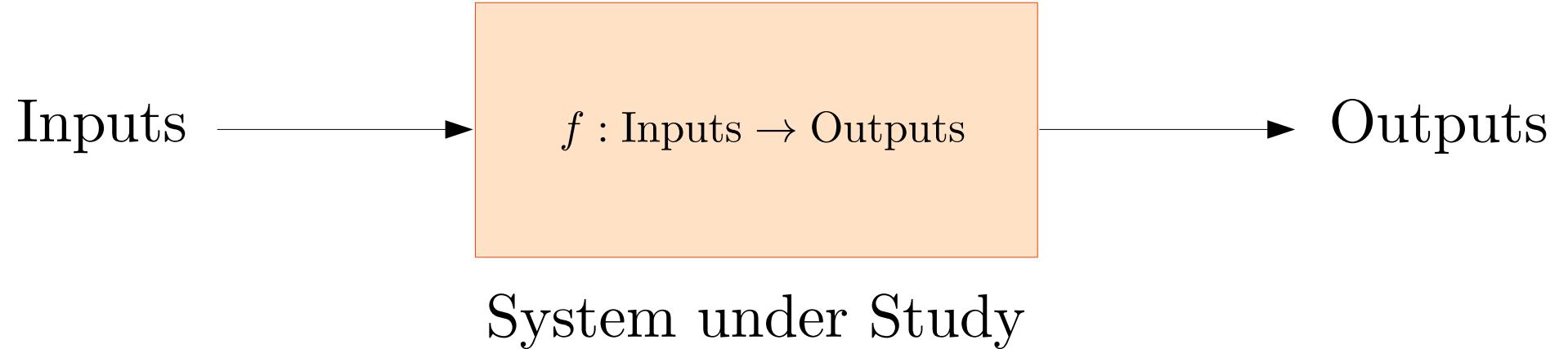


# What is Machine Learning?

---

Sanjay Arora  
AI Center of Excellence  
Red Hat

# Why and how does something work?



**Question:** what are the rules governing the system i.e. what is  $f$ ?

# The Scientific Method

Most successful approach till now:

- Observe system with different inputs and measure outputs
- Guess rules or guess f
- Make predictions from guessed rules
- Compare predictions to outputs
- If predictions “match” outputs, you **might** have the correct rules
- If predictions don’t match outputs, you are definitely wrong

# Caveats

- A scientific theory can never be **proven** to be correct. There's no proof the sun will rise tomorrow – just an overwhelming possibility.
- Comparison of the predictions and outputs is more subtle than it appears. It depends on the degree of accuracy required. Rules might predict outputs within 10% but not within 1% and that might be sufficient to explain the core elements of the system under study.
- Process not linear in time!
- The process for guessing rules is crucial.
- Simplicity is desirable. The best (in terms of predictive power) scientific theories tend to be simple.

# How to guess the rules?

No standard way but each field has its own intuitions, methods,  
ideas that worked in the past that can be tried again

One way:

Collect data and scan for patterns. How Kepler found the laws for  
planetary motion

# What is Machine Learning?

- Collect data and find patterns to discover rules.
- Let computers scan through several guesses/hypotheses and find the best ones.
- Need to invent algorithms that can scan multiple hypotheses efficiently (in time, memory, amount of data required).

# Terminology: Model

Our guess of the rules

# Terminology: Learning

Process of scanning data to find the best rules

# Algorithms

Every field has a few simple examples that help one build intuition

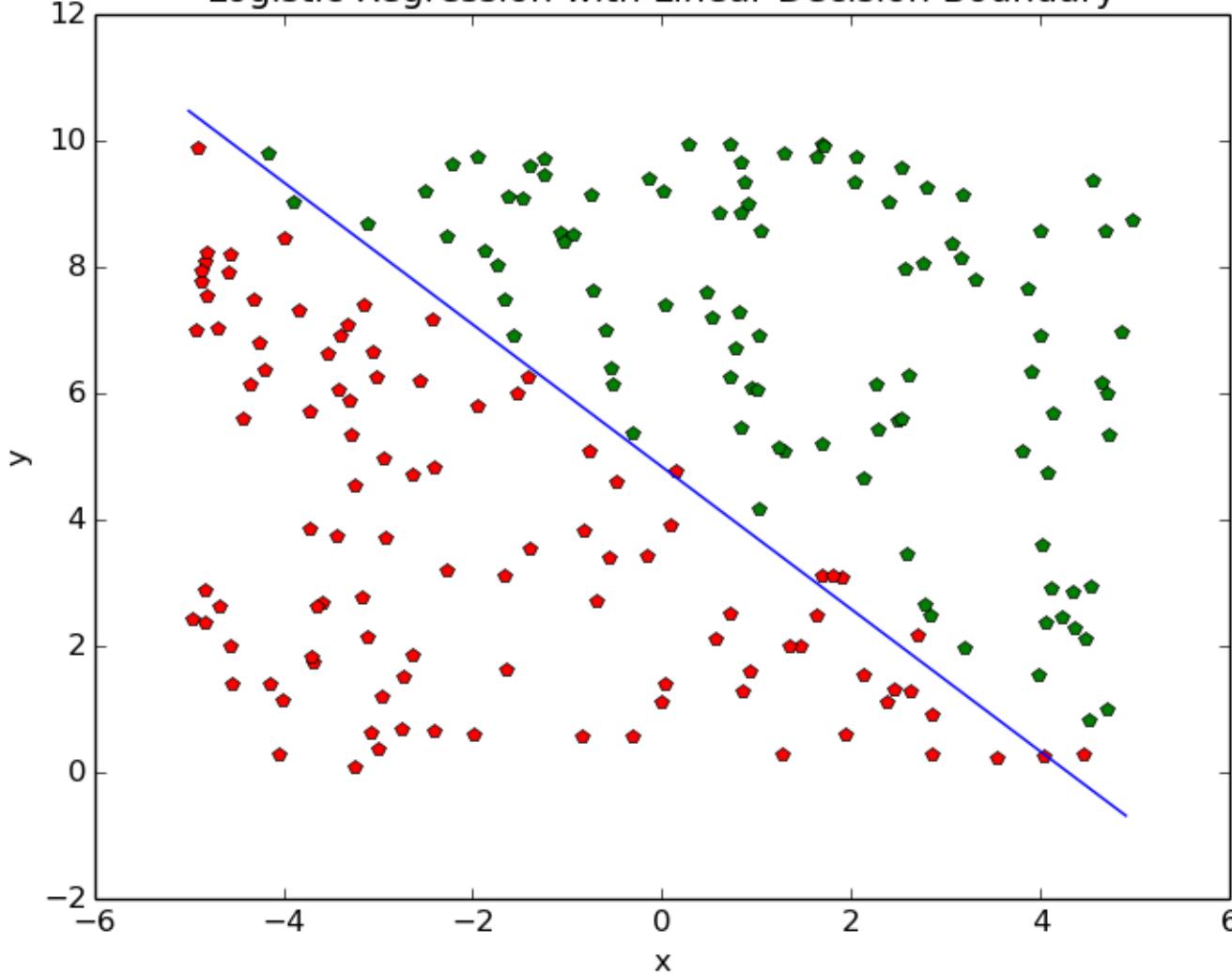
No mathematical details here → Just pictures

