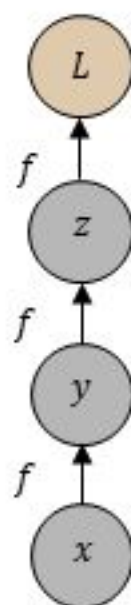
Input: $u^{(m)}, m = 1, \dots, n$.

$\text{grad_table}[u^{(n)}] \leftarrow 1$

For $j = n - 1, \dots, 1$

$$\text{grad_table}[u^{(j)}] \leftarrow \sum_{i: j \in \text{Pa}(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Output: $\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(k)}}, k = 1, \dots, n_I$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

$$u^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$u^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $u^{(n)}, \{u^{(i)}\}_{i=1}^n$.

Algorithm II (Backward-Pass):

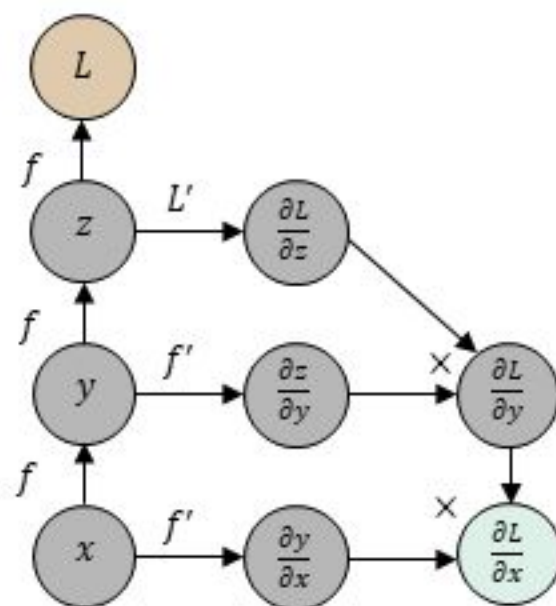
Input: $u^{(m)}, m = 1, \dots, n$.

$\text{grad_table}[u^{(n)}] \leftarrow 1$

For $j = n - 1, \dots, 1$

$$\text{grad_table}[u^{(j)}] \leftarrow \sum_{i: j \in \text{Pa}(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Output: $\frac{\partial u^{(n)}}{\partial u^{(k)}}, k = 1, \dots, n_I$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

$$u^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$u^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $u^{(n)}, \{u^{(i)}\}_{i=1}^n$.

Algorithm II (Backward-Pass):

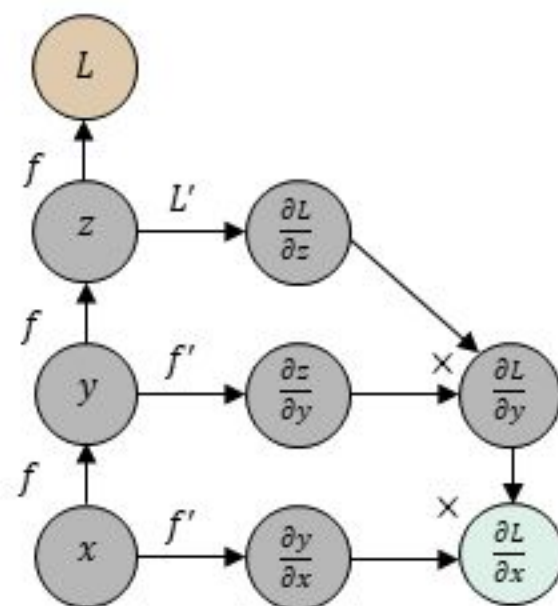
Input: $u^{(m)}, m = 1, \dots, n$.

$\text{grad_table}[u^{(n)}] \leftarrow 1$

For $j = n - 1, \dots, 1$

$$\text{grad_table}[u^{(j)}] \leftarrow \sum_{i: j \in \text{Pa}(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Output: $\frac{\partial u^{(n)}}{\partial u^{(k)}}, k = 1, \dots, n_I$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Algorithm II maintains at each stage: $\text{grad_table}[u^{(j)}] = \frac{\partial u^{(n)}}{\partial u^{(j)}}$.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

$$u^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$u^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $u^{(n)}, \{u^{(i)}\}_{i=1}^n$

Algorithm II (Backward-Pass):

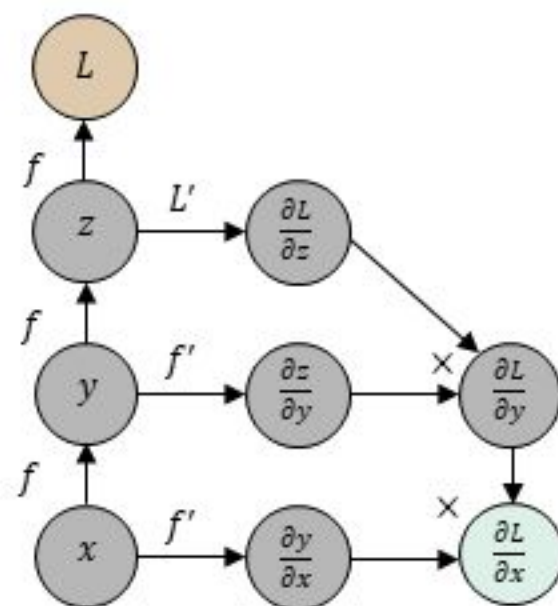
Input: $u^{(m)}, m = 1, \dots, n$.

$\text{grad_table}[u^{(n)}] \leftarrow 1$

For $j = n - 1, \dots, 1$

$$\text{grad_table}[u^{(j)}] \leftarrow \sum_{i: j \in \text{Pa}(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Output: $\frac{\partial u^{(n)}}{\partial u^{(k)}}, k = 1, \dots, n_I$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Algorithm II maintains at each stage: $\text{grad_table}[u^{(j)}] = \frac{\partial u^{(n)}}{\partial u^{(j)}}$.

Back-Propagation consists of the two steps (i) forward pass of algorithm I followed by (ii) backward pass of algorithm II.

Back-Propagation Method

Back-Propagation Method:

$$\mathbf{u}^{(k)} = f(\mathbf{u}^{(k-1)}), \mathbf{u}^{(n)} = f(f(\dots f(\mathbf{u}^{(1)}) \dots))$$

$$\frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(j)}} = \sum_{i: j \in \text{Pa}(\mathbf{u}^{(i)})} \frac{\partial \mathbf{u}^{(n)}}{\partial \mathbf{u}^{(i)}} \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{u}^{(j)}}$$

Algorithm I (Forward-Pass):

Input: $x^{(1)}, x^{(2)}, \dots, x^{(n_I)}$

for $k = 1, 2, \dots, n_I$

$$u^{(k)} \leftarrow x^{(k)}$$

for $m = n_I + 1, \dots, n$

$$\mathbb{U}^{(m)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(m)})\}$$

$$u^{(m)} \leftarrow f^{(m)}(\mathbb{U}^{(m)})$$

Output: $u^{(n)}, \{u^{(i)}\}_{i=1}^n$

Algorithm II (Backward-Pass):

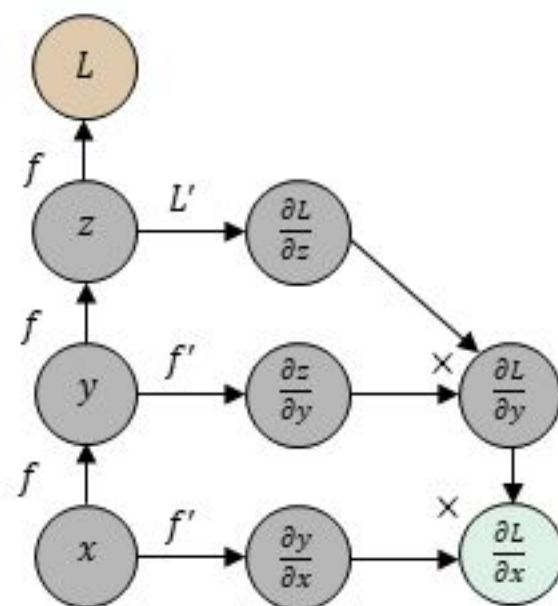
Input: $u^{(m)}, m = 1, \dots, n$.

$\text{grad_table}[u^{(n)}] \leftarrow 1$

For $j = n - 1, \dots, 1$

$$\text{grad_table}[u^{(j)}] \leftarrow \sum_{i: j \in \text{Pa}(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Output: $\frac{\partial u^{(n)}}{\partial u^{(k)}}, k = 1, \dots, n_I$.



Algorithm I computes the functional evaluations $\{u^{(i)}\}_{i=1}^n$.

Algorithm II maintains at each stage: $\text{grad_table}[u^{(j)}] = \frac{\partial u^{(n)}}{\partial u^{(j)}}$.

Back-Propagation consists of the two steps (i) forward pass of algorithm I followed by (ii) backward pass of algorithm II.

Parallelized versions and other variants also used for efficiency.



Feed-Forward Neural Networks (FFNNs)

Basic Examples of NN's

Feed-Forward Neural Networks

Neural Network Architecture

Neural network architecture with one processing layer feeding forward into the next processing layer.

Hidden Layer 1



Input Layer

Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

Neural Network Architecture

Hidden Layer 2



Hidden Layer 1



Input Layer

Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

Neural Network Architecture

Output Layer



Hidden Layer 2



Hidden Layer 1

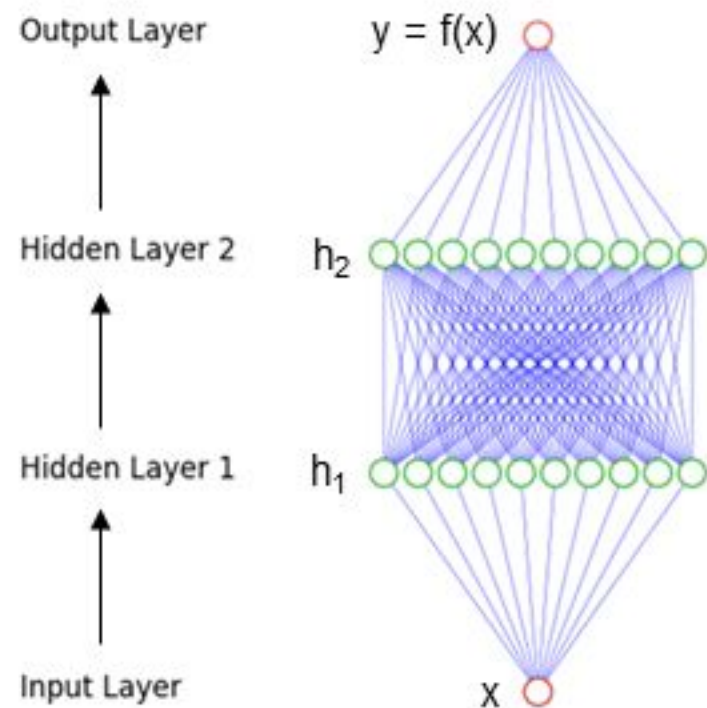


Input Layer

Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

Neural Network Architecture



Feed-Forward Neural Networks

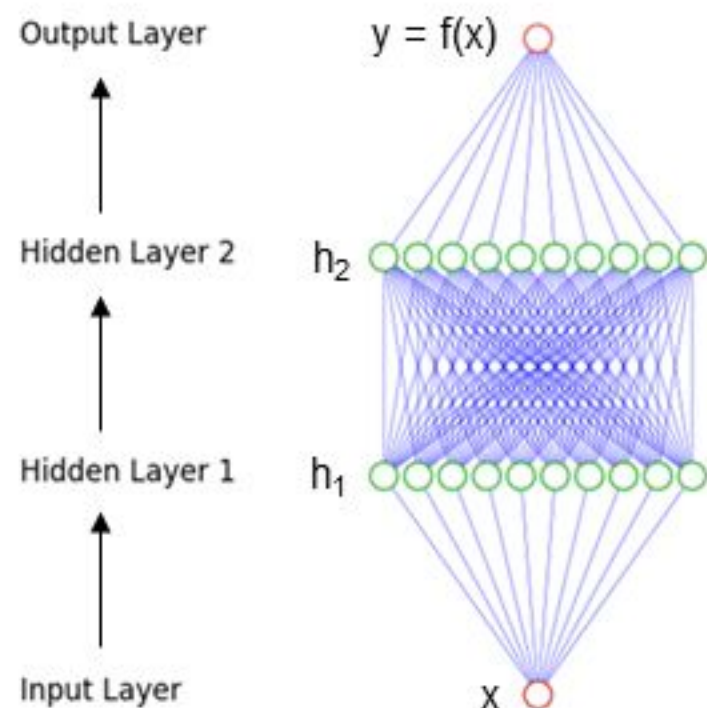
Neural network architecture with one processing layer feeding forward into the next processing layer.

Intermediate hidden processing layers of the form $g(XW + b)$.

Nonlinear transformation by some activation function $g(z)$.

