

---

# Deep Learning Foundations

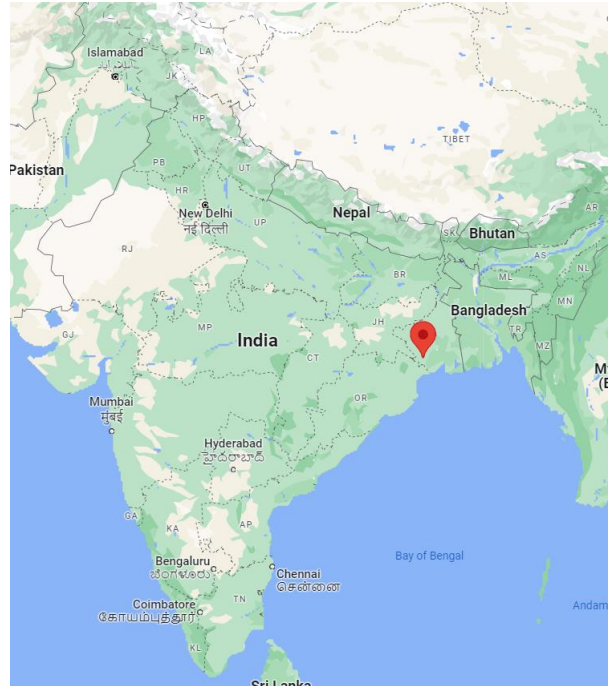
Winter term 2024/25

Sandipan Sikdar



# Who is giving the course?

---



- Assistant Professor at the Faculty of Electrical Engineering and Computer Science (Nov'22 - )
- Postdoctoral Researcher at RWTH Aachen (2018 – 2022)
- PhD, Computer Science and Engineering from Indian Institute of Technology Kharagpur

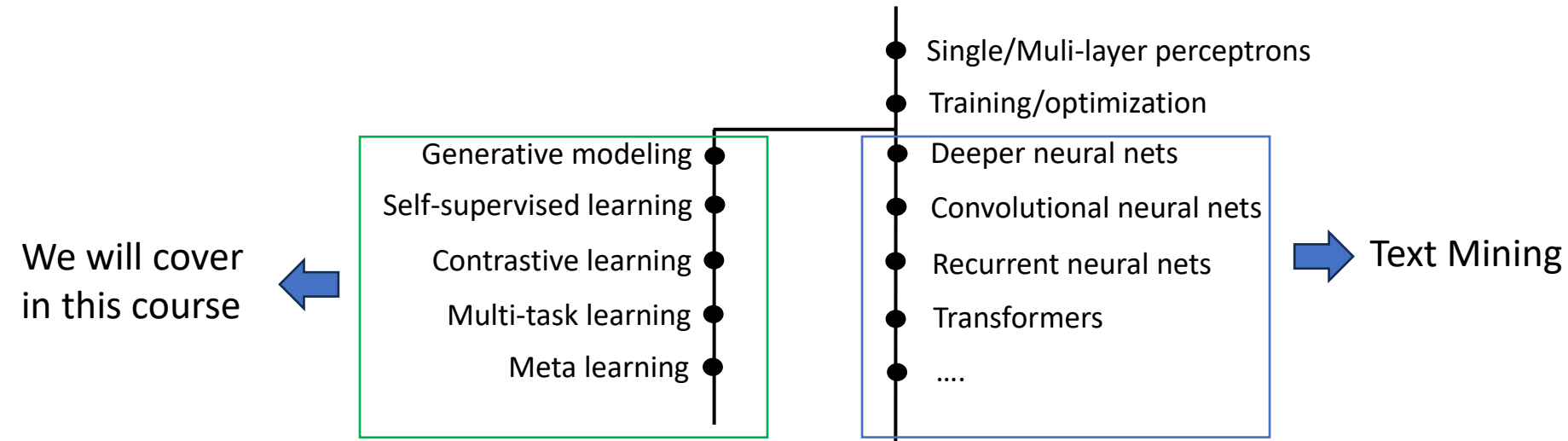
# Deep Learning Revolution

---



# Deep Learning

---



# Setup

---

- Lectures: WED 8:00 – 9:30 (2 \* 45 minutes)
- Exercises: WED 15:00 – 16:30 (2 \* 45 minutes)
- Lecture
  - Mostly theory, algorithms, findings
- Exercise
  - Presentation on how to do in practice
  - Hands on implementations (Python/Pytorch)
  - Finish exercises at home (own responsibility)

# Setup

---

- Home assignments
  - 3 home assignments (3 weeks time)
  - Teams of upto 3 students
  - No plagiarism! (Can result in exclusion from the course)
  - Tentative dates:
    - Home assignment 1 (13.11.2024)
    - Home assignment 2 (04.12.2024)
    - Home assignment 3 (08.01.2025)

# Grading

---

- Final exam:
  - Pen and paper exam
  - The final exam will be scaled to 100 points. The minimum passing score is 50 points.
  - Bonus points for home assignments: Points from the home assignments will be scaled to 10 points and added to the final score.

# Course Plan

Week	Date	Theme	Topics
1	16.10	Intro + Preliminaries	Neural network basics
2	23.10	No lectures	
3	30.10	Neural Networks	Training neural networks
4	06.11	Deeper architectures	CNNs/RNNs/Transformers
5	13.11	Generative models I	Preliminaries, latent variable models
6	20.11	Generative models II	Variational Auto Encoder (VAE)
7	27.11	Generative models III	Diffusion models
8	04.12	Generative models IV	Flow Models
9	11.12	Generative models V	GANs
10	18.12	Learning paradigm	Self-supervised and contrastive learning
11	08.01	Learning paradigm	Multi-task and Meta learning
12	15.01	Adversarial Robustness	Poisoning attacks
13	22.01	Adversarial Robustness	Attacks and defenses
14	29.01	Summary and outlook	