# Problem Statement

Given features :  $x_1, x_2, \dots, x_n$

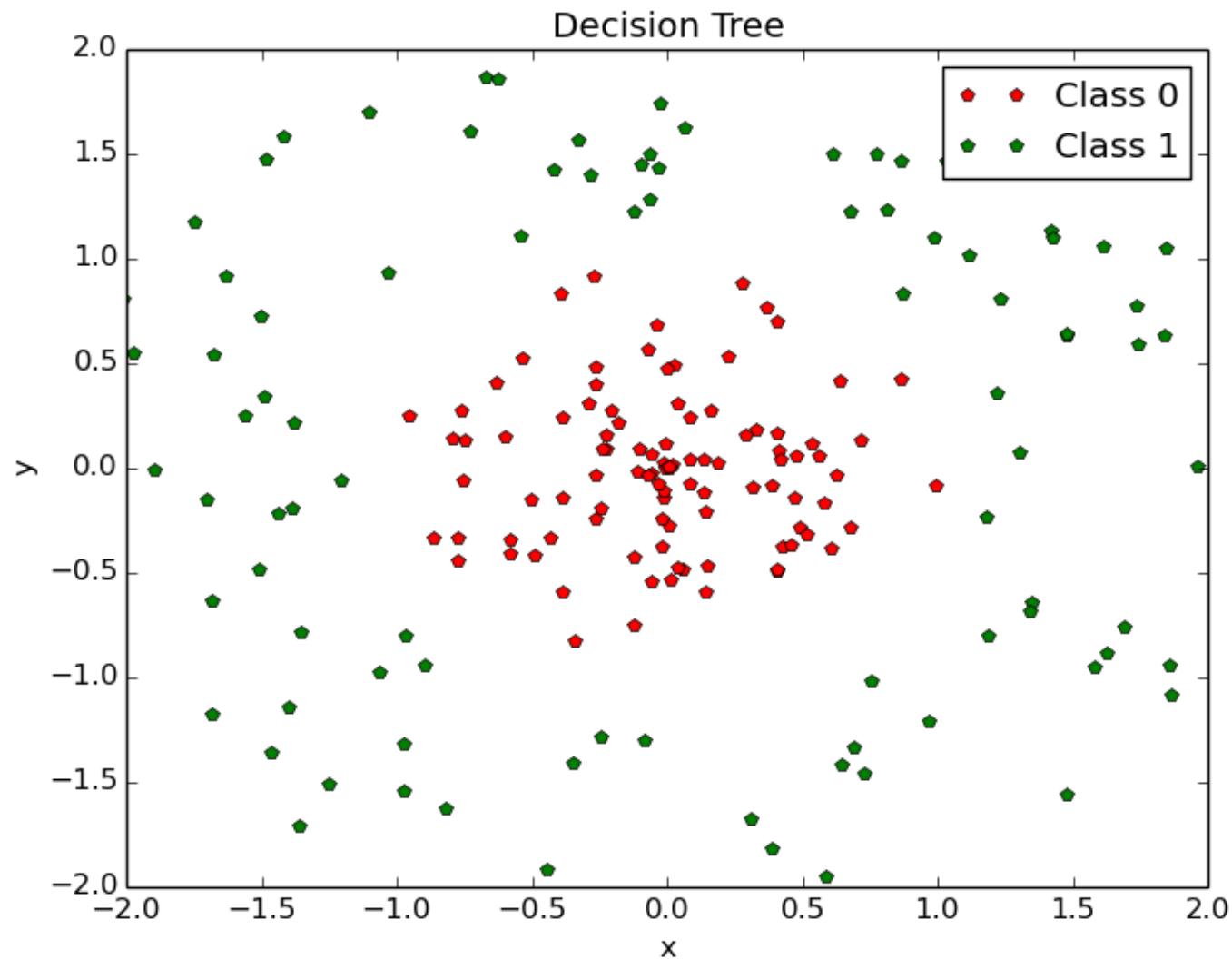
predict class or label,  $y = 0$  or  $1$

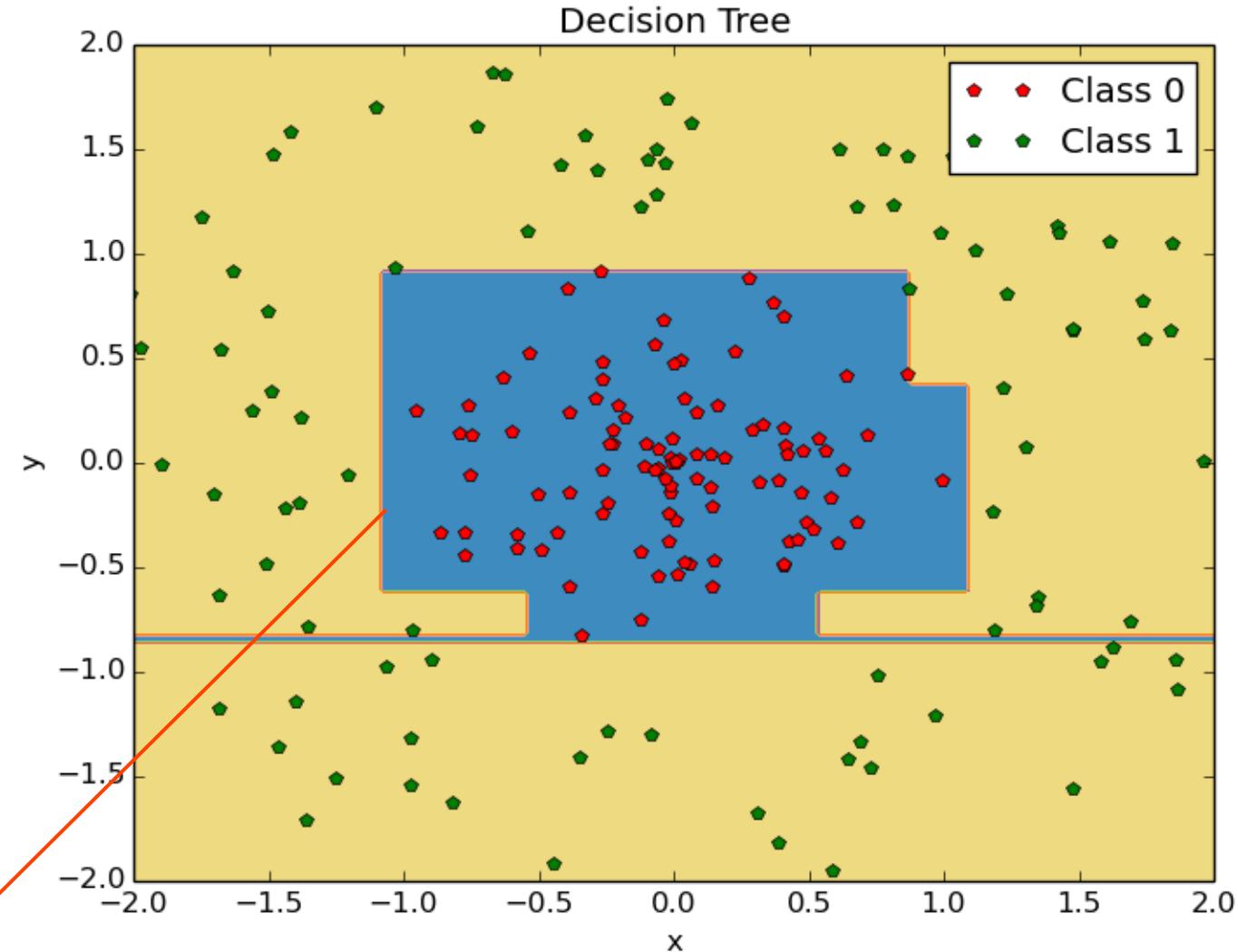
Logistic Regression with Linear Decision Boundary