Last processing layer typically is linear $XW + b$.

Neural Network Architecture



Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

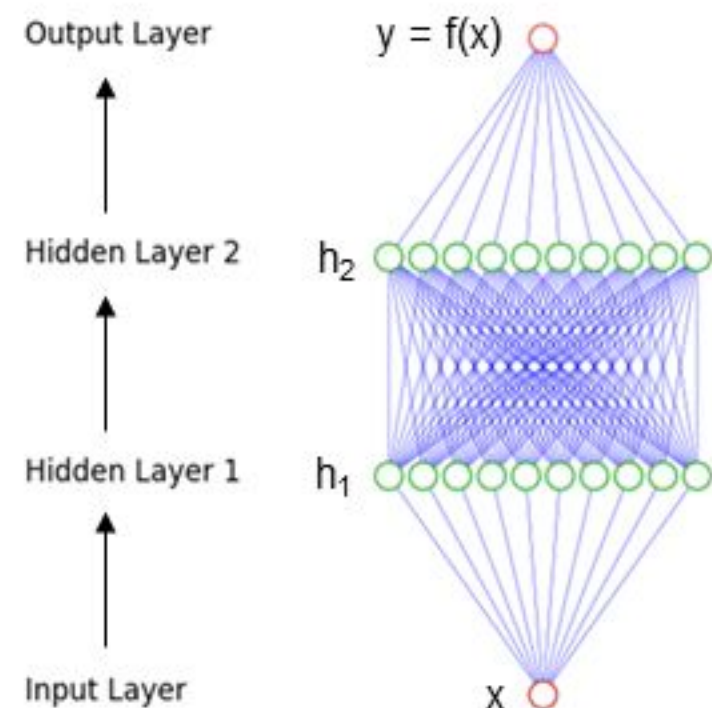
Intermediate hidden processing layers of the form $g(XW + b)$.

Nonlinear transformation by some activation function $g(z)$.

Last processing layer typically is linear $XW + b$.

Feed-Forward Neural Network (FFNN) provide model for $y = f(x)$.

Neural Network Architecture



Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

Intermediate hidden processing layers of the form $g(XW + b)$.

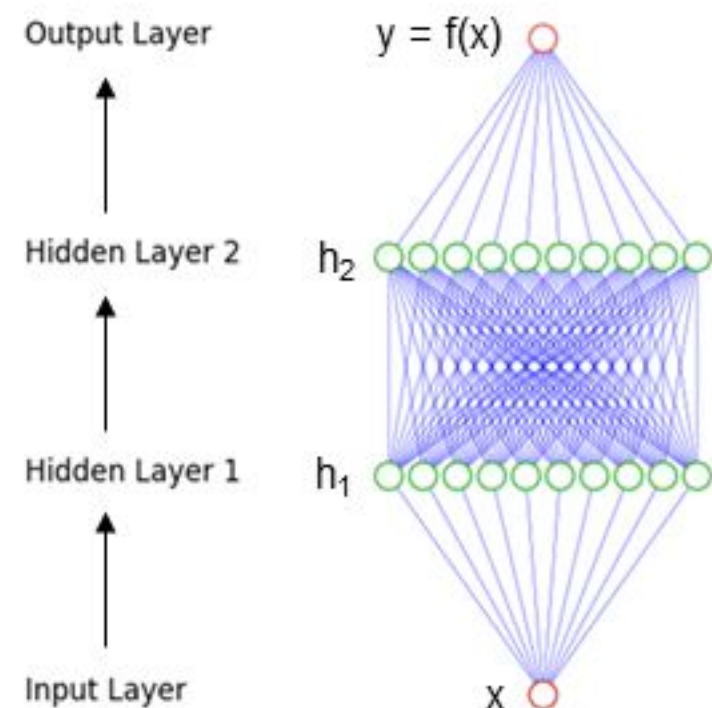
Nonlinear transformation by some activation function $g(z)$.

Last processing layer typically is linear $XW + b$.

Feed-Forward Neural Network (FFNN) provide model for $y = f(x)$.

Learning involves adjusting weights W and bias b of layers.

Neural Network Architecture



Feed-Forward Neural Networks

Neural network architecture with one processing layer feeding forward into the next processing layer.

Intermediate hidden processing layers of the form $g(XW + b)$.

Nonlinear transformation by some activation function $g(z)$.

Last processing layer typically is linear $XW + b$.

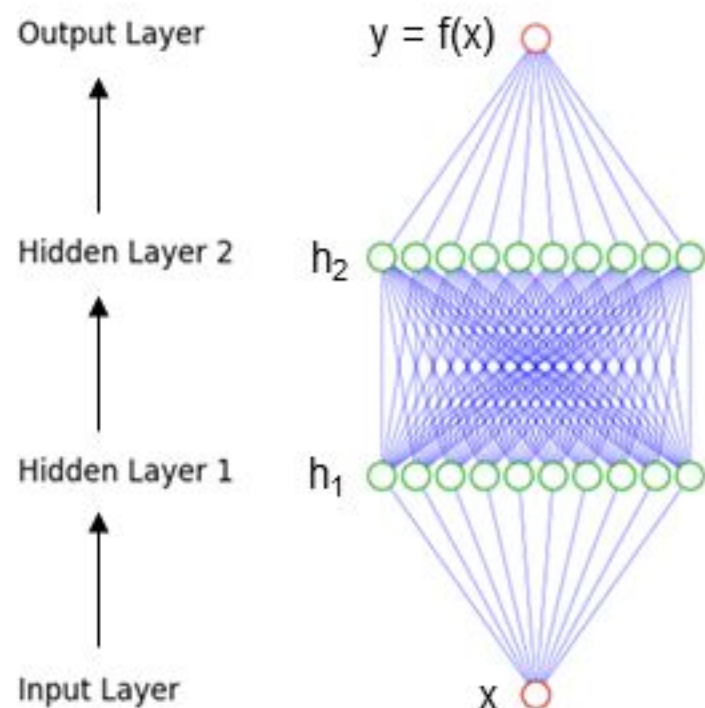
Feed-Forward Neural Network (FFNN) provide model for $y = f(x)$.

Learning involves adjusting weights W and bias b of layers.

Stochastic Gradient Descent (SGD) currently widely used for optimization of weights and bias.

Implicit regularization by choice of batch size and learning rates.

Neural Network Architecture



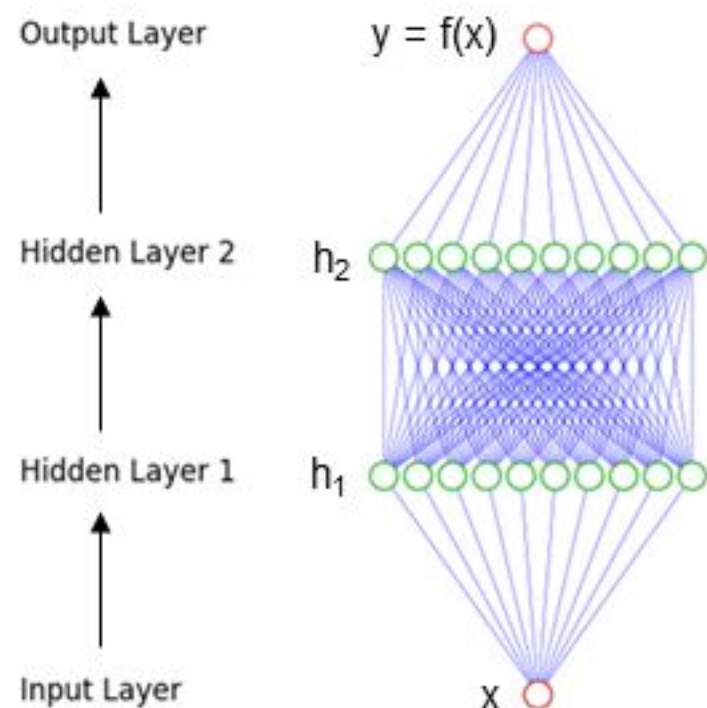
Feed-Forward Neural Networks: Example

Example: Approximate the function $y = \sin(x)$ using FFNN.

Architecture: 2-layers with 10-hidden ReLu nodes per layer.

This NN architecture spans piecewise linear functions (10 nodes).

Neural Network Architecture



Feed-Forward Neural Networks: Example

Example: Approximate the function $y = \sin(x)$ using FFNN.

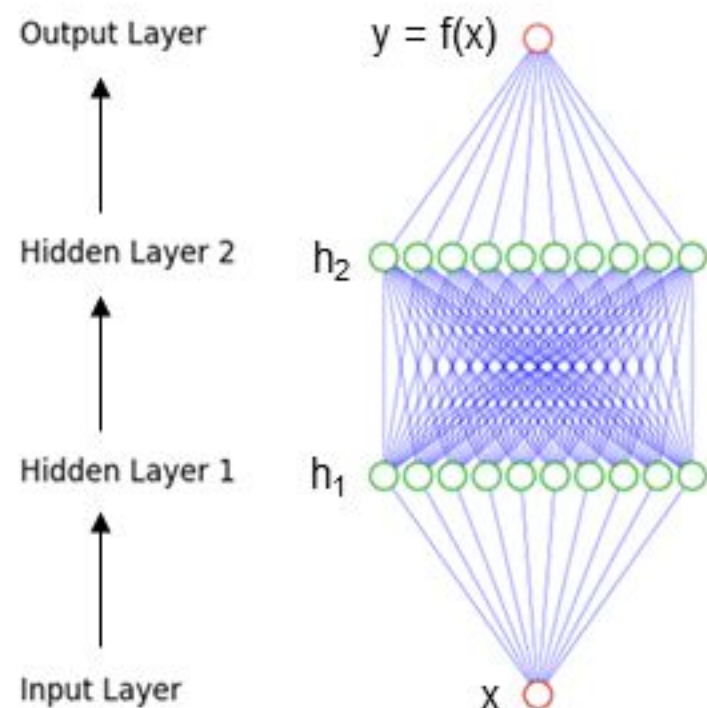
Architecture: 2-layers with 10-hidden ReLu nodes per layer.

This NN architecture spans piecewise linear functions (10 nodes).

Explicitly: $f(X) = g(g(g(X \cdot W^{(1)} + b^{(1)}) \cdot W^{(2)} + b^{(2)}) \cdot W^{(3)} + b^{(3)})$
with $g(z) = \max(0, z)$.

A notion of "loss" required to assess level of success in fit.

Neural Network Architecture



Feed-Forward Neural Networks: Example

Example: Approximate the function $y = \sin(x)$ using FFNN.

Architecture: 2-layers with 10-hidden ReLu nodes per layer.

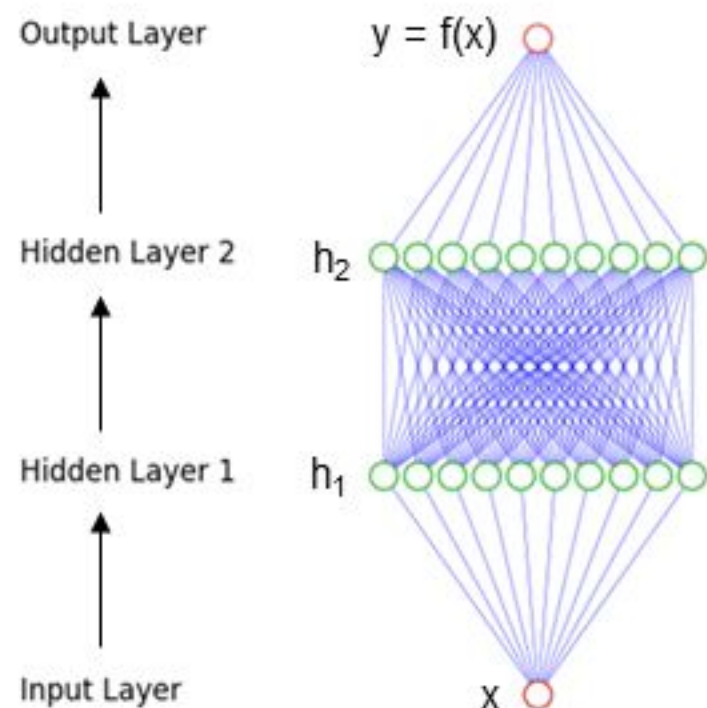
This NN architecture spans piecewise linear functions (10 nodes).

Explicitly: $f(X) = g(g(g(X \cdot W^{(1)} + b^{(1)}) \cdot W^{(2)} + b^{(2)}) \cdot W^{(3)} + b^{(3)})$
with $g(z) = \max(0, z)$.

A notion of "loss" required to assess level of success in fit.

Least-squares loss function $\ell(\{x_i, y_i\}) = \sum_i (f(x_i) - y_i)^2$.

Neural Network Architecture



Feed-Forward Neural Networks: Example

Example: Approximate the function $y = \sin(x)$ using FFNN.

Architecture: 2-layers with 10-hidden ReLu nodes per layer.

This NN architecture spans piecewise linear functions (10 nodes).