---

# Preliminaries

# Supervised Learning

---

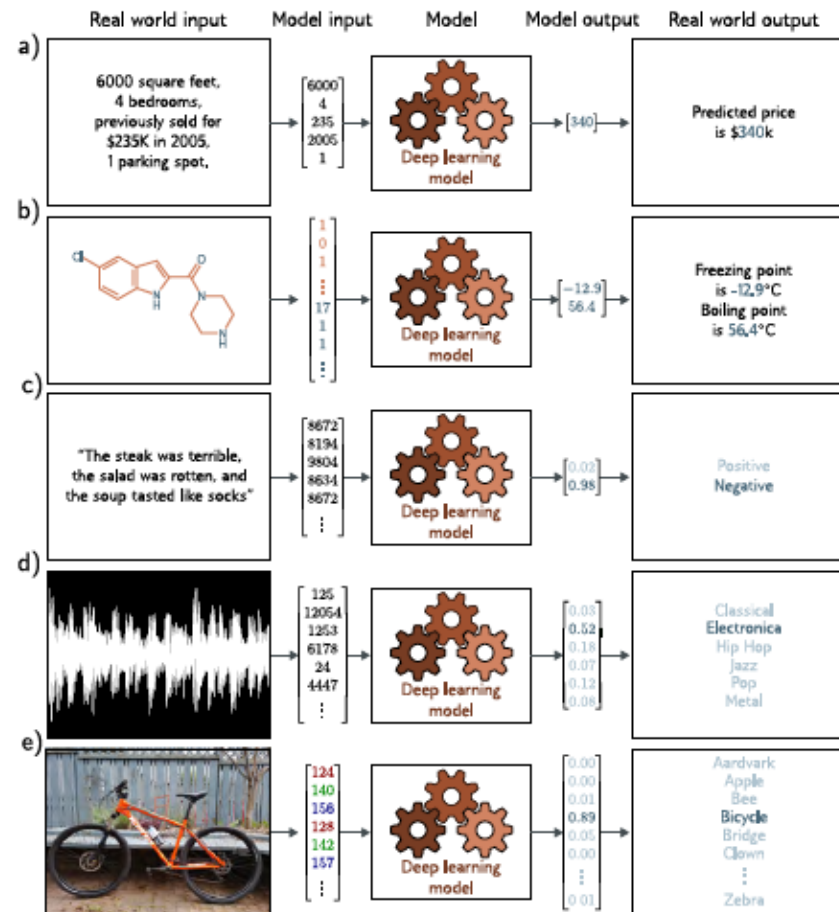
- We are given a set of training examples of the form  $\{(x^1, y^1), \dots, (x^n, y^n)\}$
- $x^i$  represents the feature vector for the  $i$ th training example
- $y^i$  represents the ground truth
  - Classification:  $y^i \in [c]$
  - Regression:  $y^i \in R$

# Supervised Learning

---

- We are given a set of training examples of the form  $\{(x^1, y^1), \dots, (x^n, y^n)\}$
- $x^i$  represents the feature vector for the  $i$ th training example
- $y^i$  represents the ground truth
- Goal: Find  $f_W: R^d \rightarrow R, f_W(x^i) \approx y^i$
- $W$ : model parameters

# Supervised Learning



# Supervised Learning

---

- Goal: Find  $f_W: R^d \rightarrow R, f_W(x^i) \approx y^i$
- How do we find  $W$ ?

$$\min_W \frac{1}{n} \sum_{i=1}^n l(f_W(x^i), y^i)$$

- $l \rightarrow$  loss/cost function
- Empirical risk minimization

# Loss Function

---

- A function that measures, for each value of the  $W$ s how close  $f_W(x^i)$ 's are to the corresponding  $y^i$ s
- Regression:

$$l(f_W(x^i), y^i) = \frac{1}{2} (f_W(x^i) - y^i)^2$$

# Loss Function

---

- A function that measures, for each value of the  $W$ s how close  $f_W(x^i)$ 's are to the corresponding  $y^i$ s
- Classification:

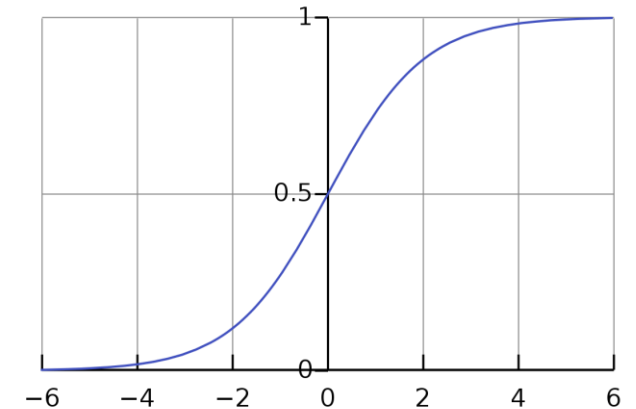
$$y^i = \begin{cases} 1 \\ 0 \end{cases}$$

# Loss Function

---

- Convert the output of the model to probabilities

- Sigmoid function:  $g(x) = \frac{1}{1+e^{(-x)}}$



$$\begin{cases} P(y^i = 1 | x^i, W) = g(f_W(x^i)) \\ P(y^i = 0 | x^i, W) = 1 - g(f_W(x^i)) \end{cases}$$



# Loss Function

---

$$\begin{cases} P(y^i = 1 | x^i, W) = g(f_W(x^i)) \\ P(y^i = 0 | x^i, W) = 1 - g(f_W(x^i)) \end{cases}$$



$$P(y^i | x^i, W) = g(f_W(x_i))^{y^i} \left( 1 - g(f_W(x^i)) \right)^{1-y^i}$$

# Loss Function

---

$$\max_W \prod_{i=1}^n P(y^i | x^i, W)$$

$$= \max_W \sum_{i=1}^n \log P(y^i | x^i, W)$$

$$= \min_W - \sum_{i=1}^n \log P(y^i | x^i, W)$$

# Loss Function

---

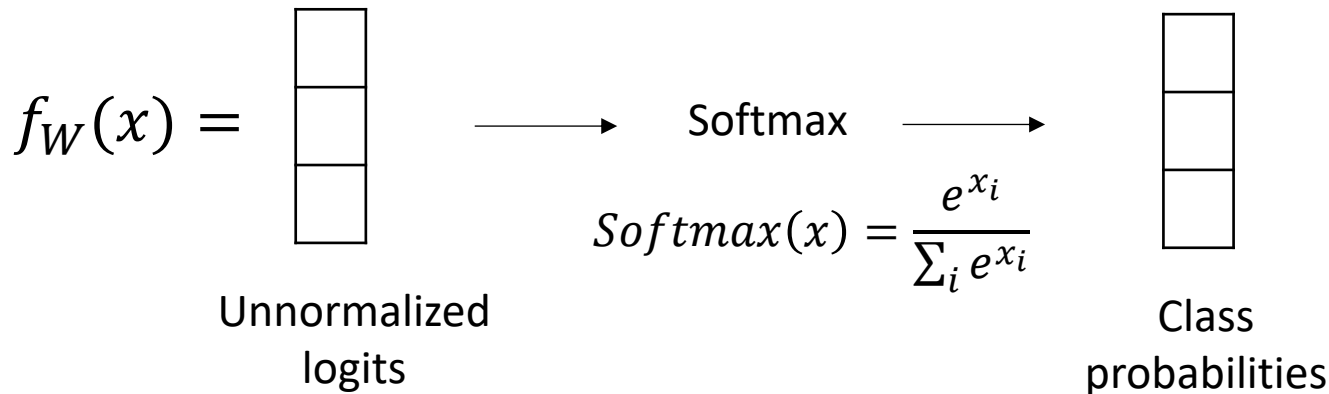
$$\min_W - \sum_{i=1}^n \log P(y^i | x^i, W)$$
$$= \min_W - \sum_{i=1}^n y^i \log(g(f_W(x))) + (1 - y^i) \log(1 - g(f_W(x^i)))$$

Cross entropy loss

# Loss Function

---

- Generalizing cross entropy to multiple class
- Model output -> vector of the dimension of the number of classes
- Use softmax instead of sigmoid