Logistic Regression = “Draw a line separating red from green”

(Linear) Support Vector Machine =  
“Draw a road separating red from green  
with highest width”



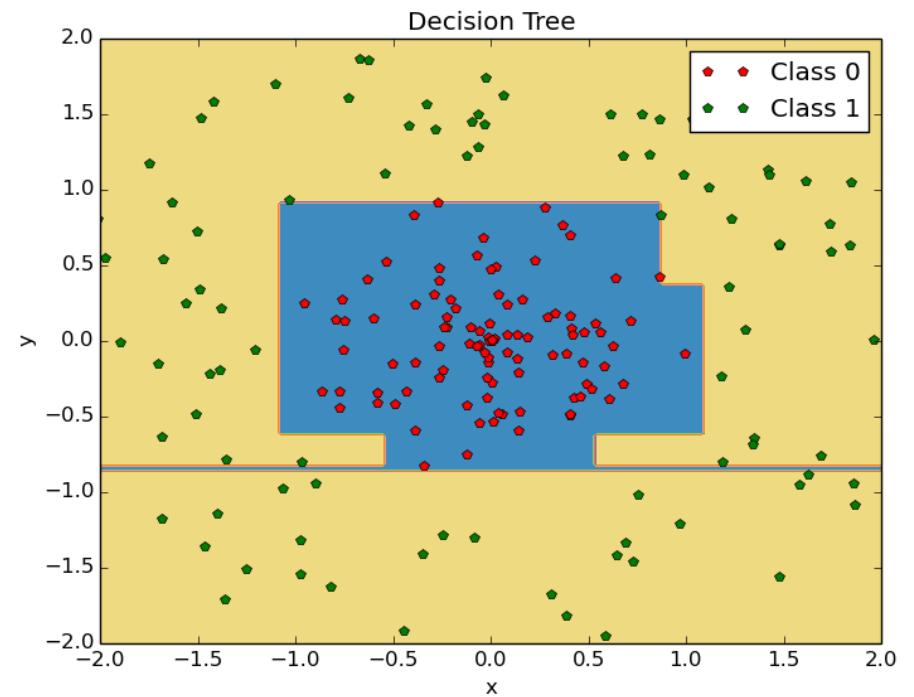
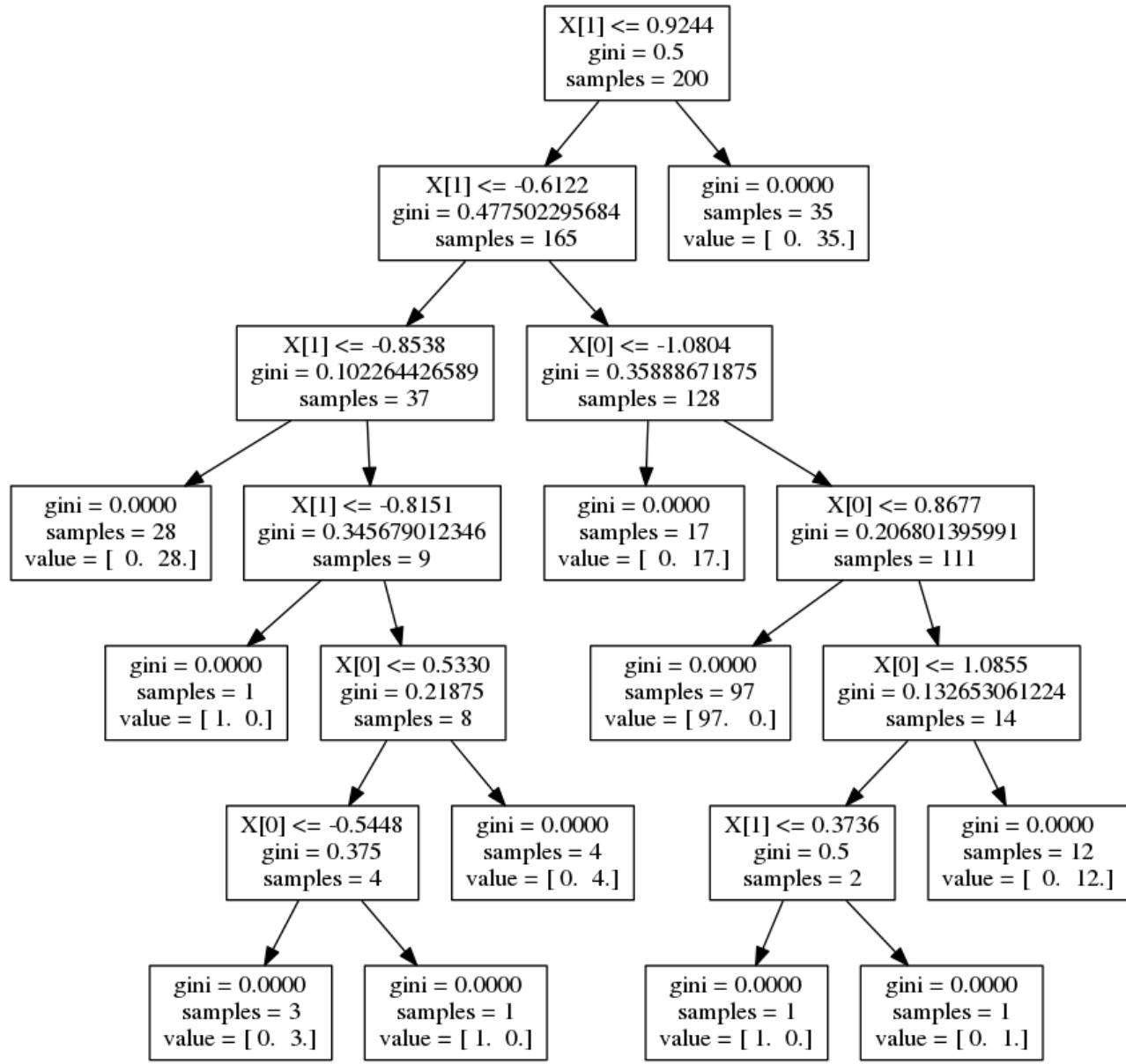