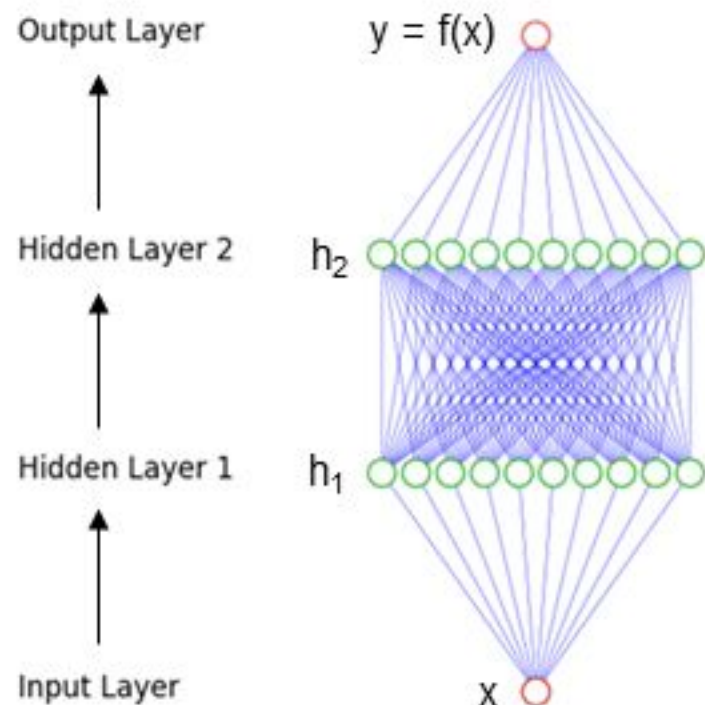
Explicitly: $f(X) = g(g(g(X \cdot W^{(1)} + b^{(1)}) \cdot W^{(2)} + b^{(2)}) \cdot W^{(3)} + b^{(3)})$
with $g(z) = \max(0, z)$.

A notion of "loss" required to assess level of success in fit.

Least-squares loss function $\ell(\{x_i, y_i\}) = \sum_i (f(x_i) - y_i)^2$.

Learning W, b proceeds by stochastic gradient descent.

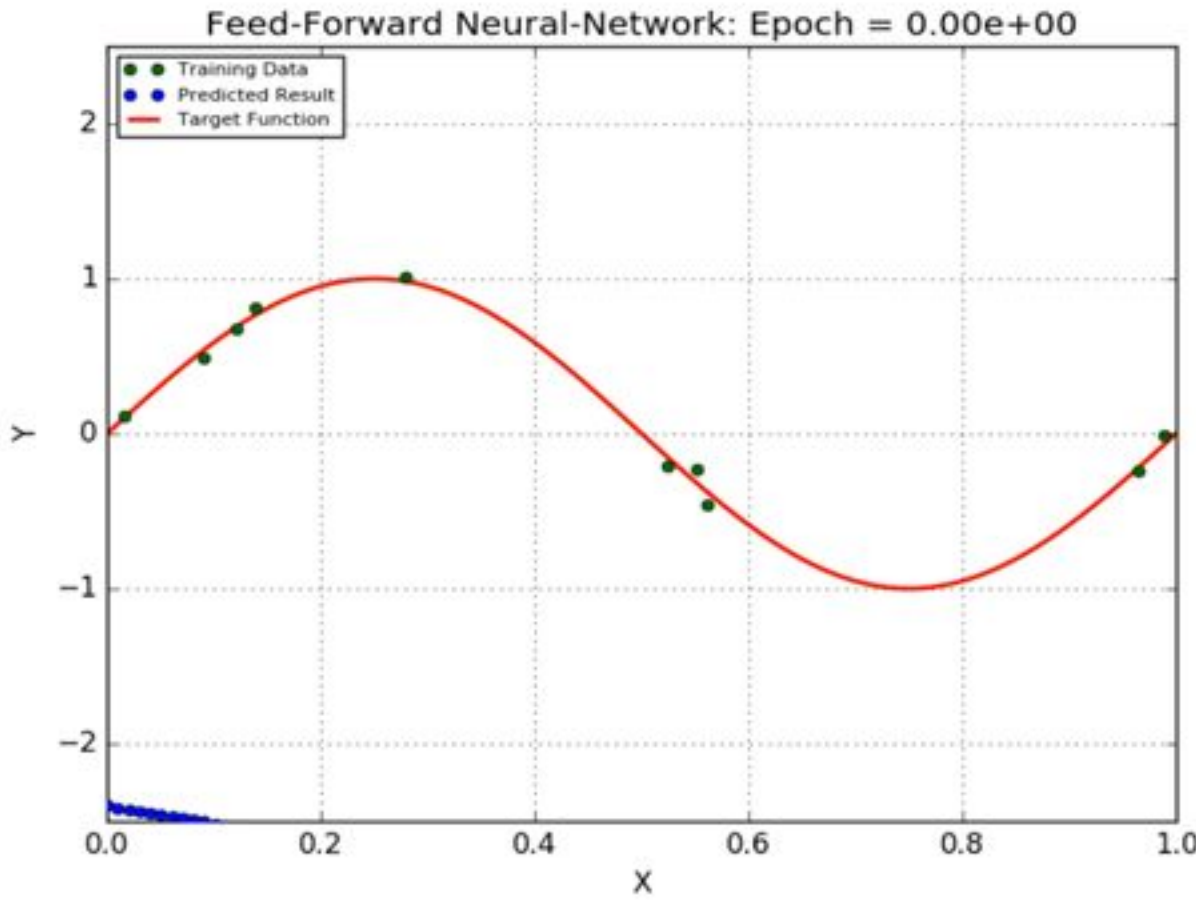
Neural Network Architecture



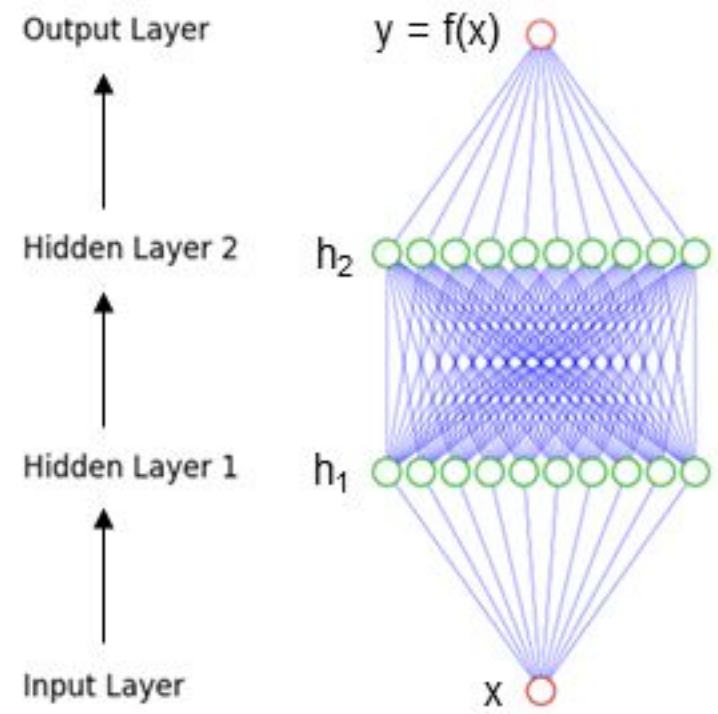
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 10 and learning rate 10^{-4} .

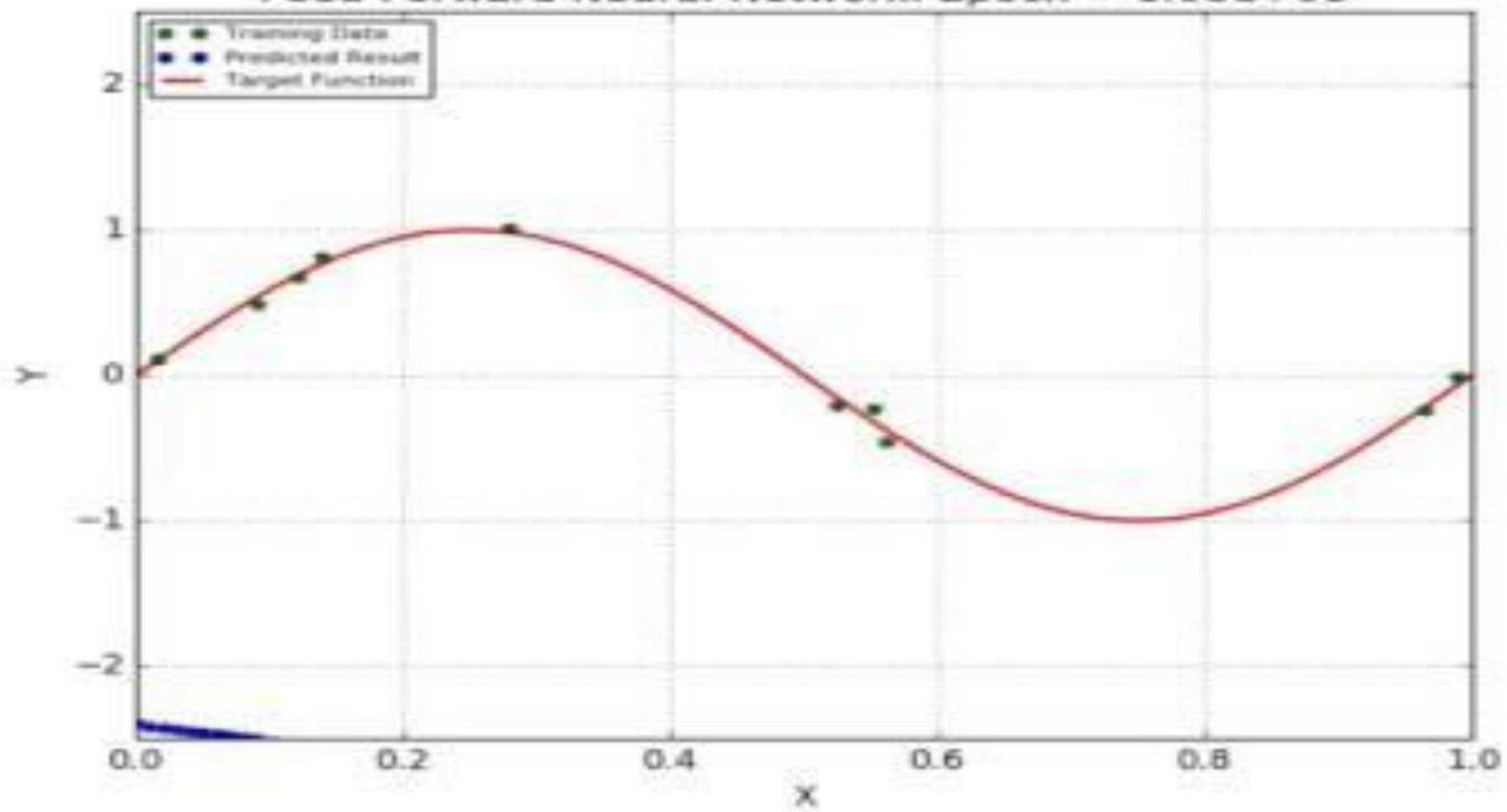


Neural Network Architecture



Epoch = 1 step of SGD throughout these examples.