# Loss Function

---

- Class:  $c_j: j \in \{0,1,2\}$
- Each class represented by a 1-hot vector

$$c_j = \begin{cases} 0, [0,0,1] \\ 1, [0,1,0] \\ 2, [1,0,0] \end{cases}$$

- Cross-entropy loss for multiclass classification is then computed as -

$$-\sum_{j=0}^{|c|} c[j] \log(f_W(x^i)[j])$$

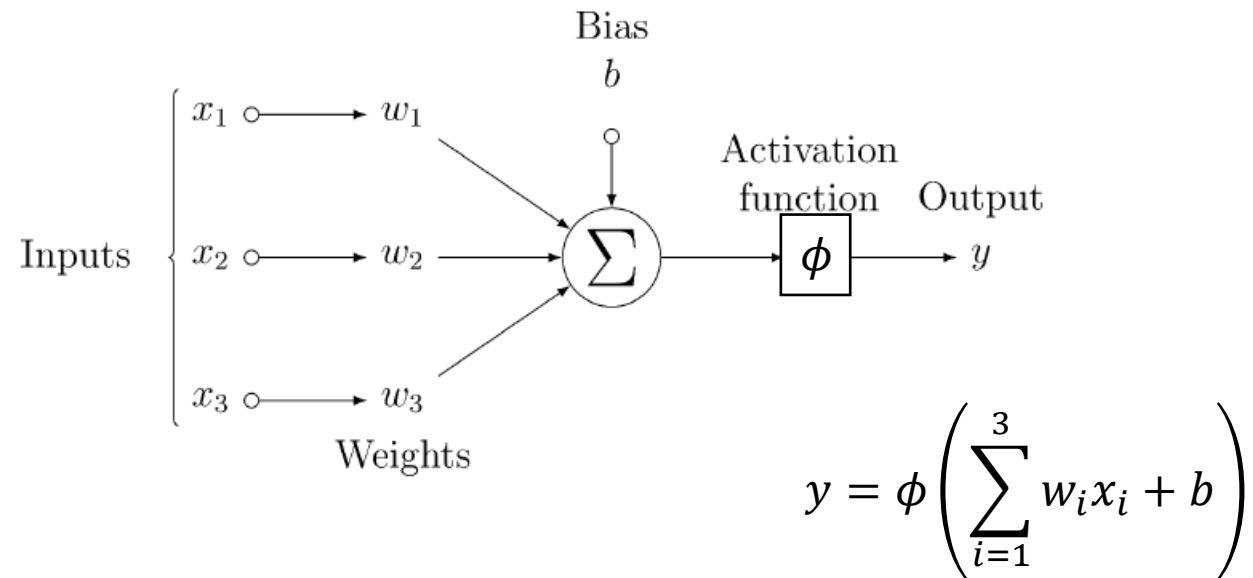
# Choice of model

---

- $f_W: R^d \rightarrow R, f_W(x^i) \approx y^i$
- What should  $f_W$  look like?
- Linear function
- $f_W(x) = W^T x + b$
- + Convex optimization
- - Limited representation

# Choice of model

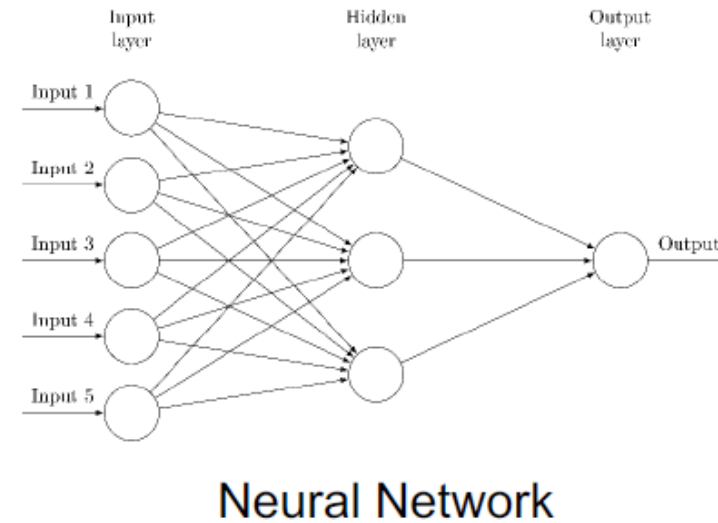
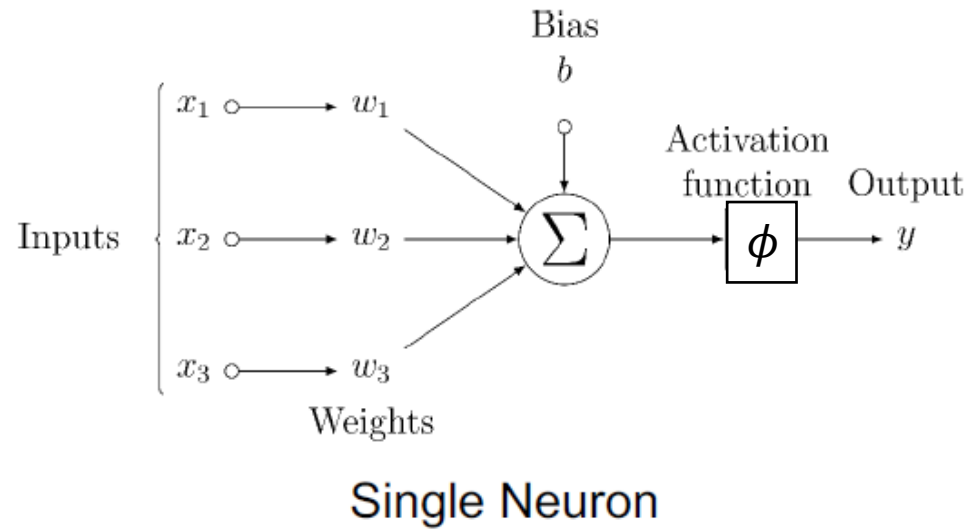
- Non-linear functions