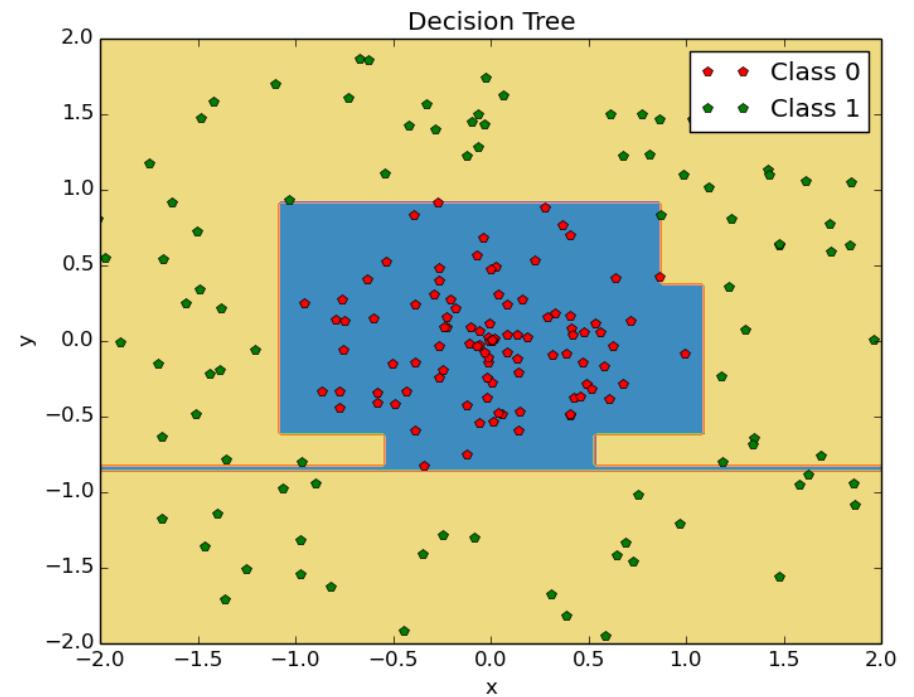
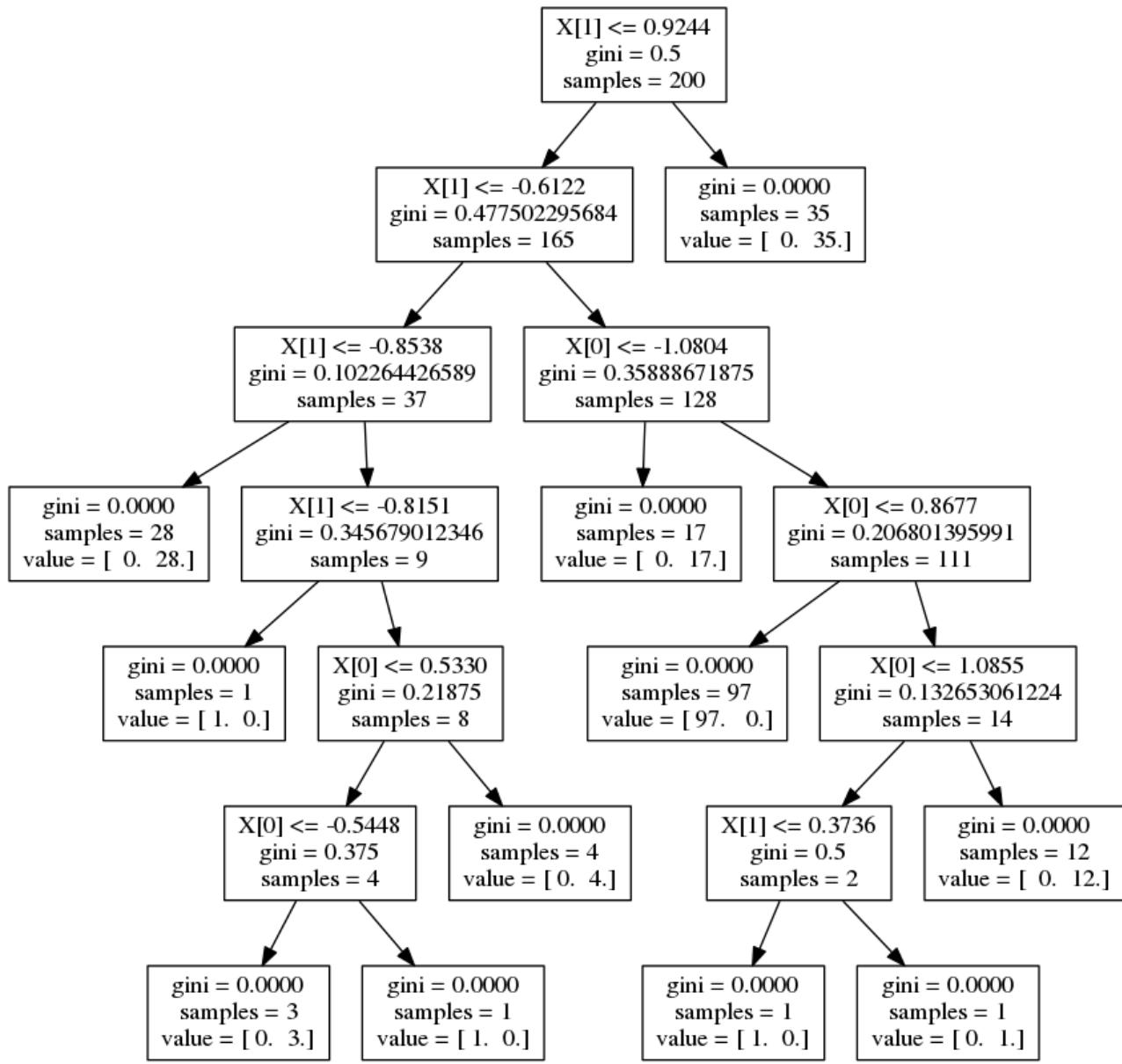
14