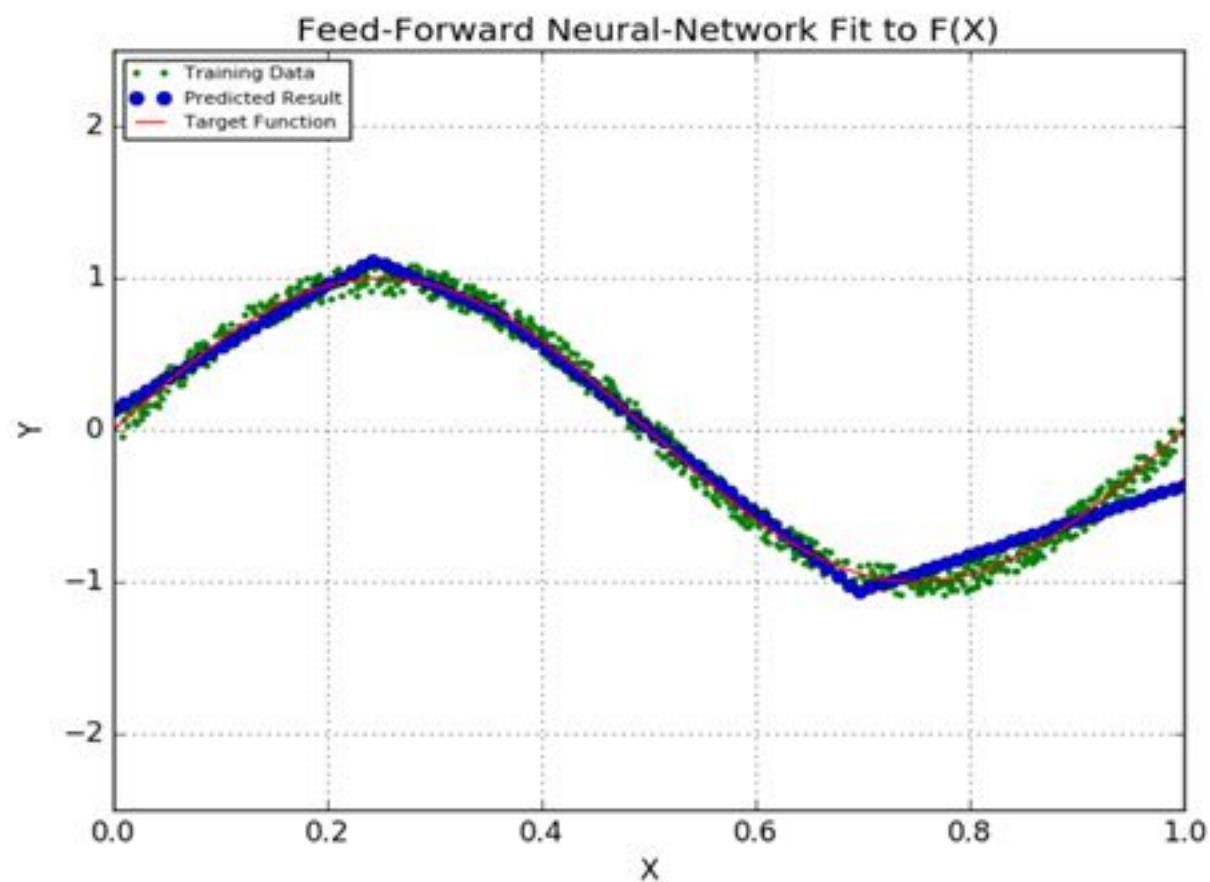
Feed-Forward Neural-Network: Epoch = 0.00e+00



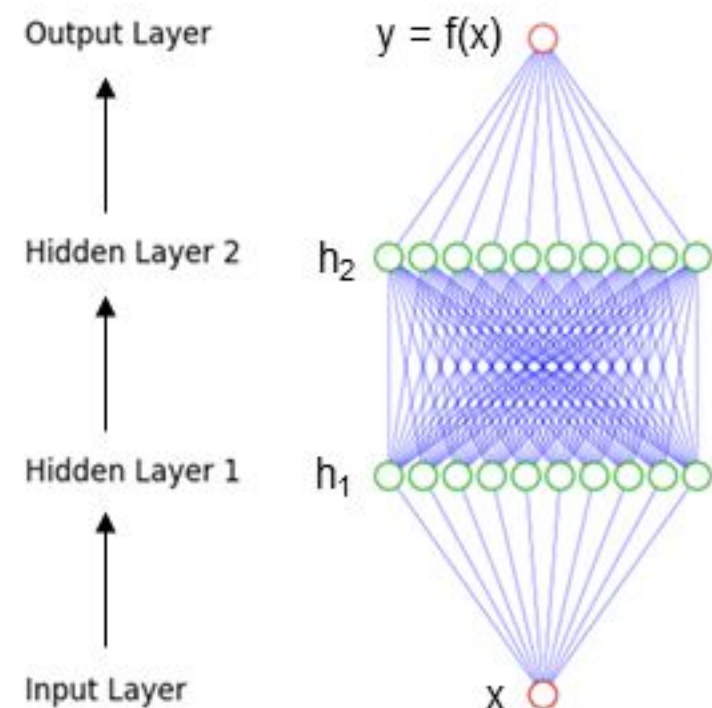
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 10 and learning rate 10^{-4} .



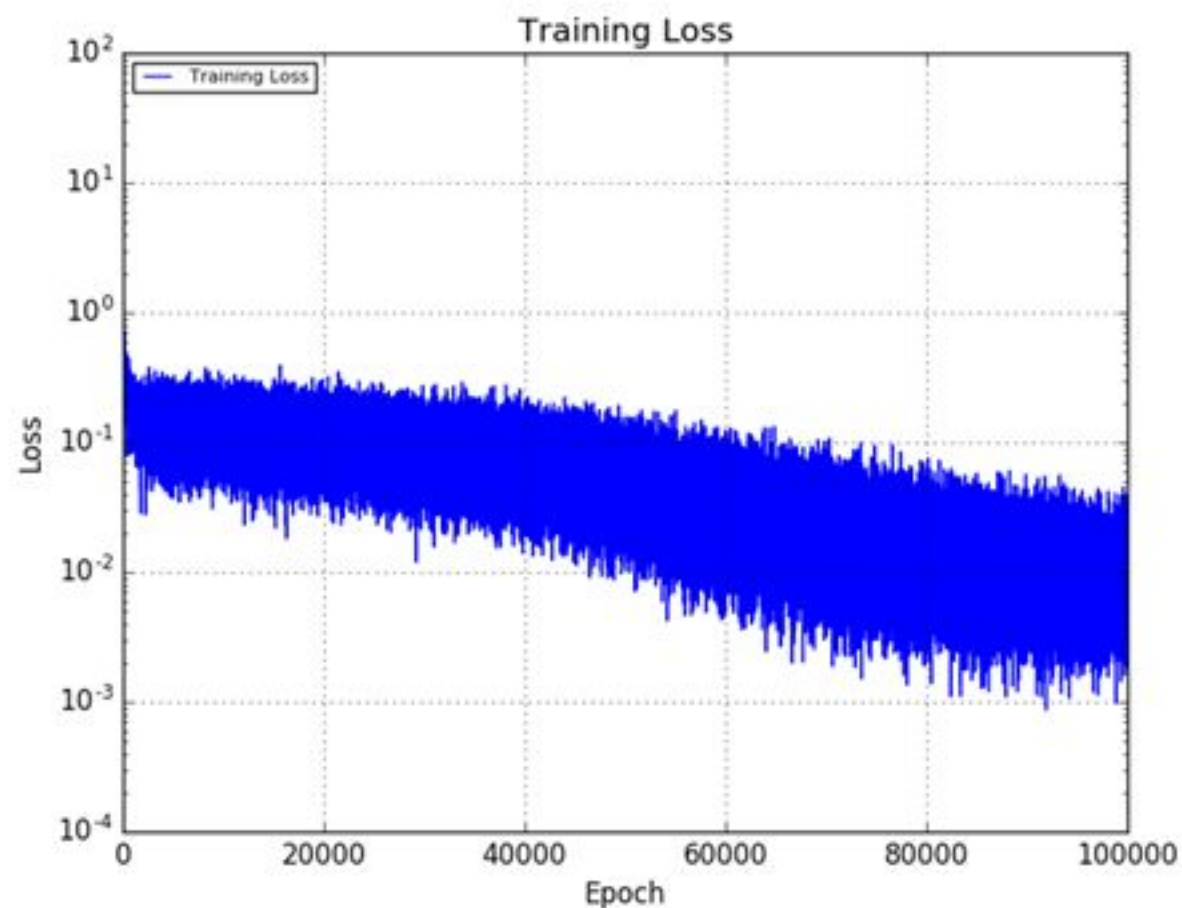
Neural Network Architecture



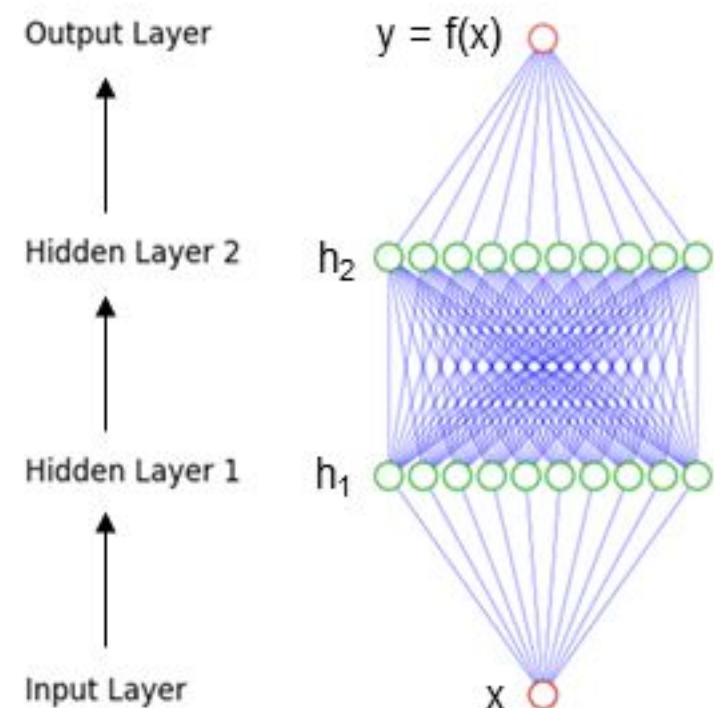
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 10 and learning rate 10^{-4} .



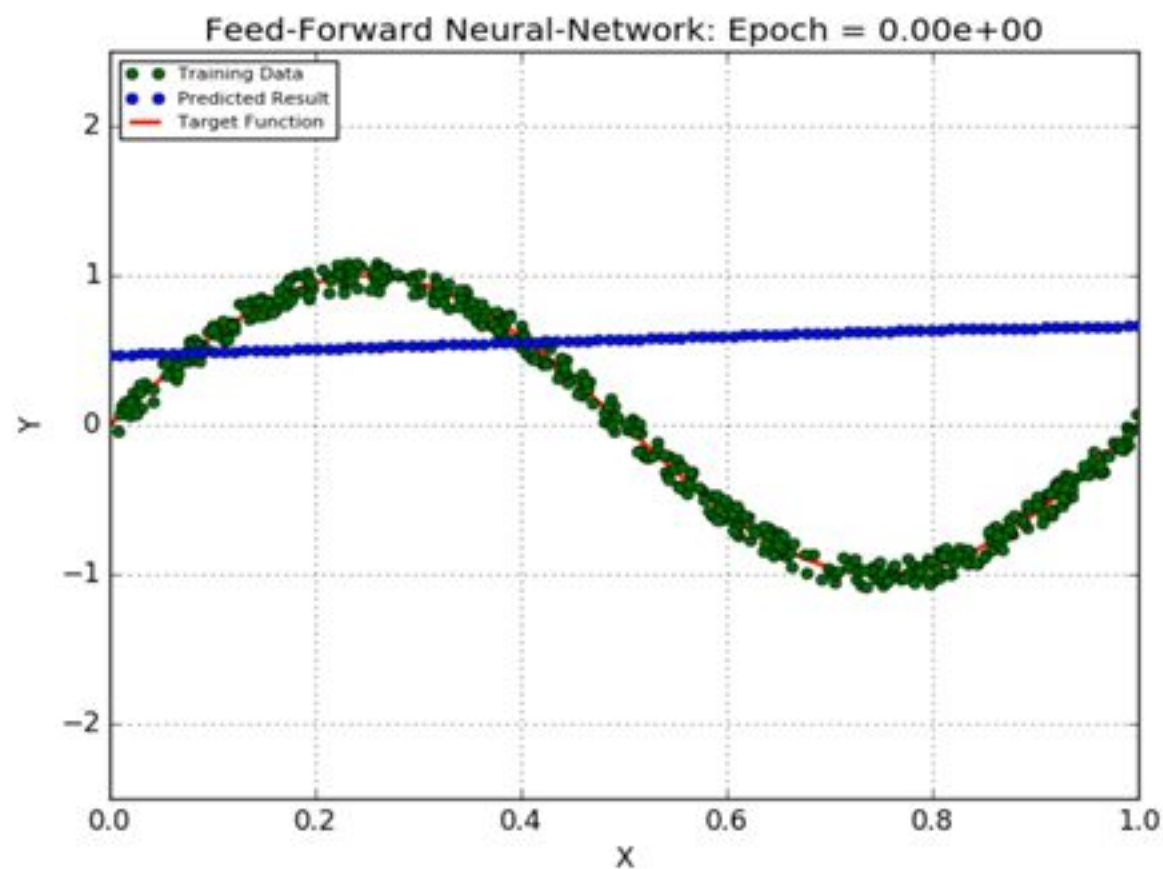
Neural Network Architecture



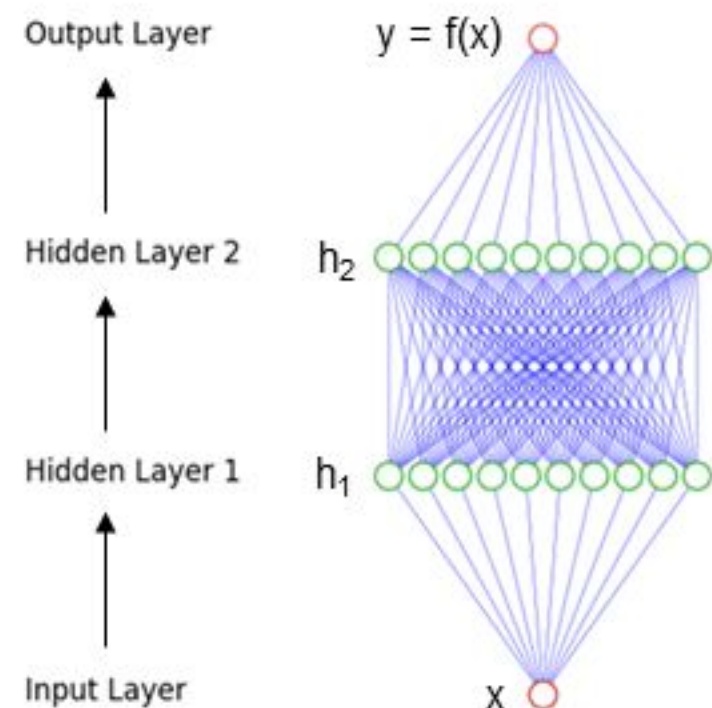
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