Simplest neural network

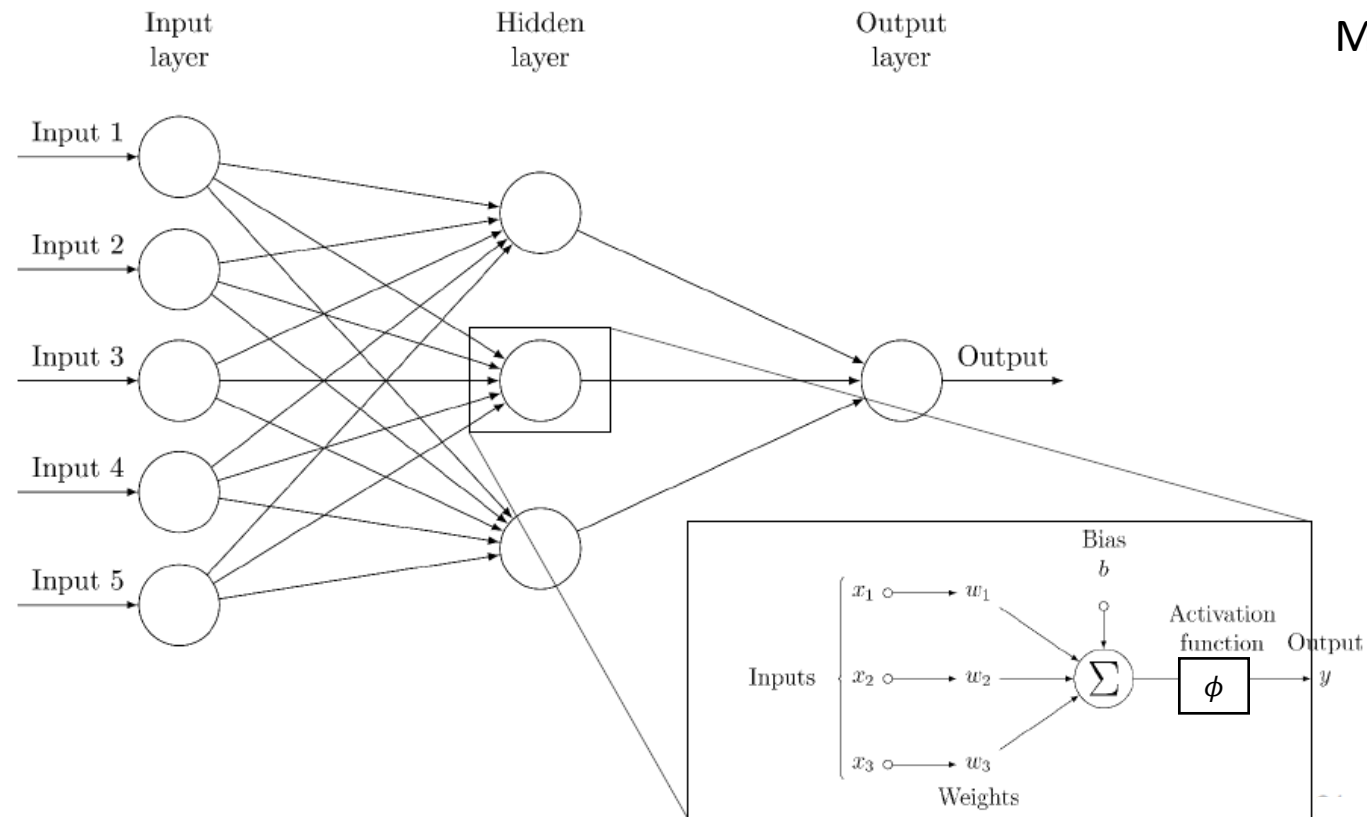
# Neural Network

- Cascaded layers





# Neural Network

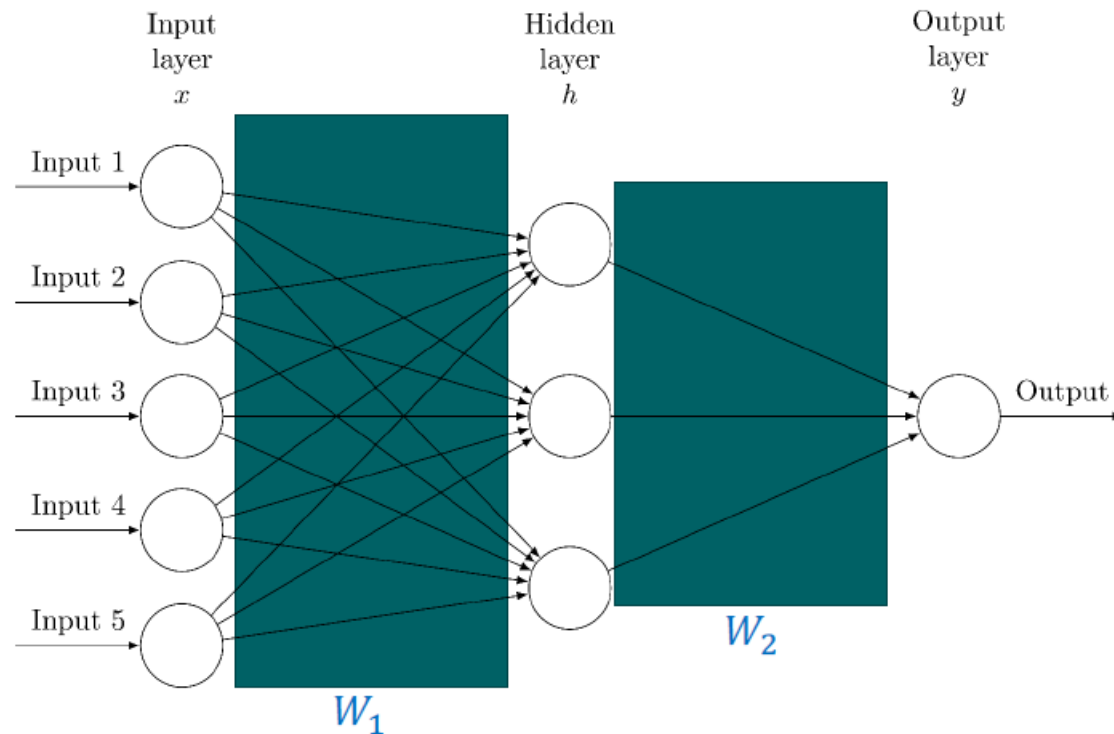


# Neural Network – Forward Pass

---

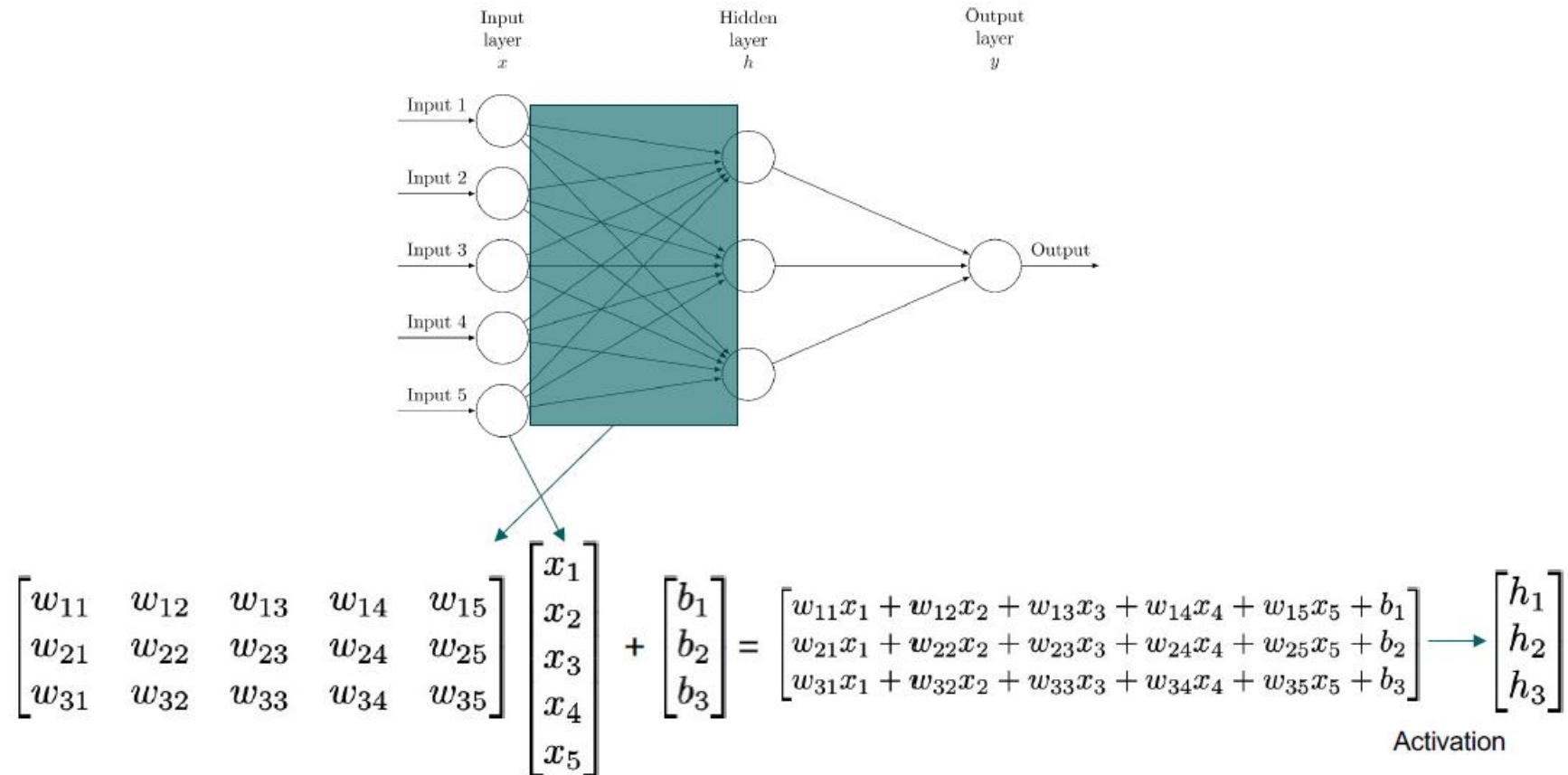
- Given
  - A neural network with a fixed set of weights
  - An input sample to be classified
- How to get the classification?
- A layer is applied by multiplying its weight matrix with its input vector

# Neural Network – Forward Pass



$$h = \phi(W_1x + b_1)$$
$$y = \phi(W_2h + b_2)$$

# Neural Network – Forward Pass



# Chained Linear Classifiers

---

- Remember
  - A linear classifier provides class scores by calculating  $y = \phi(Wx + b)$
  - A neural network is a chain of linear classifiers with activation functions
- A neural network is some function like this:

$$y = \phi(W_3\phi(W_2\phi(W_1x + b_1) + b_2) + b_3)$$

- We can chain this as deep as we want

# Activation Functions

---

- After each layer, an **activation function** is used
- An activation function can be **any function**
- However, **non-linearity** is required for a more expressive model