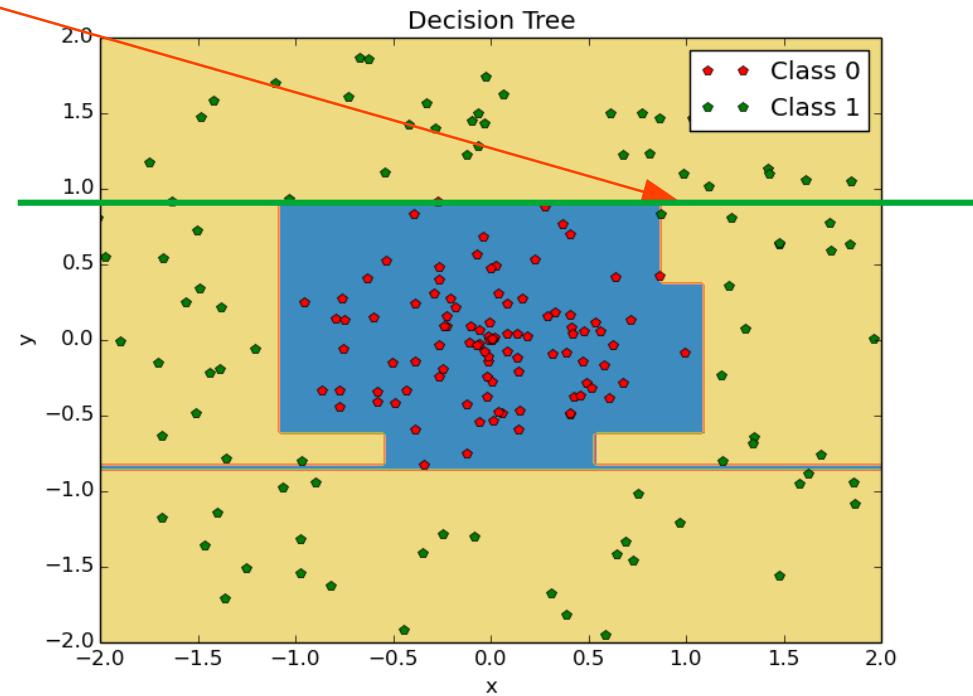
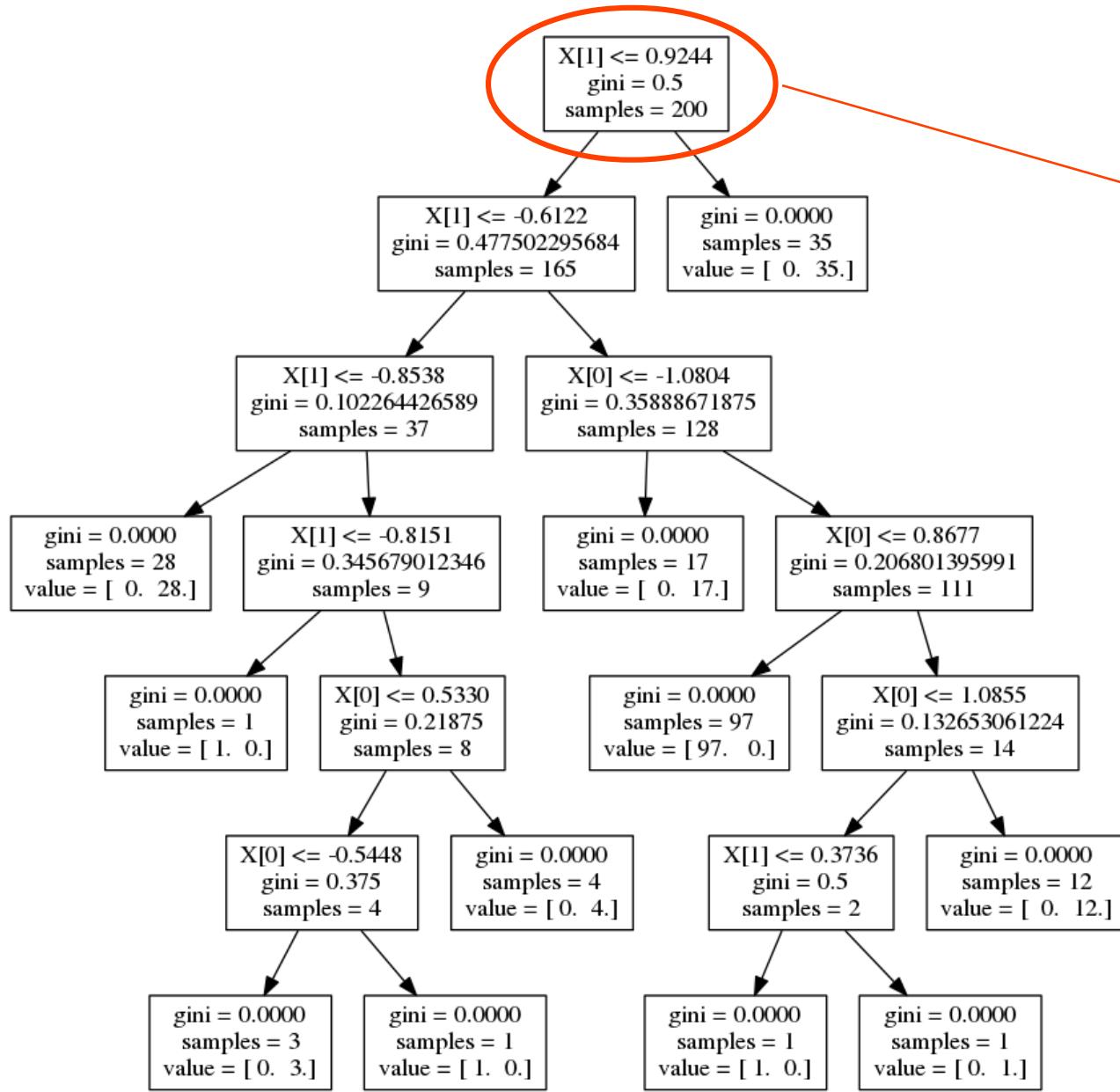
Any if-else statement on x or y