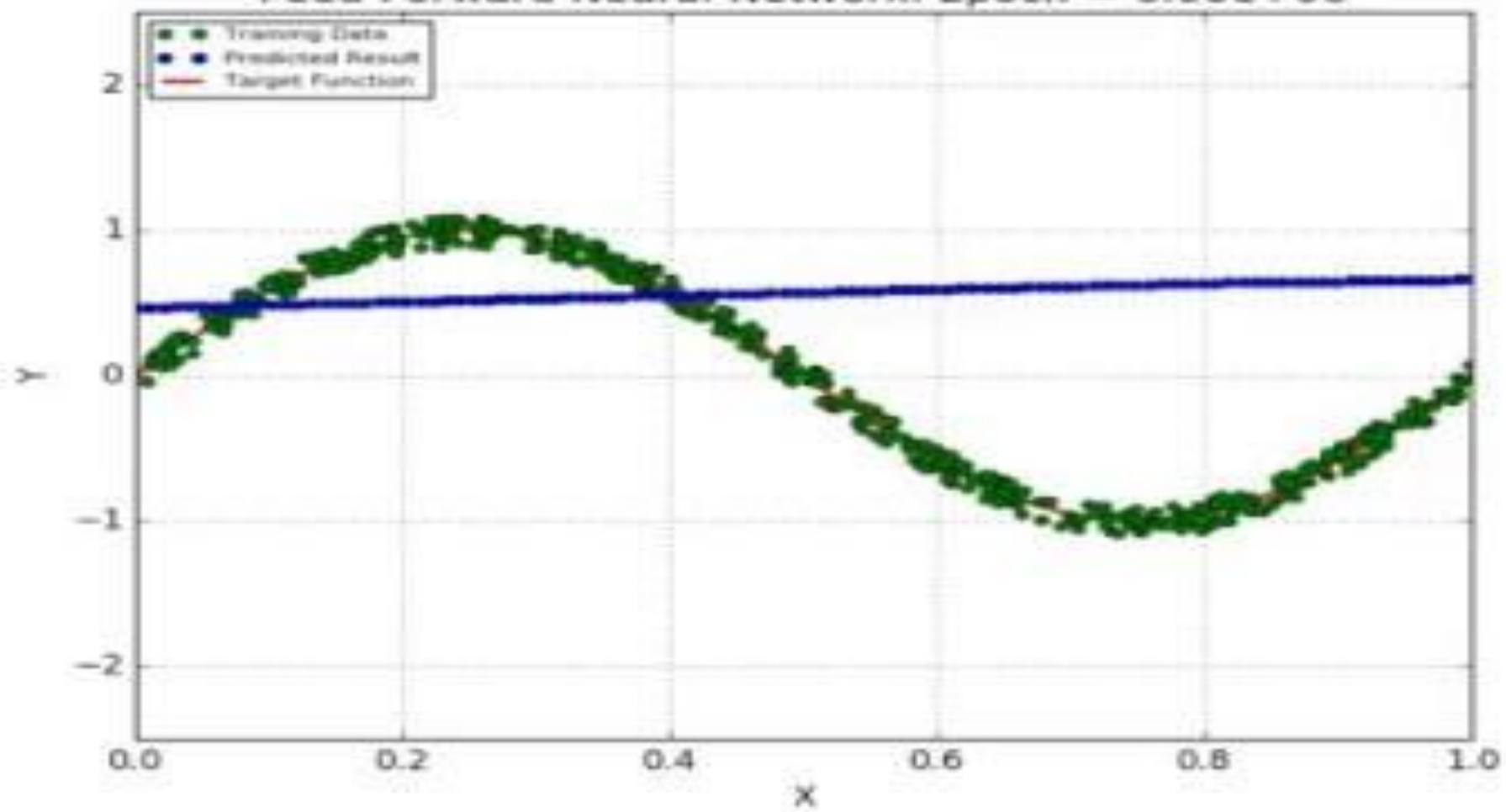
Stochastic Gradient Descent (SGD) used with batch size 500 and learning rate 10^{-4} .



Neural Network Architecture



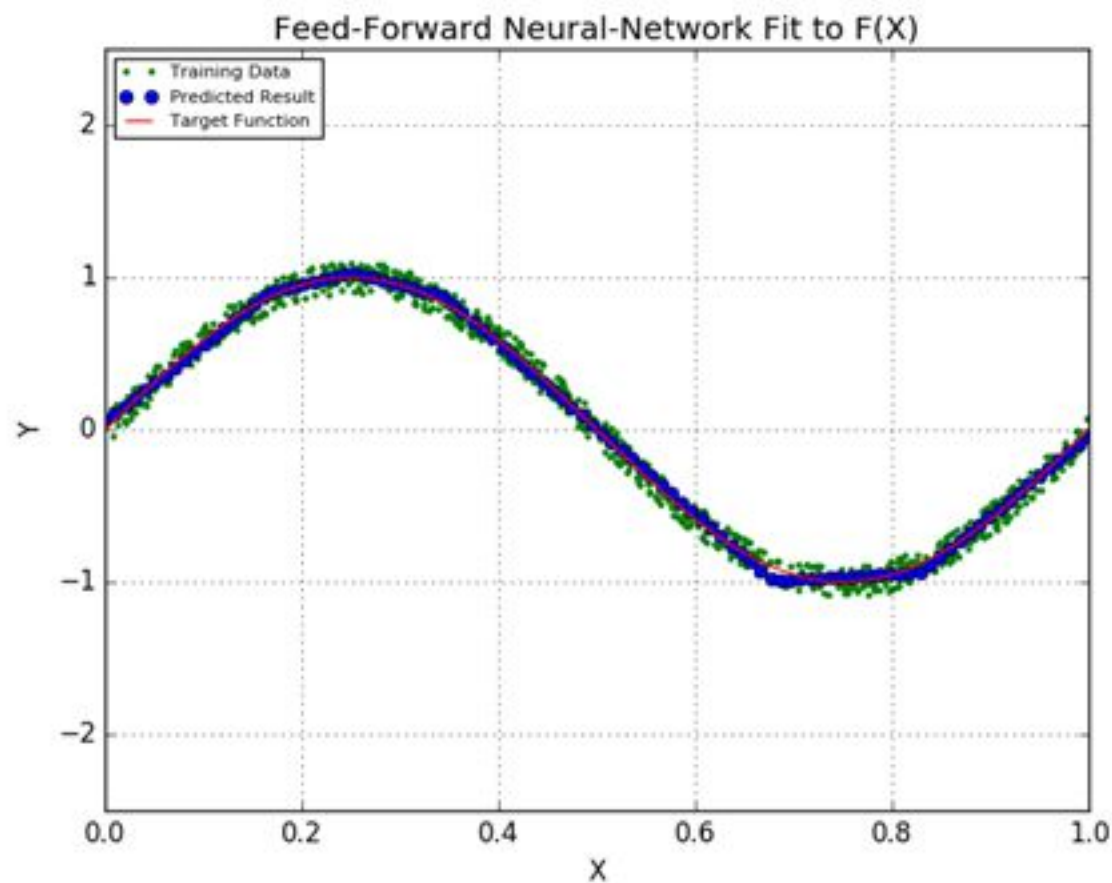
Feed-Forward Neural-Network: Epoch = 0.00e+00



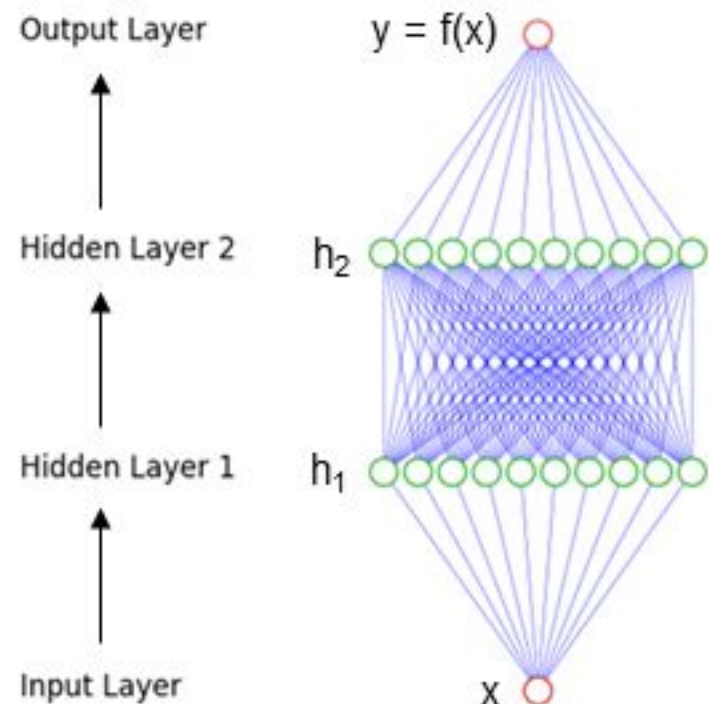
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 500 and learning rate 10^{-4} .