- Why do we need this?
  - Remember: Linear classifiers can separate data linearly
  - A neural network with a linear activation function is still a linear classifier
  - Non-linear activation functions enable us to learn non-linear relations

# Using Linear Activations

---

- Consider a linear activation  $\phi(a) = a$

$$\begin{aligned}y &= \phi(W_3\phi(W_2\phi(W_1x + b_1) + b_2) + b_3) \\&= W_3(W_2(W_1x + b_1) + b_2) + b_3 \\&= W_3(W_2W_1x + W_2b_1 + b_2) + b_3 \\&= W_3W_2W_1x + W_3W_2b_1 + W_3b_2 + b_3 \\&= \mathbf{W}x + \mathbf{b}\end{aligned}$$

constant term

# Chained Linear Classifiers

---

$$y = \phi(W_3\phi(W_2\phi(W_1x + b_1) + b_2) + b_3)$$



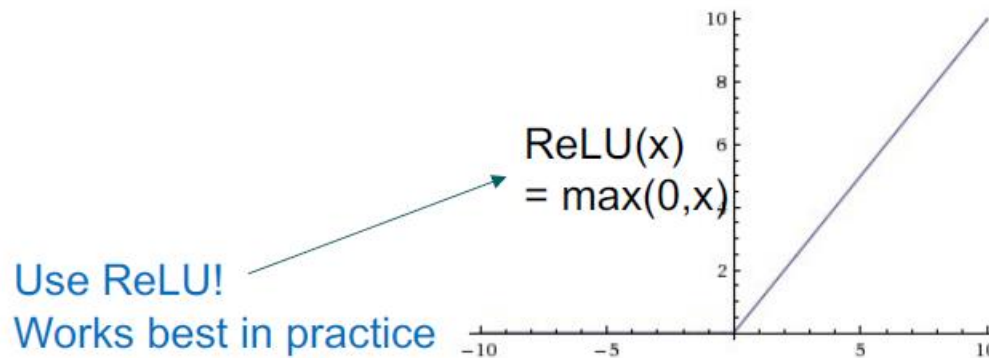
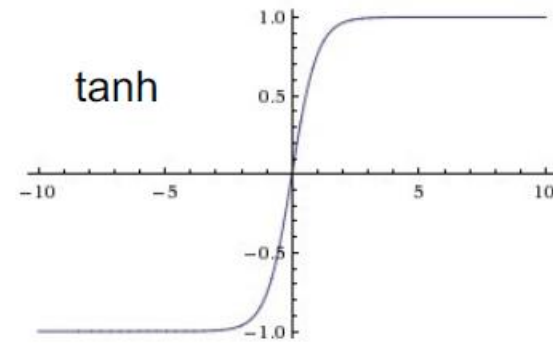
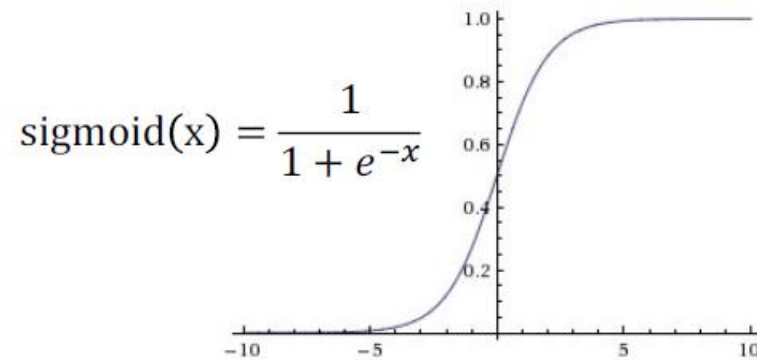
$$y = \textcolor{green}{softmax}(W_3 \textcolor{green}{tanh}(W_2 \textcolor{green}{tanh}(W_1x + b_1) + b_2) + b_3)$$

Activation function



# Activation Functions


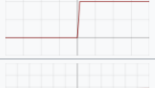


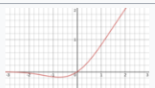





- Some common activation functions:



Some more:

- Heaviside
- Leaky ReLU
- SeLU
- ...