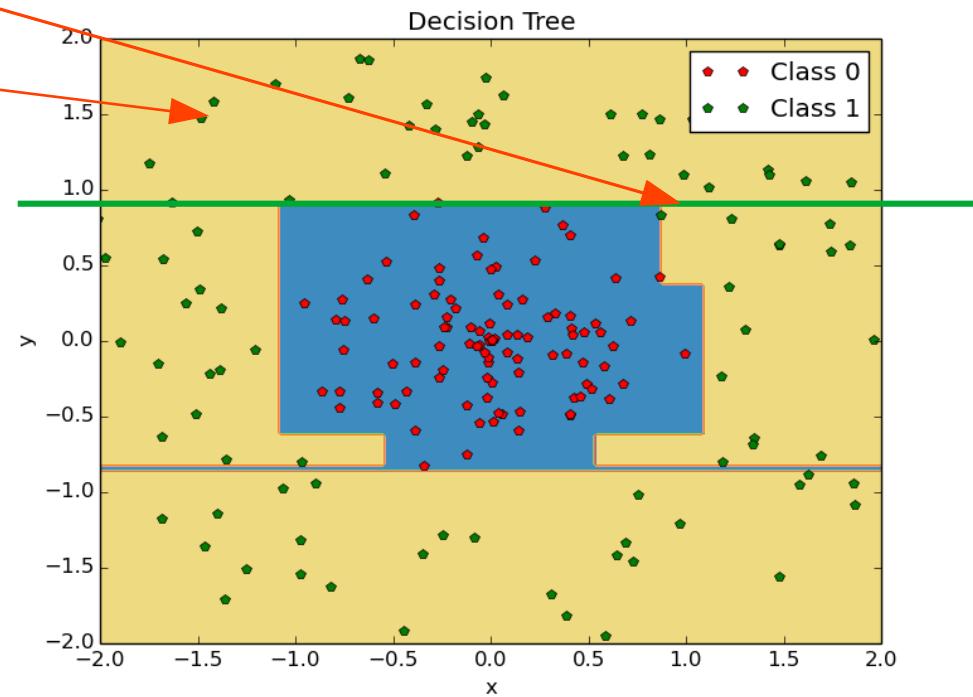
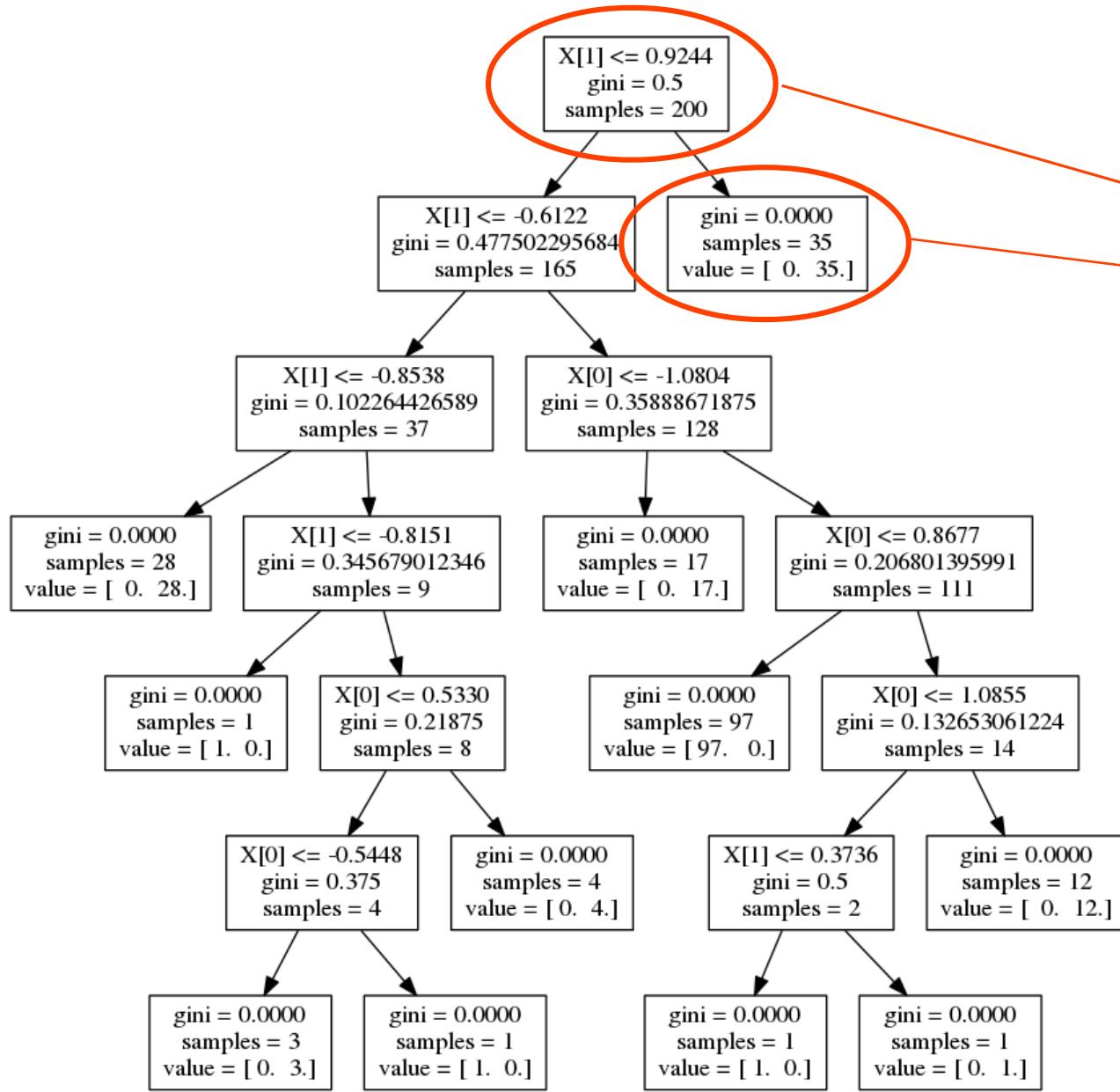