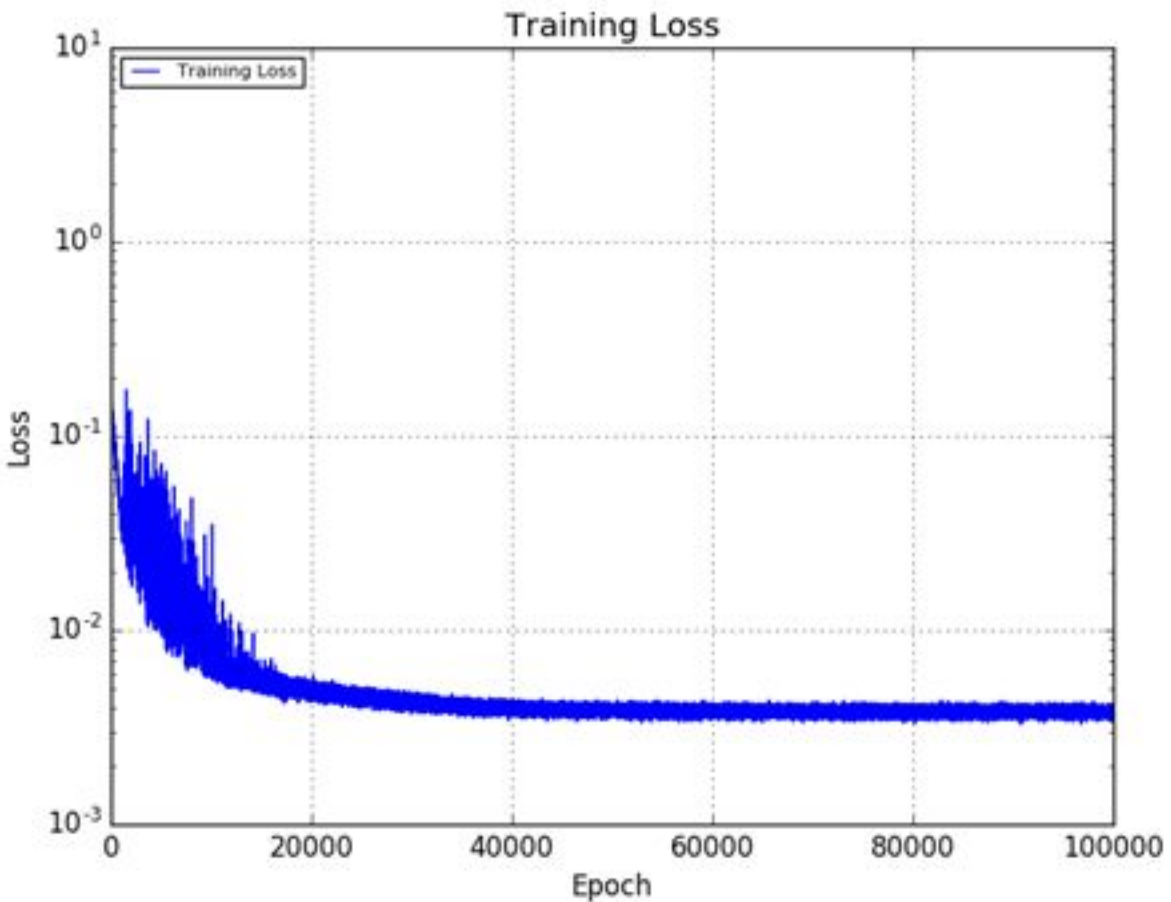
Neural Network Architecture



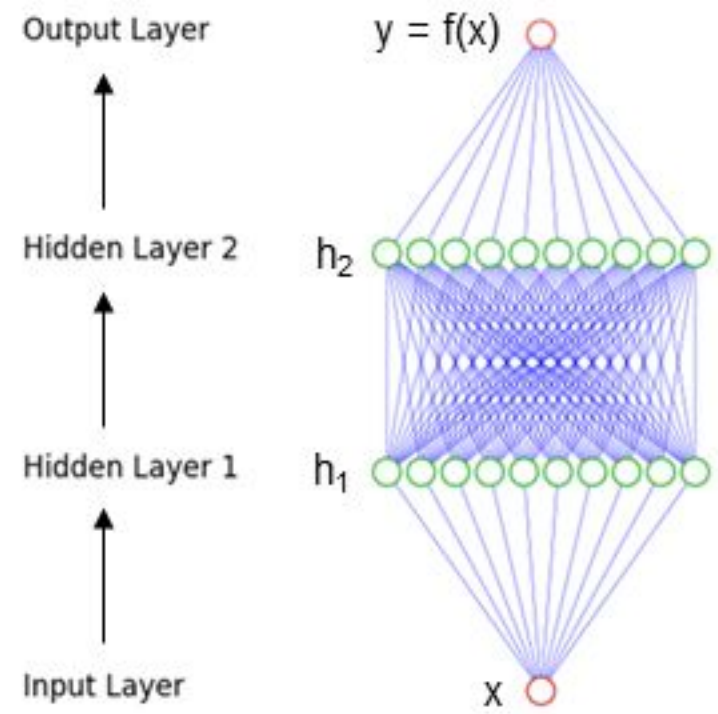
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(x_i) + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 500 and learning rate 10^{-4} .



Neural Network Architecture



Additional fine-tuning of hyper-parameters should be done to enhance efficiency of training.

Feed-Forward Neural Networks: Example

Example: Approximate the function $y = \sin(6\pi x) + 2x^2$ using FFNN.

Architecture: 2-layers with 100-hidden ReLu nodes per layer.

This NN architecture spans piecewise linear functions (100 nodes).

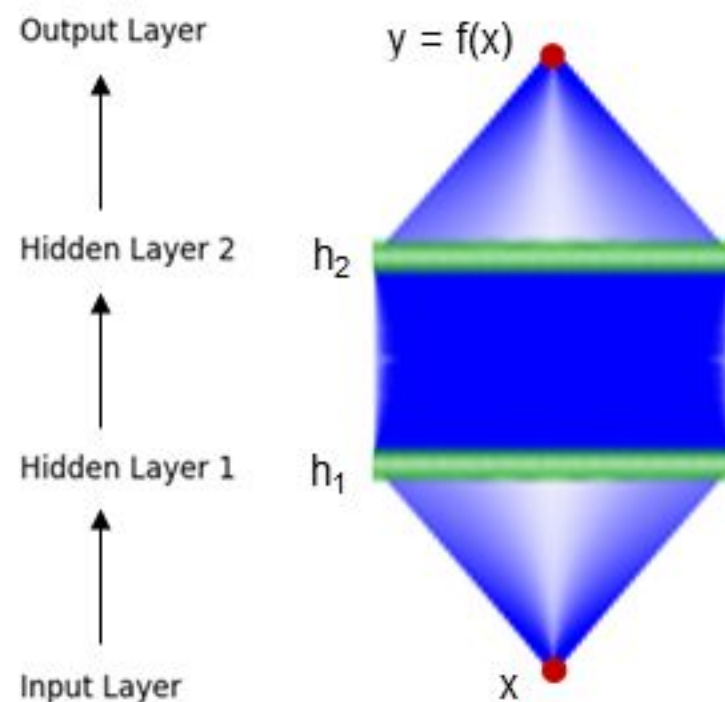
Explicitly: $f(X) = g(g(g(X \cdot W^{(1)} + b^{(1)}) \cdot W^{(2)} + b^{(2)}) \cdot W^{(3)} + b^{(3)})$
with $g(z) = \max(0, z)$.

A notion of "loss" required to assess level of success in fit.

Least-squares loss function $\ell(\{x_i, y_i\}) = \sum_i (f(x_i) - y_i)^2$.

Learning W, b proceeds by stochastic gradient descent.

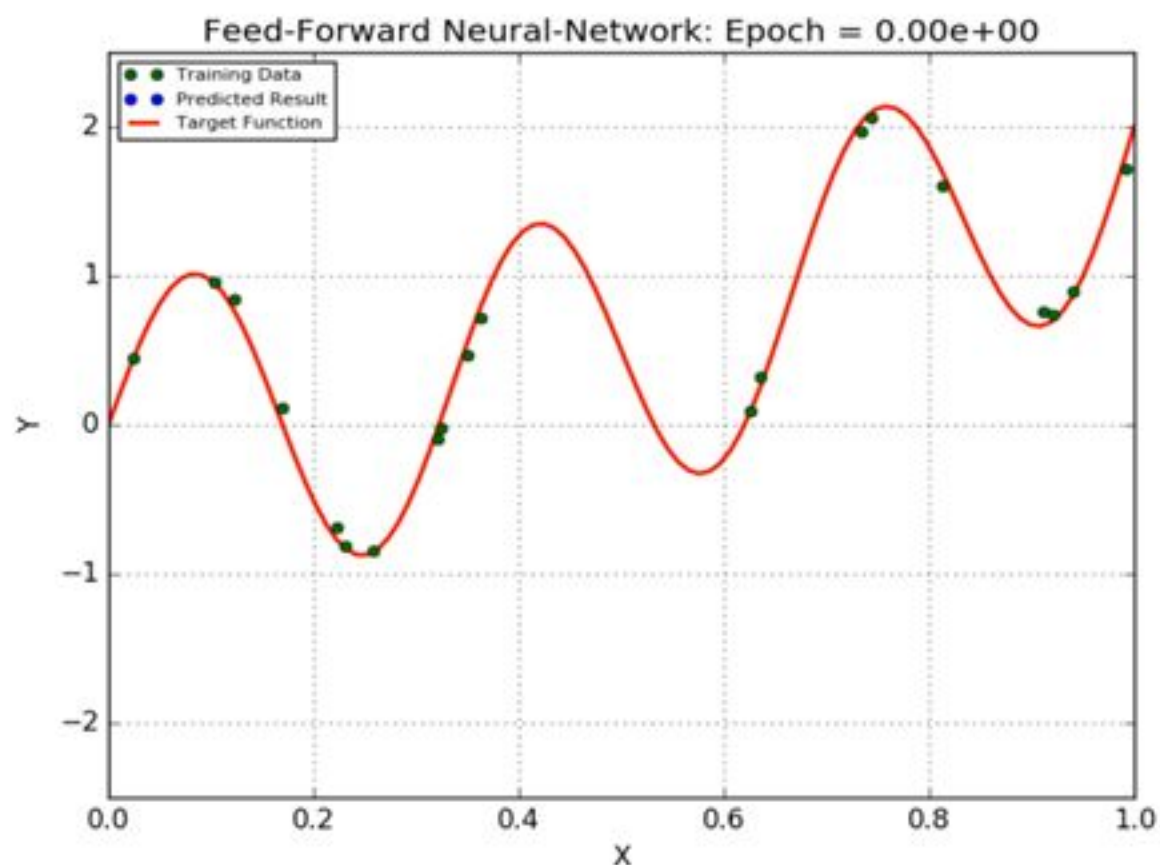
Neural Network Architecture



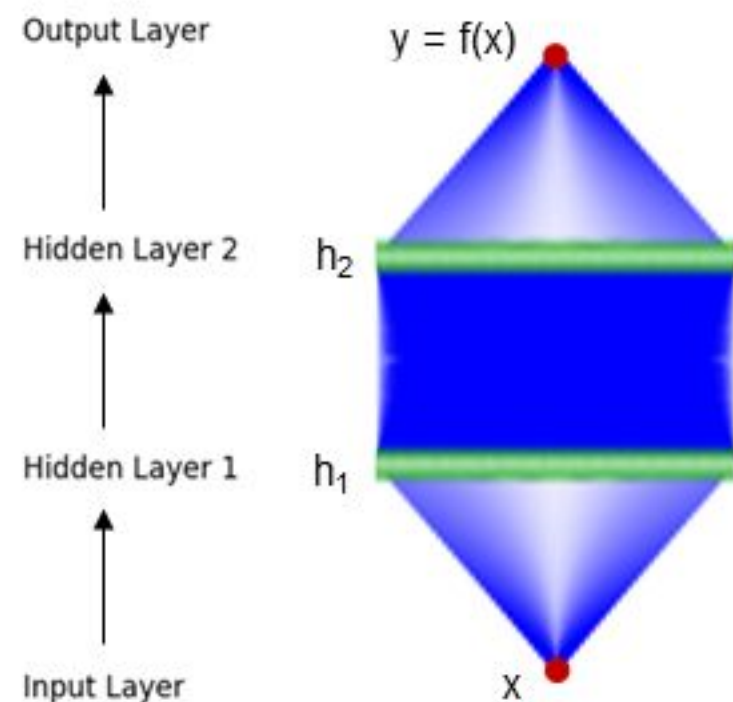
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(6\pi x_i) + 2x^2 + \xi_i$.

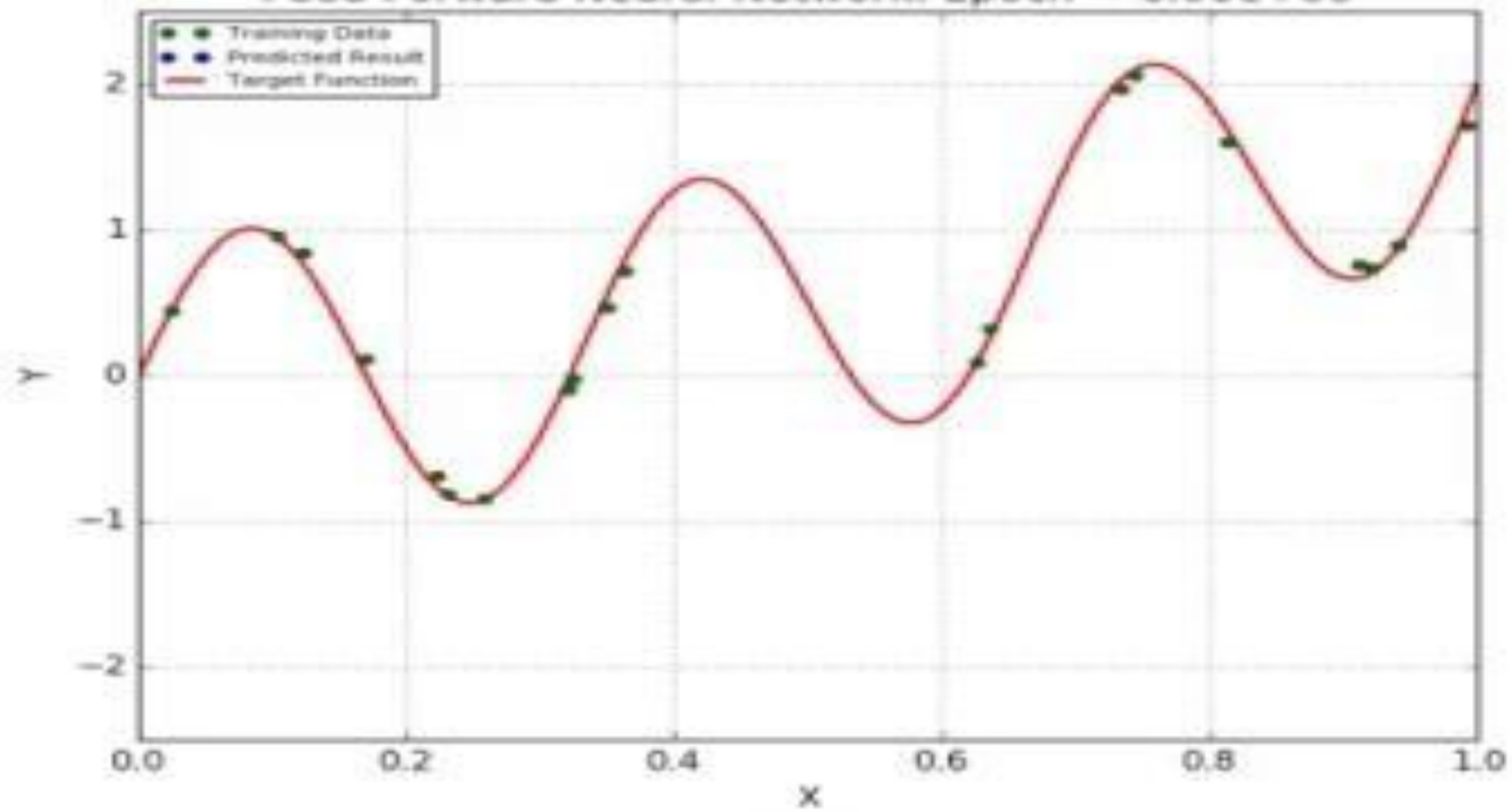
Stochastic Gradient Descent (SGD) used with batch size 20 and learning rate 10^{-4} .