# Activation Functions

Name	Plot	Function, $g(x)$	Derivative of $g$ , $g'(x)$	Range	Order of continuity
Identity		$x$	1	$(-\infty, \infty)$	$C^\infty$
Binary step		$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	0	$\{0, 1\}$	$C^{-1}$
Logistic, sigmoid, or soft step		$\sigma(x) \doteq \frac{1}{1 + e^{-x}}$	$g(x)(1 - g(x))$	$(0, 1)$	$C^\infty$
Hyperbolic tangent (tanh)		$\tanh(x) \doteq \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$1 - g(x)^2$	$(-1, 1)$	$C^\infty$
Rectified linear unit (ReLU) <sup>[9]</sup>		$(x)^+ \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max(0, x) = x \mathbf{1}_{x>0}$	$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	$[0, \infty)$	$C^0$
Gaussian Error Linear Unit (GELU) <sup>[5]</sup>		$\frac{1}{2}x \left( 1 + \operatorname{erf} \left( \frac{x}{\sqrt{2}} \right) \right)$ $= x\Phi(x)$	$\Phi(x) + x\phi(x)$	$(-0.17 \dots, \infty)$	$C^\infty$
Softplus <sup>[9]</sup>		$\ln(1 + e^x)$	$\frac{1}{1 + e^{-x}}$	$(0, \infty)$	$C^\infty$
Exponential linear unit (ELU) <sup>[10]</sup>		$\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter $\alpha$	$\begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ 1 & \text{if } x = 0 \text{ and } \alpha = 1 \end{cases}$	$(-\alpha, \infty)$	$\begin{cases} C^1 & \text{if } \alpha = 1 \\ C^0 & \text{otherwise} \end{cases}$
Scaled exponential linear unit (SELU) <sup>[11]</sup>		$\lambda \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameters $\lambda = 1.0507$ and $\alpha = 1.67326$	$\lambda \begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$(-\lambda\alpha, \infty)$	$C^0$
Leaky rectified linear unit (Leaky ReLU) <sup>[12]</sup>		$\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$\begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	$(-\infty, \infty)$	$C^0$

And many more!!

# Universal Approximation Theorem

---

- “A neural network can approximate any continuous function”

**Theorem 1.** *Let  $\sigma$  be any continuous discriminatory function. Then finite sums of the form*

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\overset{w_j}{\downarrow} y_j^T x + \overset{b_j}{\downarrow} \theta_j) \quad (2)$$

*are dense in  $C(I_n)$ . In other words, given any  $f \in C(I_n)$  and  $\varepsilon > 0$ , there is a sum,  $G(x)$ , of the above form, for which*

$$|G(x) - f(x)| < \varepsilon \quad \text{for all } x \in I_n.$$

Cybenko G, “Approximation by Superpositions of a Sigmoidal Function”

# Universal Approximation Theorem

---

- The theorem states that  $\exists N$  and  $\exists w_j, \theta_j, \alpha_j$  for  $j = 1, 2, \dots, N$  such that

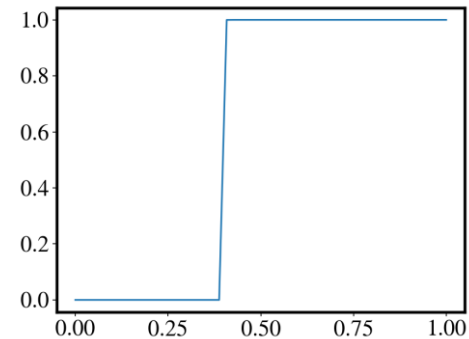
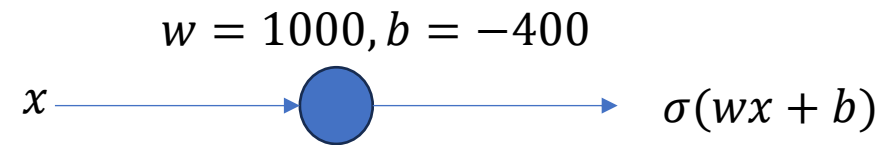
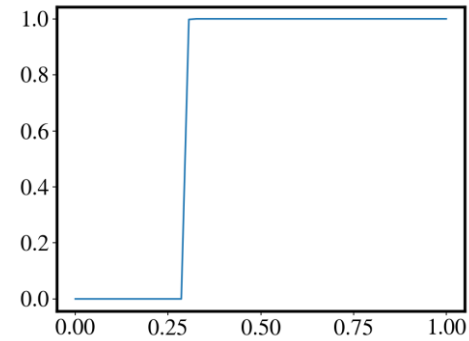
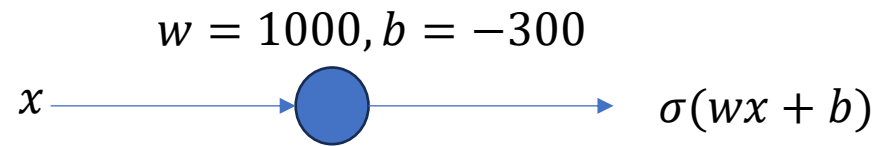
$$|G(x) - f(x)| < \epsilon$$