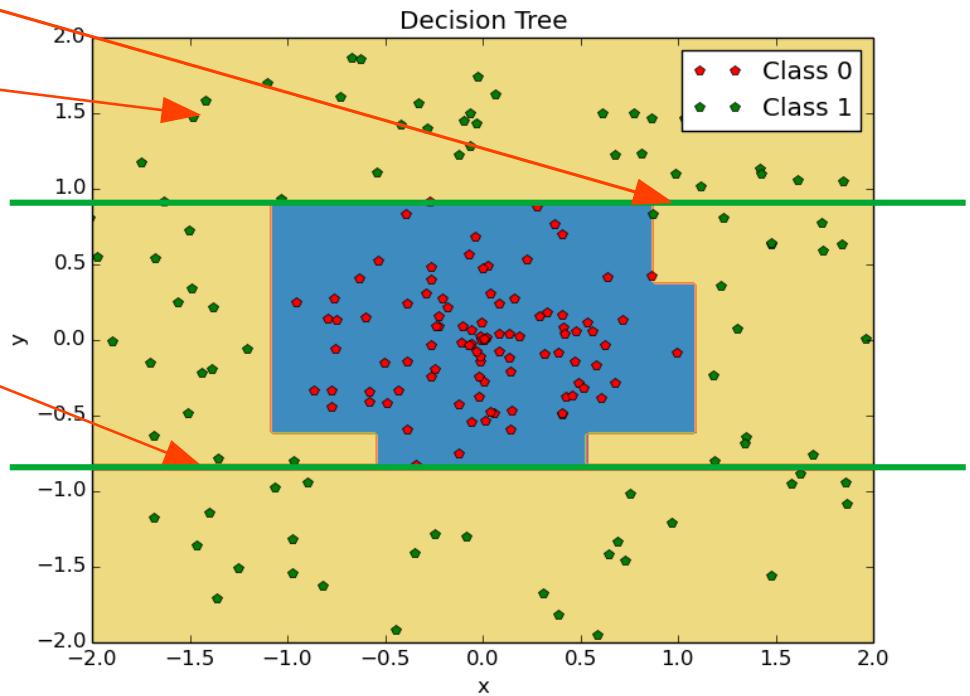
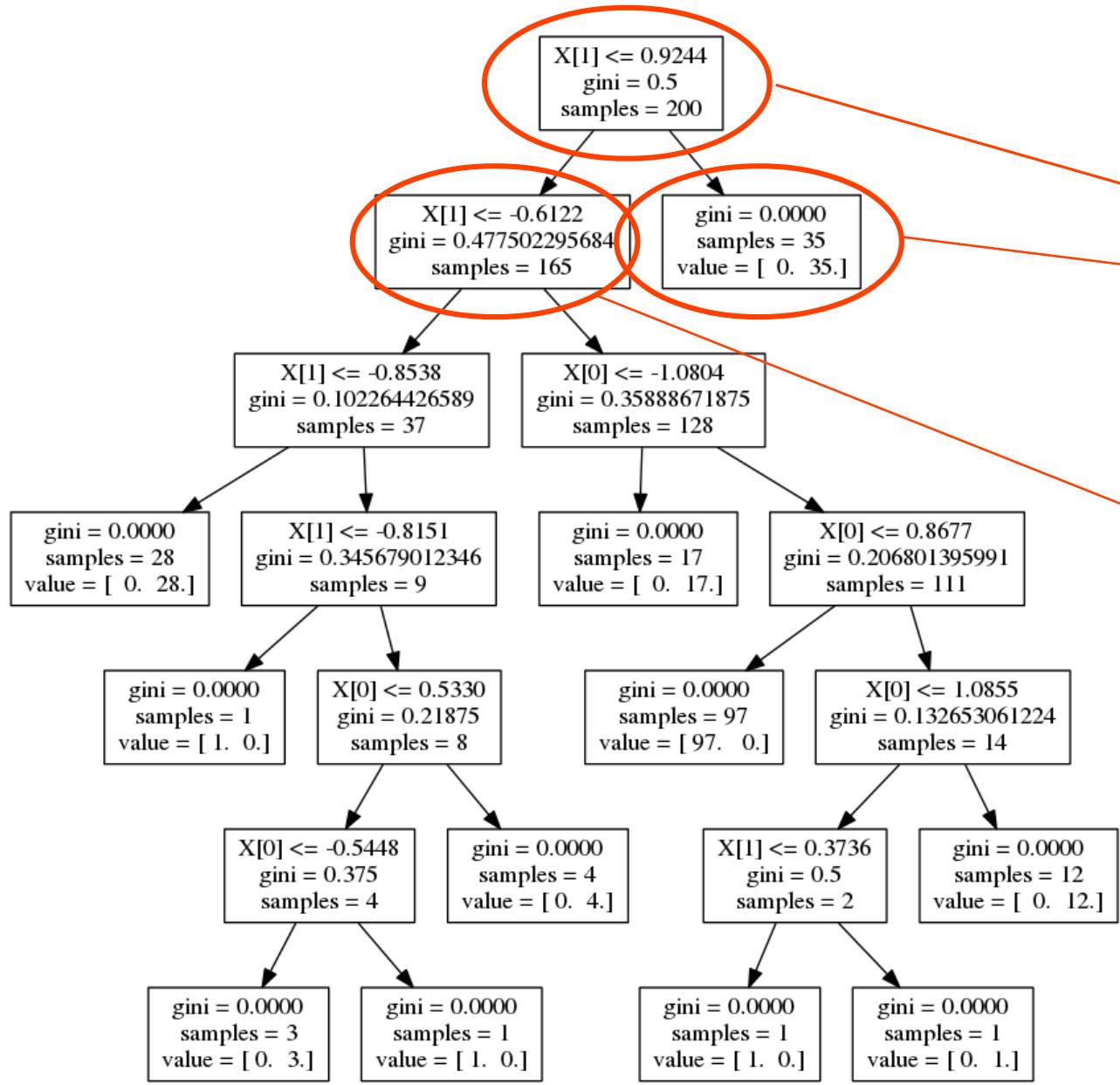
Horizontal or vertical line