Neural Network Architecture



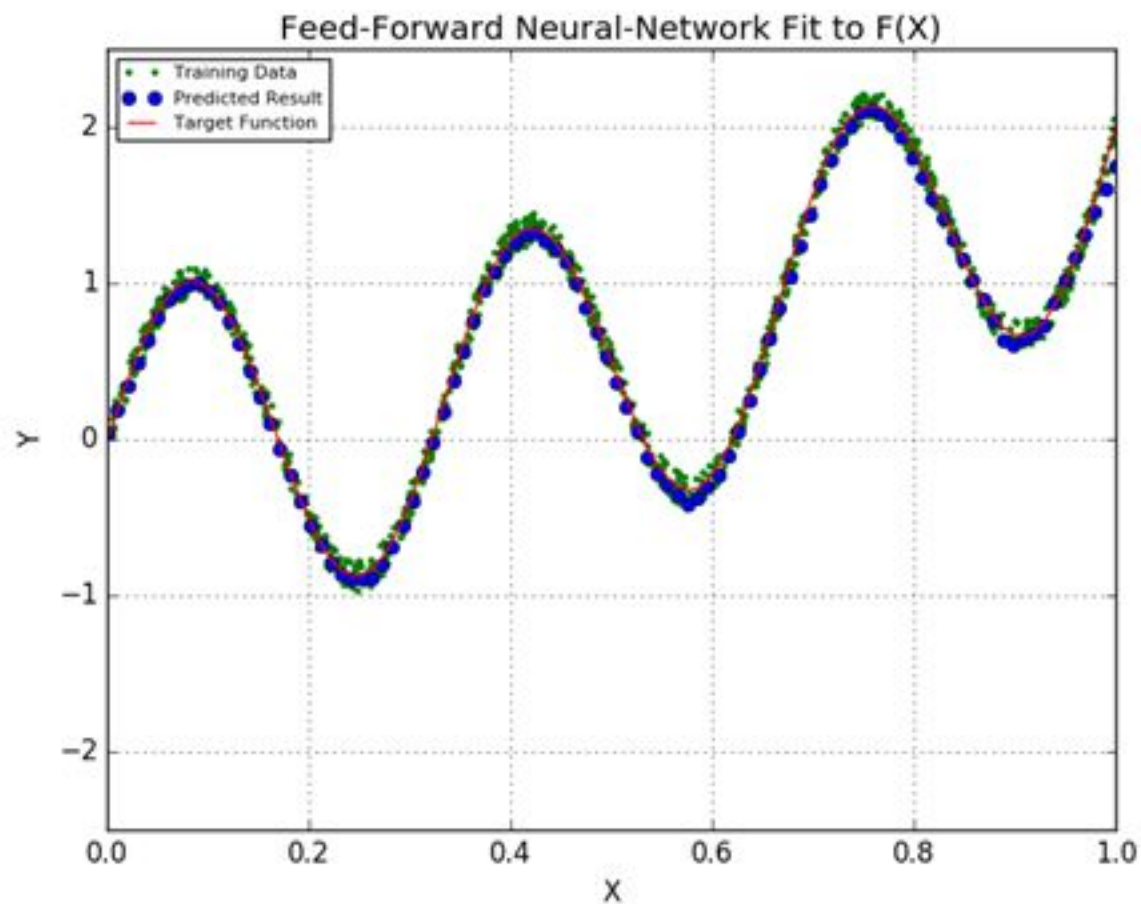
Feed-Forward Neural-Network: Epoch = 0.00e+00



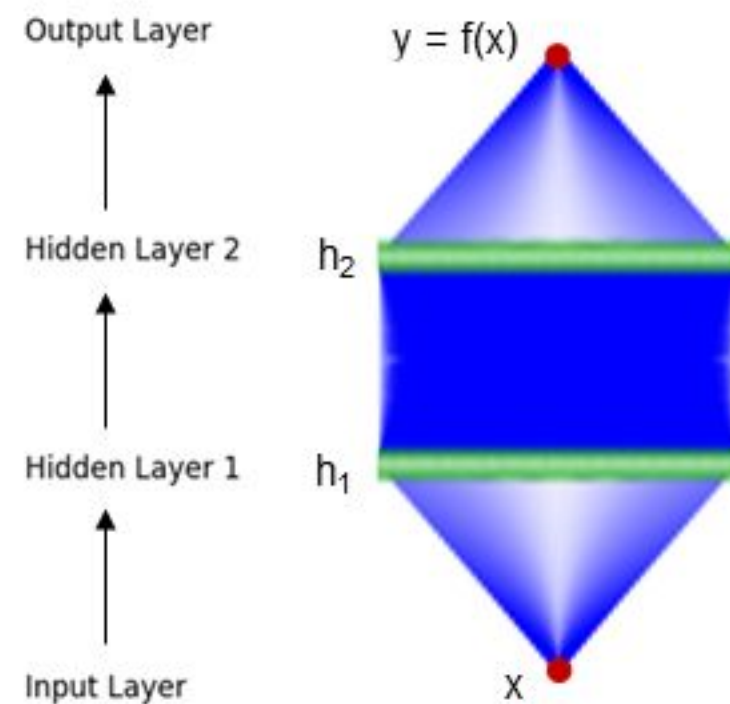
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(6\pi x_i) + 2x^2 + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 20 and learning rate 10^{-4} .



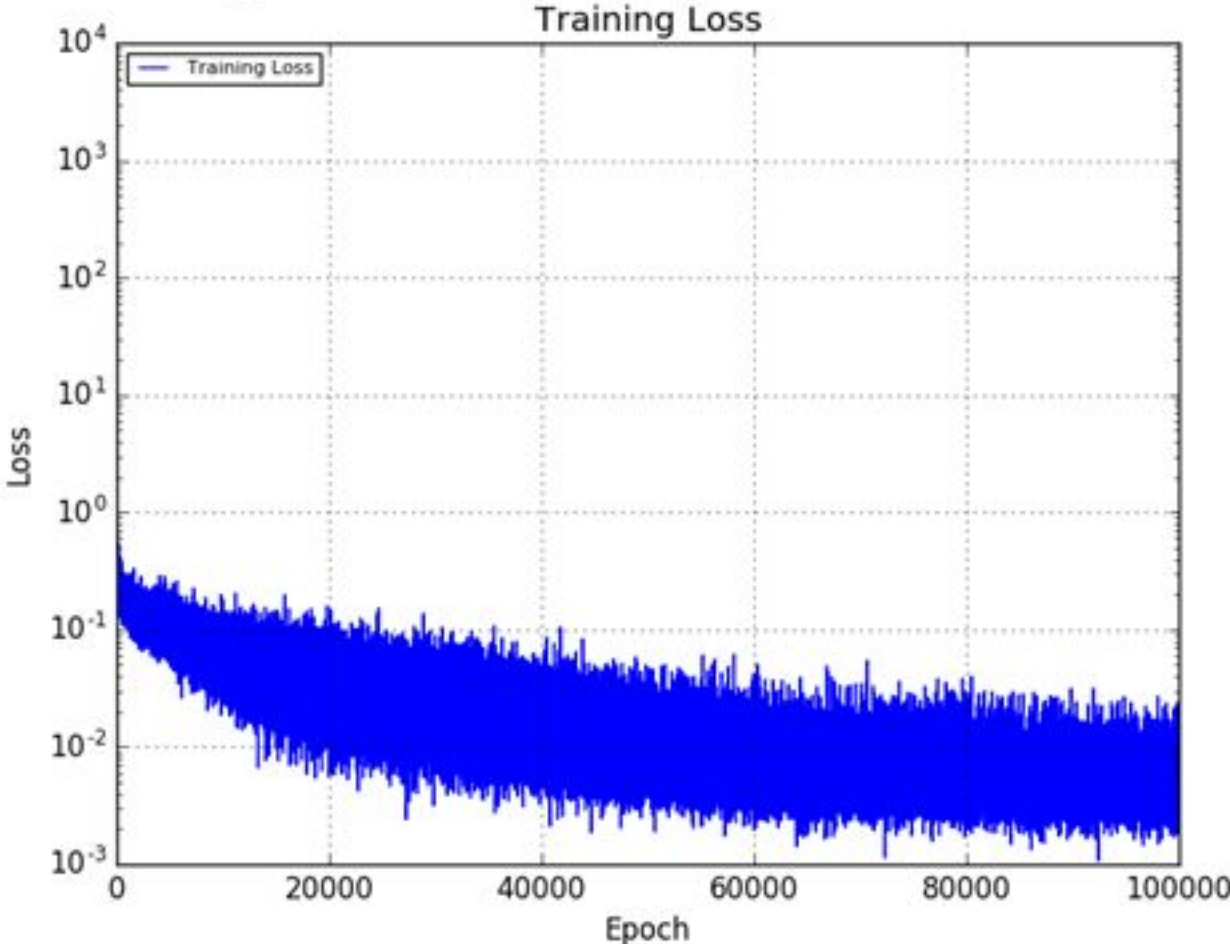
Neural Network Architecture



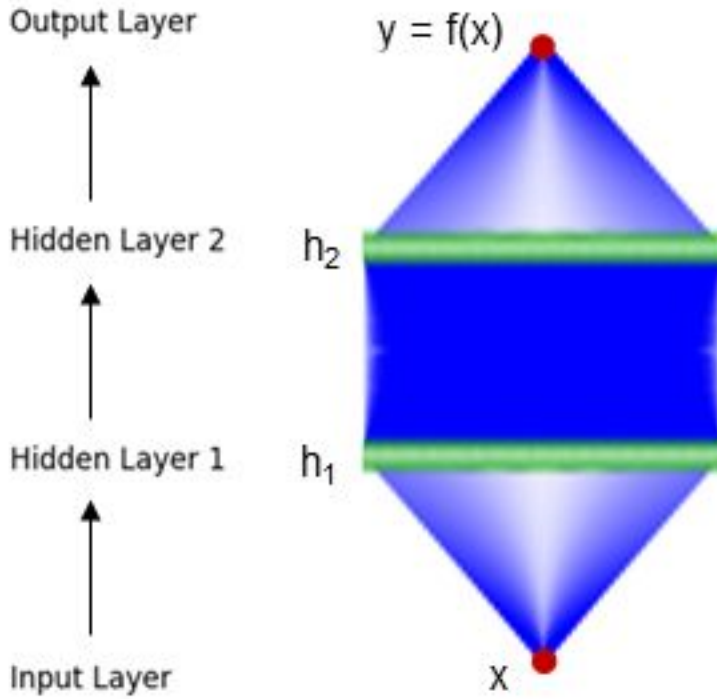
Feed-Forward Neural Networks: Example

Trained FFNN on set of 1000 samples $y_i = \sin(6\pi x_i) + 2x^2 + \xi_i$.

Stochastic Gradient Descent (SGD) used with batch size 20 and learning rate 10^{-4} .