Sigmoid:  $\sigma(w_j x + b_j)$



# Universal Approximation Theorem

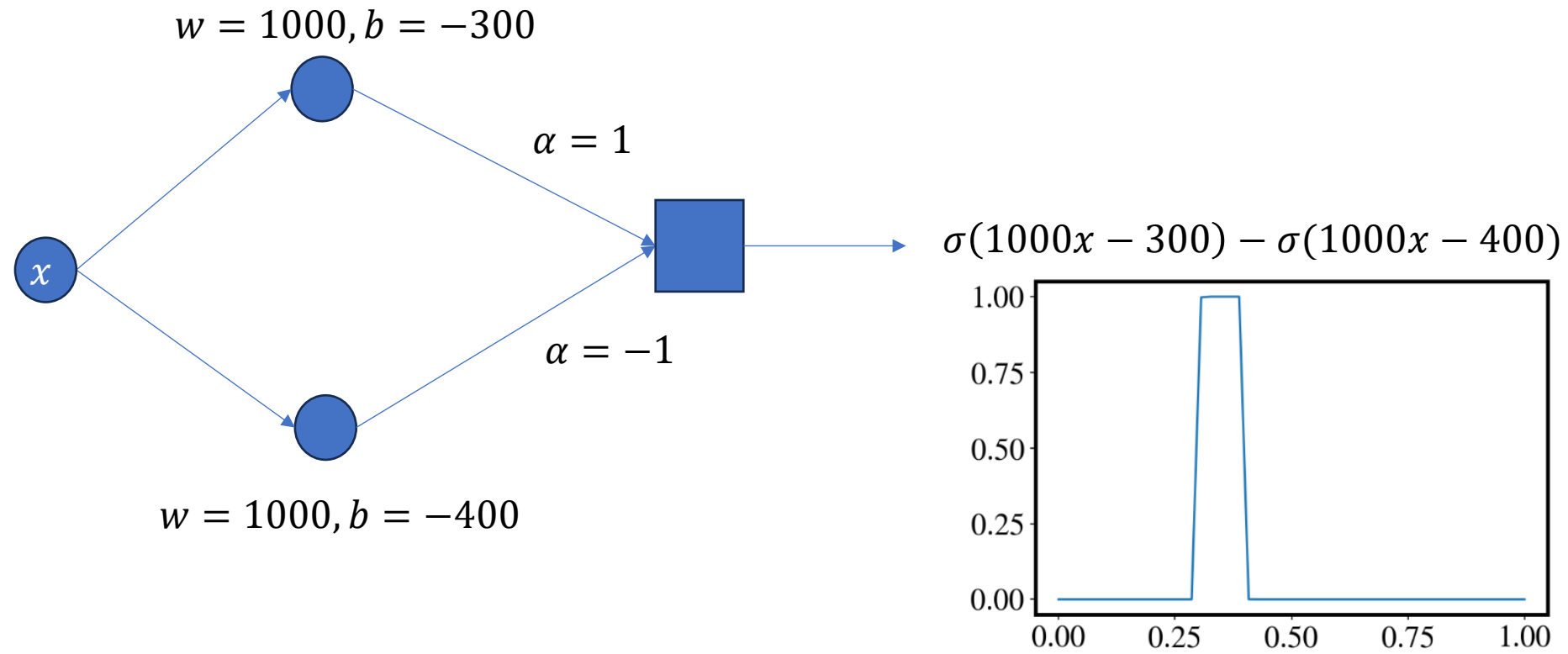
---



Piece

# Universal Approximation Theorem

---

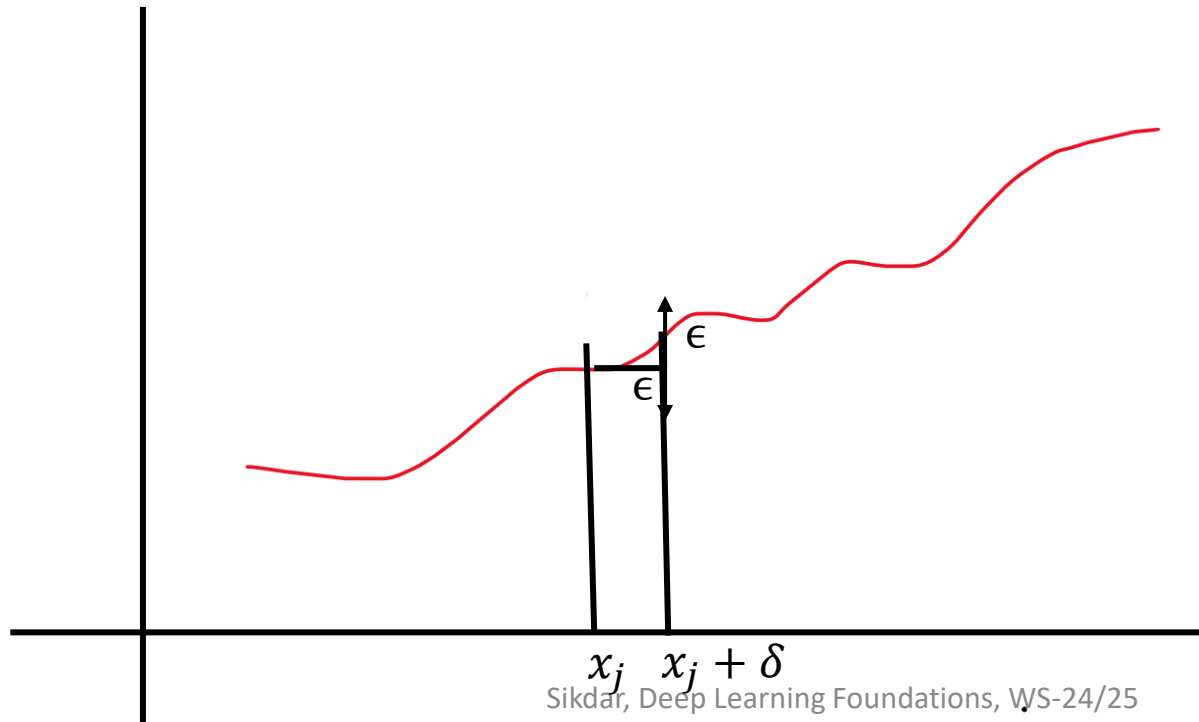


# Universal Approximation Theorem

---

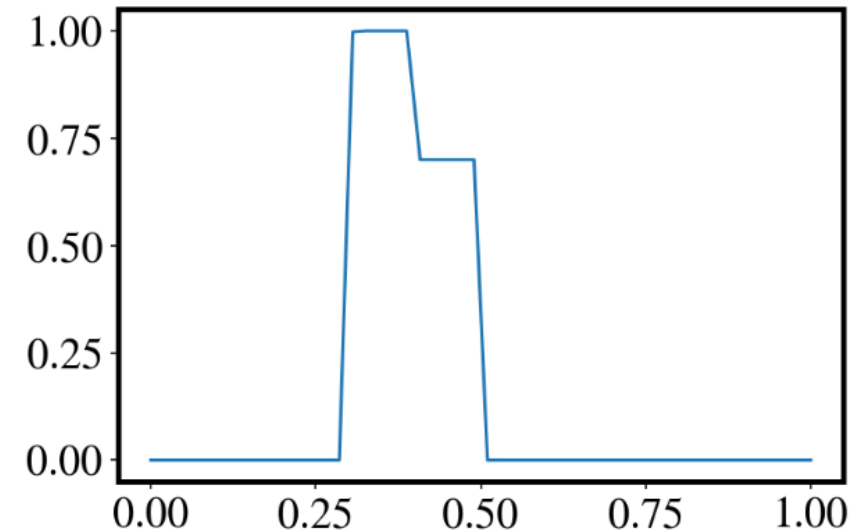
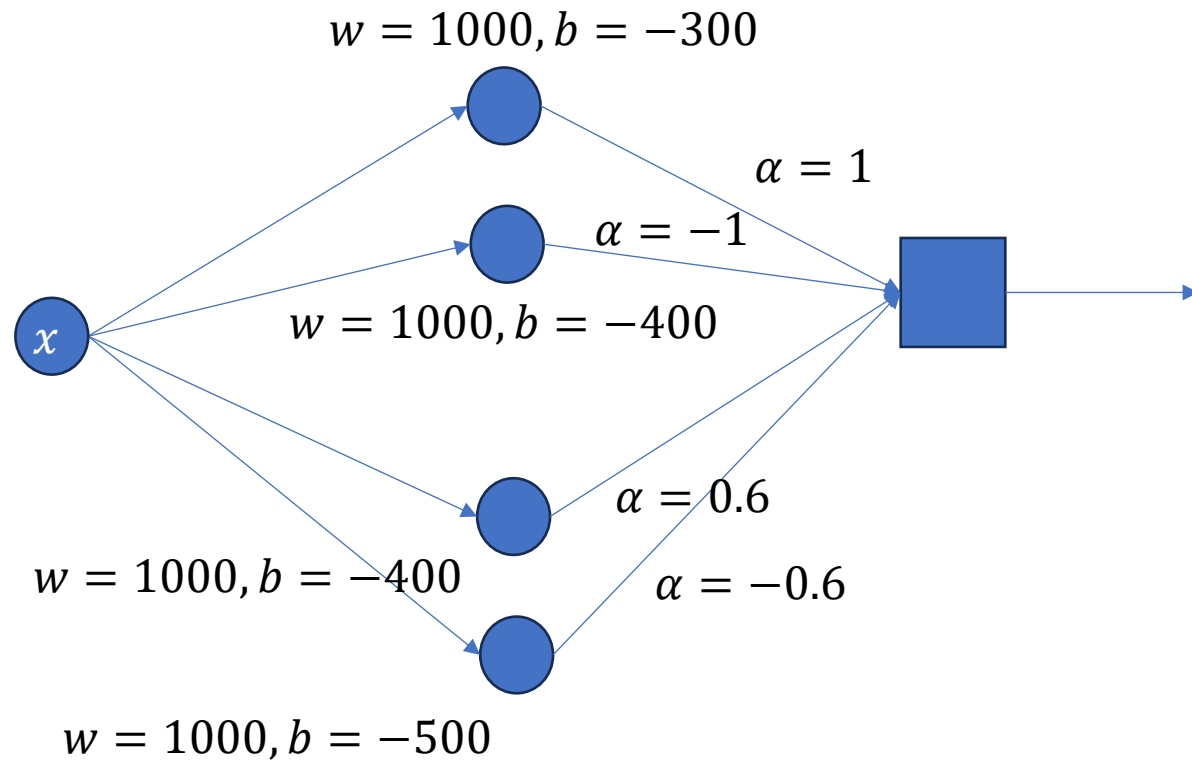
- The theorem states that  $\exists N$  and  $\exists w_j, b_j, \alpha_j$  for  $j = 1, 2, \dots, N$  such that