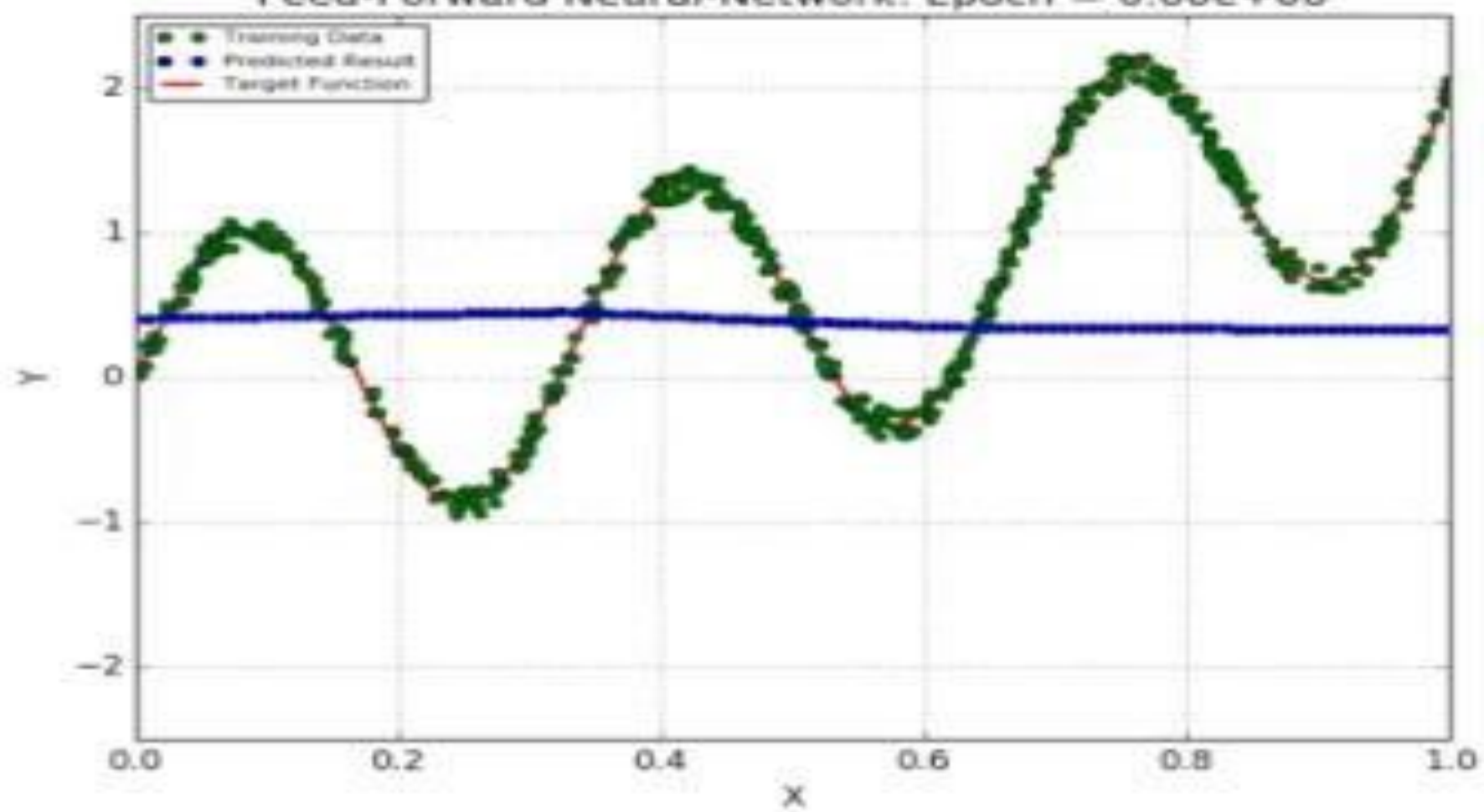


Neural Network Architecture



Additional fine-tuning of hyper-parameters should be done to enhance efficiency and robustness of training.

Feed-Forward Neural-Network: Epoch = 0.00e+00



Summary

Neural Networks: Summary

Deep Neural Network (DNN)

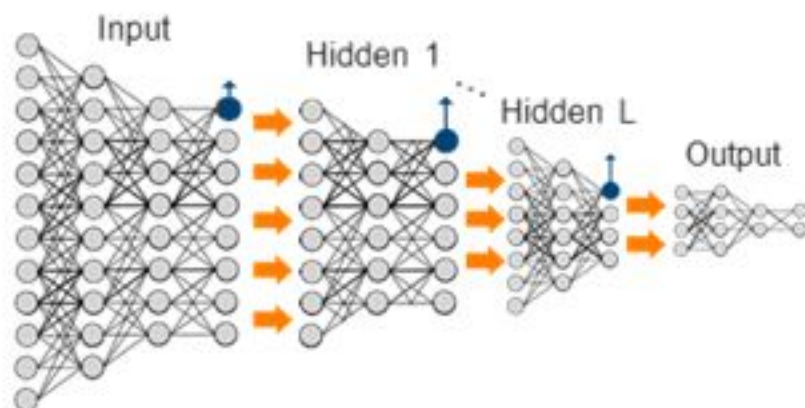
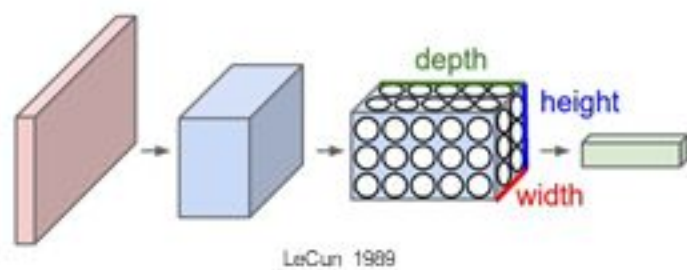
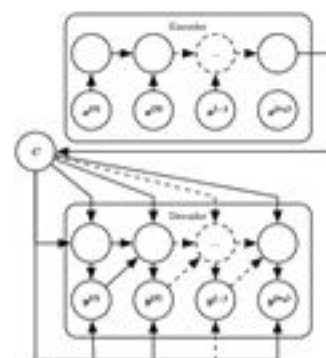


Image Classifier: Convolutional Neural Network (CNN)



Natural Language Processing
Recurrent Neural Network (RNN)



Alpha-Go



Google's Self-Driving Car



- **Neural Networks** are providing state-of-the-art results in many fields: (computer vision, natural language processing, reinforcement learning).

Neural Networks: Summary

Deep Neural Network (DNN)

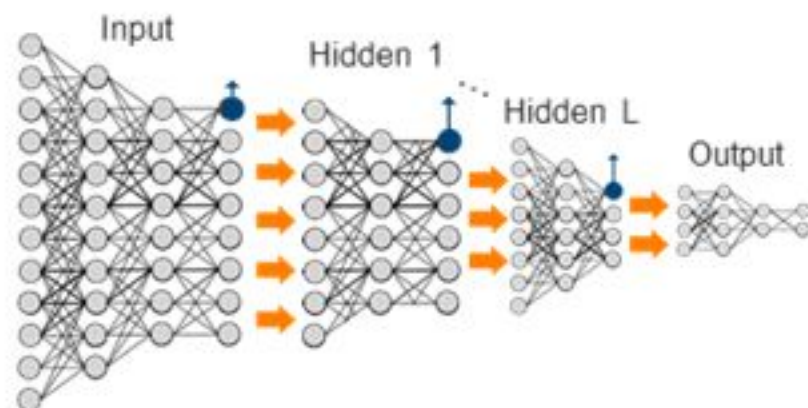
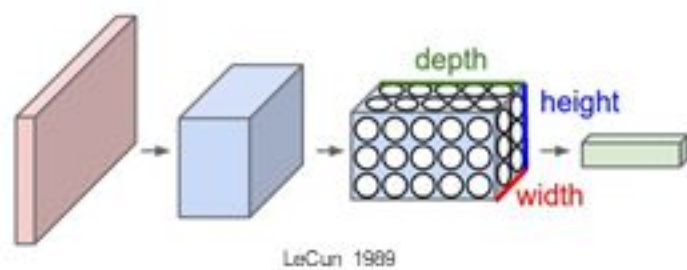
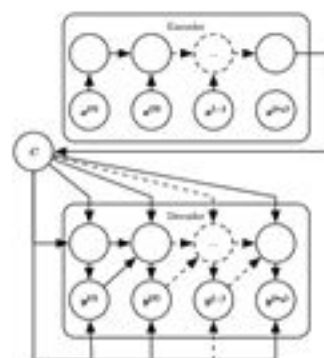


Image Classifier: Convolutional Neural Network (CNN)



Natural Language Processing
Recurrent Neural Network (RNN)



Alpha-Go



Google's Self-Driving Car



- **Neural Networks** are providing state-of-the-art results in many fields: (computer vision, natural language processing, reinforcement learning).
- **Powerful approximation properties:** target functions approximated by compositions, perform well in high dimensional spaces, many variants.

Neural Networks: Summary

Deep Neural Network (DNN)

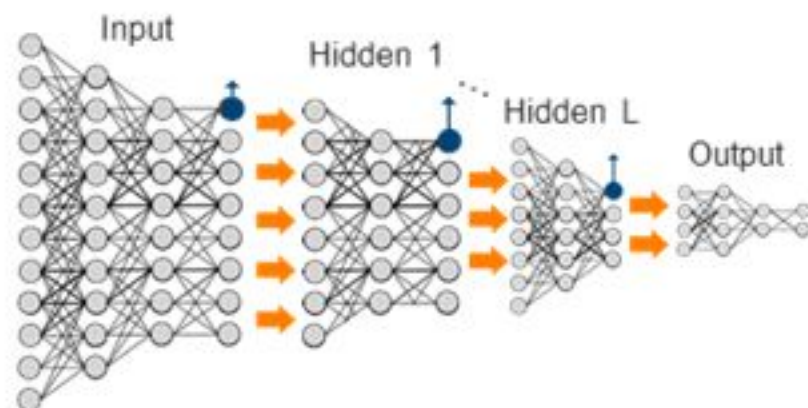
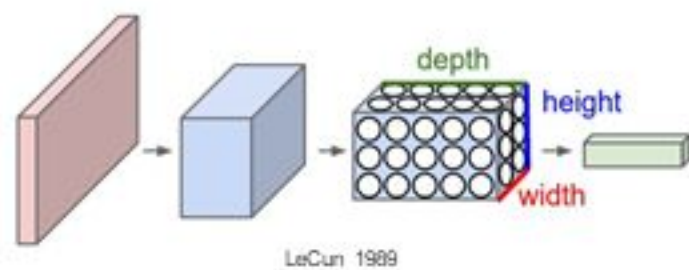
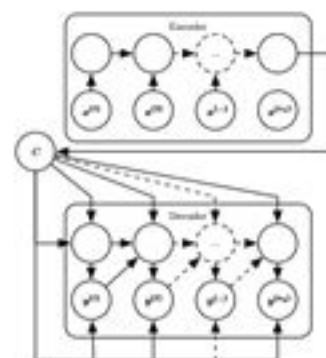


Image Classifier: Convolutional Neural Network (CNN)



Natural Language Processing
Recurrent Neural Network (RNN)



Alpha-Go



Google's Self-Driving Car



- **Neural Networks** are providing state-of-the-art results in many fields: (computer vision, natural language processing, reinforcement learning).
- **Powerful approximation properties:** target functions approximated by compositions, perform well in high dimensional spaces, many variants.
- **With appropriate learning protocols,** despite richness of NN's, seems they can be well-enough regularized to not overfit the training data.

Neural Networks: Summary

Deep Neural Network (DNN)

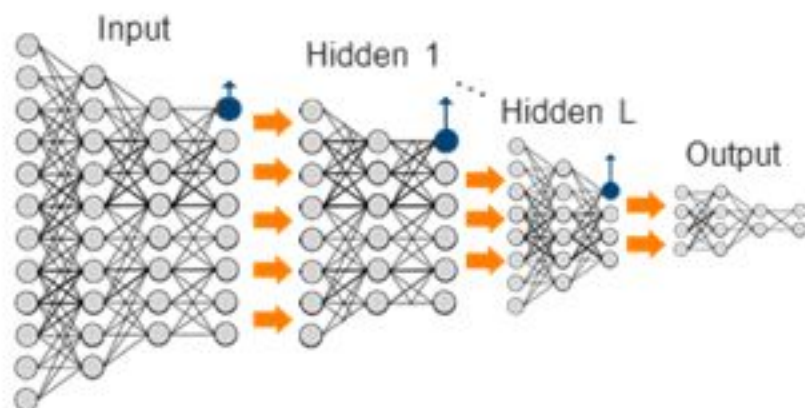
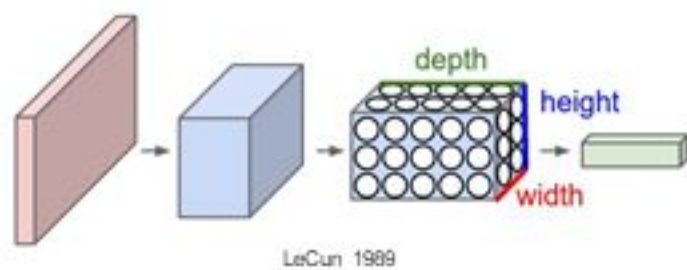
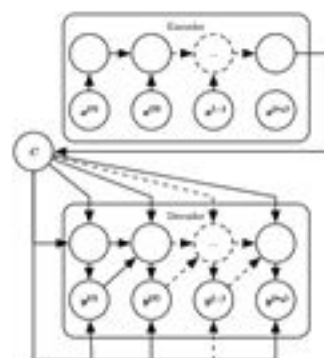


Image Classifier: Convolutional Neural Network (CNN)



Natural Language Processing
Recurrent Neural Network (RNN)



Alpha-Go



Google's Self-Driving Car



- **Neural Networks** are providing state-of-the-art results in many fields: (computer vision, natural language processing, reinforcement learning).
- **Powerful approximation properties:** target functions approximated by compositions, perform well in high dimensional spaces, many variants.
- **With appropriate learning protocols,** despite richness of NN's, seems they can be well-enough regularized to not overfit the training data.
- **Current research to better understand NN's:** choice of architectures, training protocols, approximation properties, reliability, interpretability, ...