$$|G(x) - f(x)| < \epsilon$$



# Universal Approximation Theorem

---



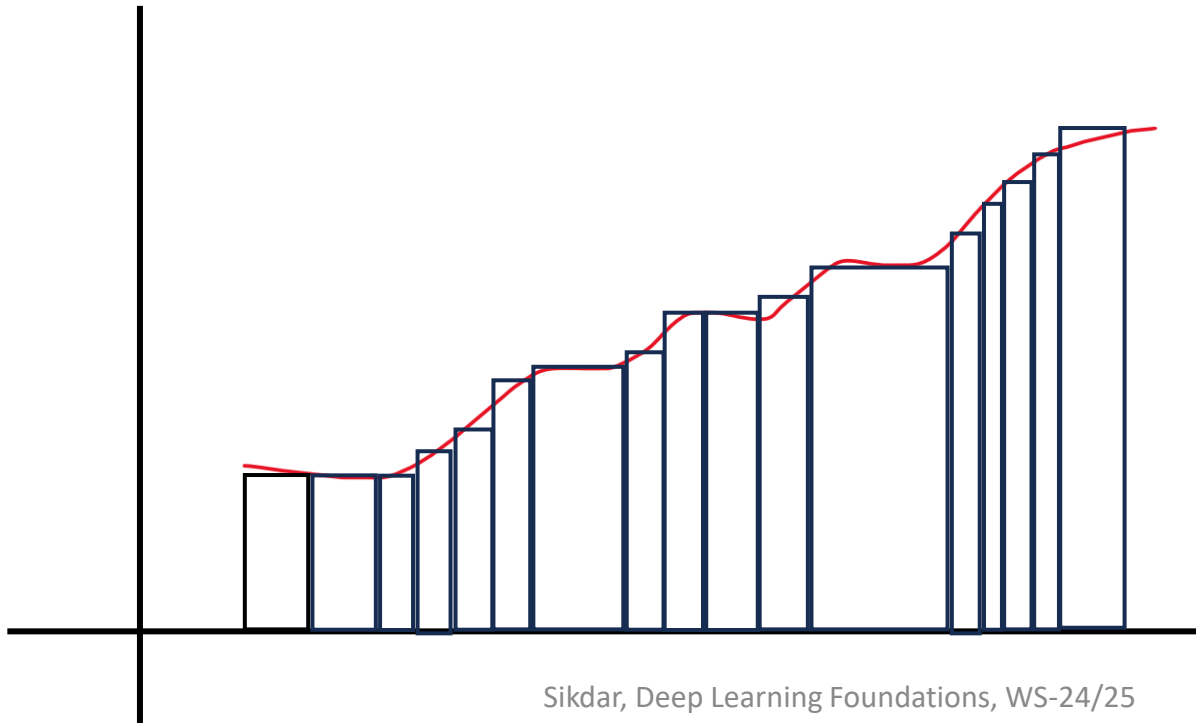


# Universal Approximation Theorem

---

- The theorem states that  $\exists N$  and  $\exists w_j, \theta_j, \alpha_j$  for  $j = 1, 2, \dots, N$  such that

$$|G(x) - f(x)| < \epsilon$$



# Training/Optimization

---

$$\min_W \frac{1}{n} \sum_{i=1}^n l(f_W(x^i), y^i)$$

- Use a search algorithm that starts with some initial guess and repeatedly changes  $W$  such that the loss gets smaller and smaller until we reach a point where the loss is minimized
- Gradient descent algorithm

# Training/Optimization

---

$$L(W) = \frac{1}{n} \sum_{i=1}^n l(f_W(x^i), y^i)$$

- Gradient descent algorithm

$$W_j = W_j - \alpha \frac{\partial}{\partial W_j} L(W)$$

Learning rate

# Training/Optimization

---

- Computing the gradient

$$\begin{aligned}\frac{\partial}{\partial \theta_j} L(W) &= \frac{\partial}{\partial W_j} \frac{1}{2} (f_W(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (f_W(x) - y) \cdot \frac{\partial}{\partial W_j} (f_W(x) - y)\end{aligned}$$

If we consider the linear model:  $f_W(x) = W^T x$

$$\begin{aligned}&= (f_W(x) - y) \cdot \frac{\partial}{\partial W_j} (\sum_i W_i x_i - y) \\ &= (f_W(x) - y) x_j\end{aligned}$$

# Training/Optimization

---

```
Loop {  
    for i=1 to n, {  
         $W_j = W_j - \alpha(f_W(x^i) - y^i)x_j^i$  (for every j)  
    }  
}
```

# Evaluation

---

- How do we know that the trained model is good?
- Split the data into train/validation/test



- Training only on the train split
- Periodically check performance on validation split
- Select model with the best performance in validation split
- Evaluate on the test split

# Evaluation

---

- Measuring performance
- Regression:
  - Mean squared error, Mean absolute error
- Classification
  - Accuracy: fraction of cases where the predicted labels are correct

# Evaluation

---

- K-fold cross validation
- Simple training-and-test can be subject to random fluctuations
- Better: (ten-fold) cross-validation
  - Split labeled randomly data into 10 parts
  - Use 9 parts for training, last part for testing
  - Repeat 10 times, each part is used for testing once
  - Average results
- Applied if limited training data is available and training time is reasonable



# Summary

---

- Supervised learning
  - Linear/shallow neural models
  - Loss
  - Optimization
- Universal approximation theorem
- Next time
  - Optimization and regularization of deep neural models