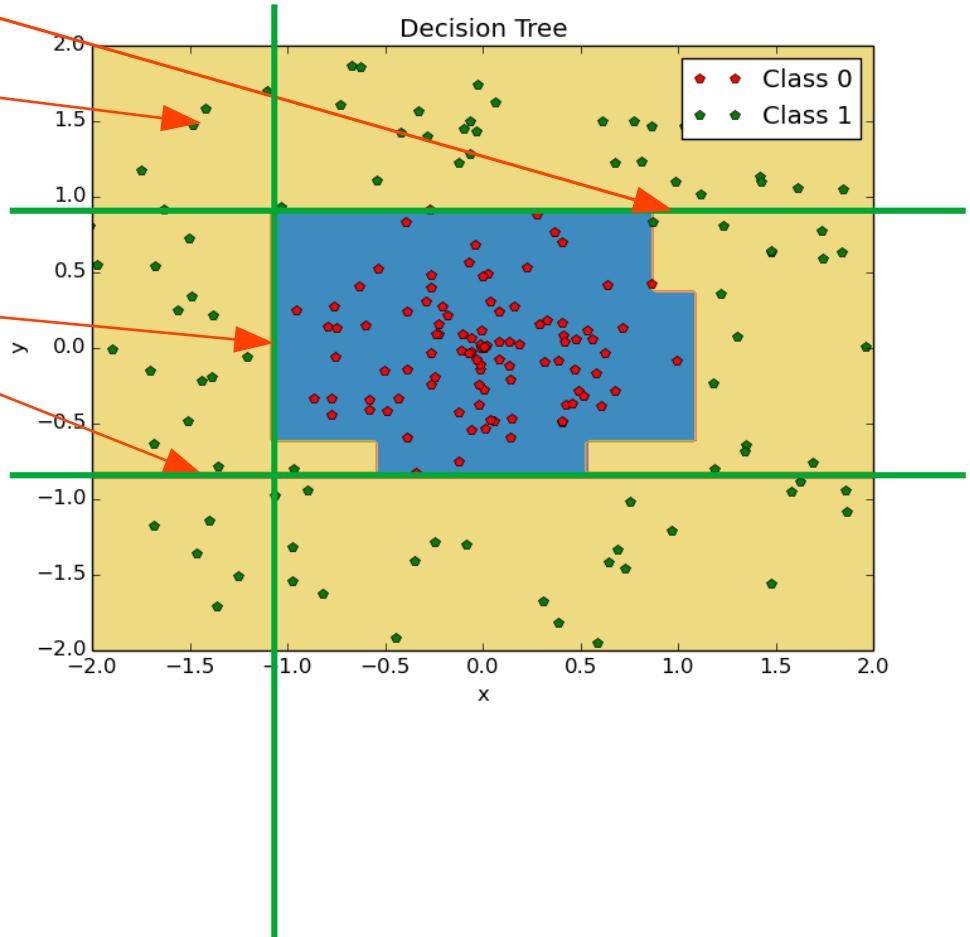
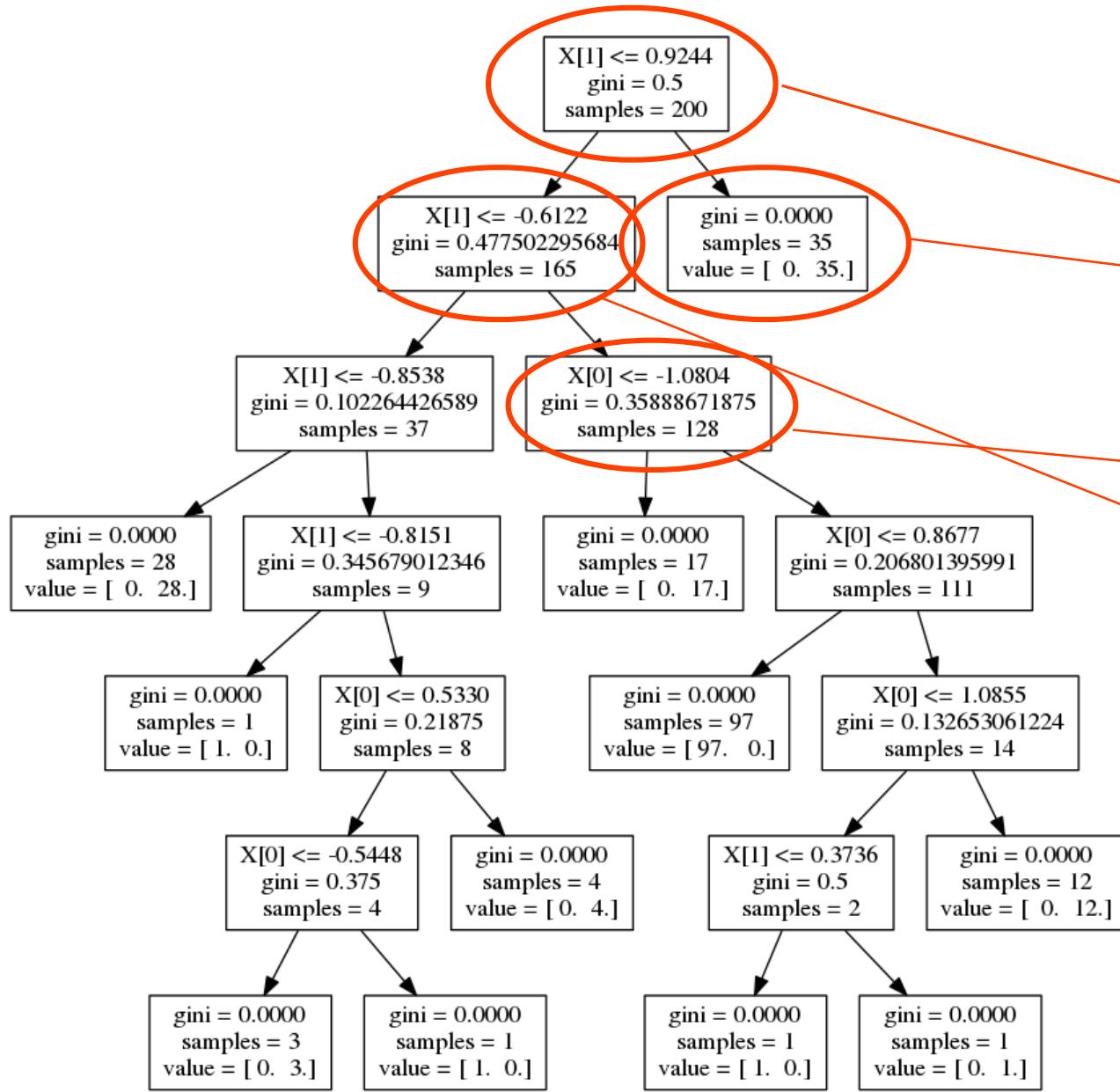




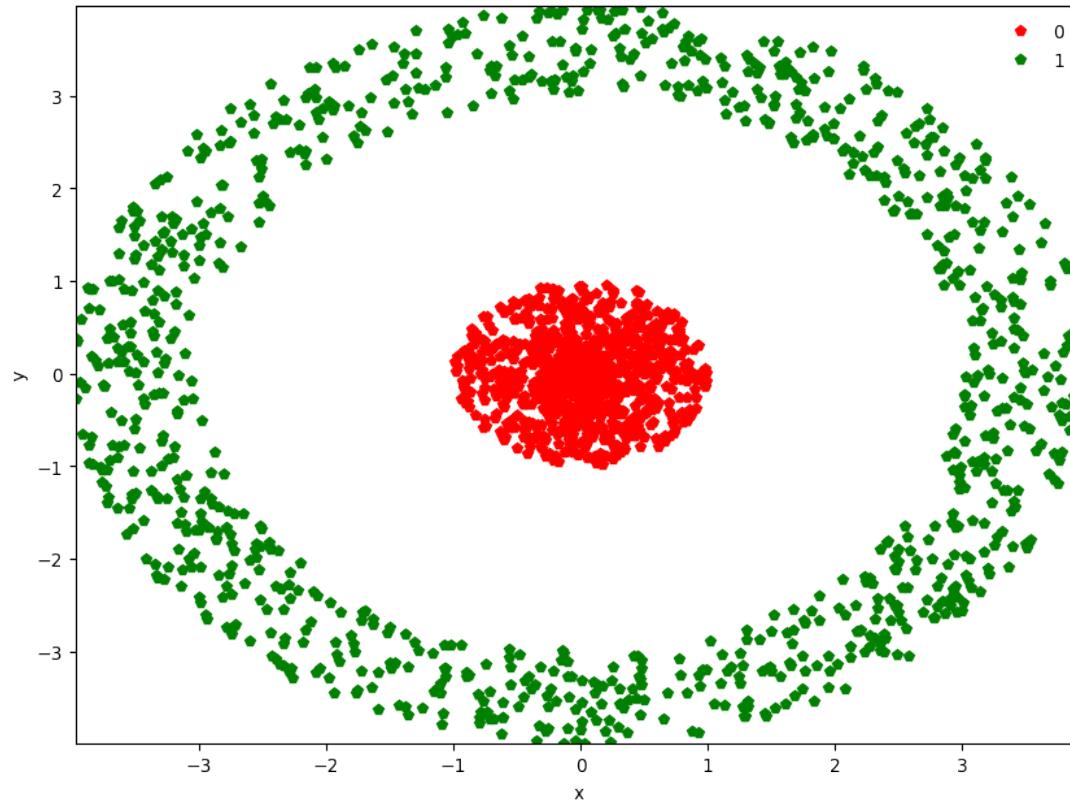






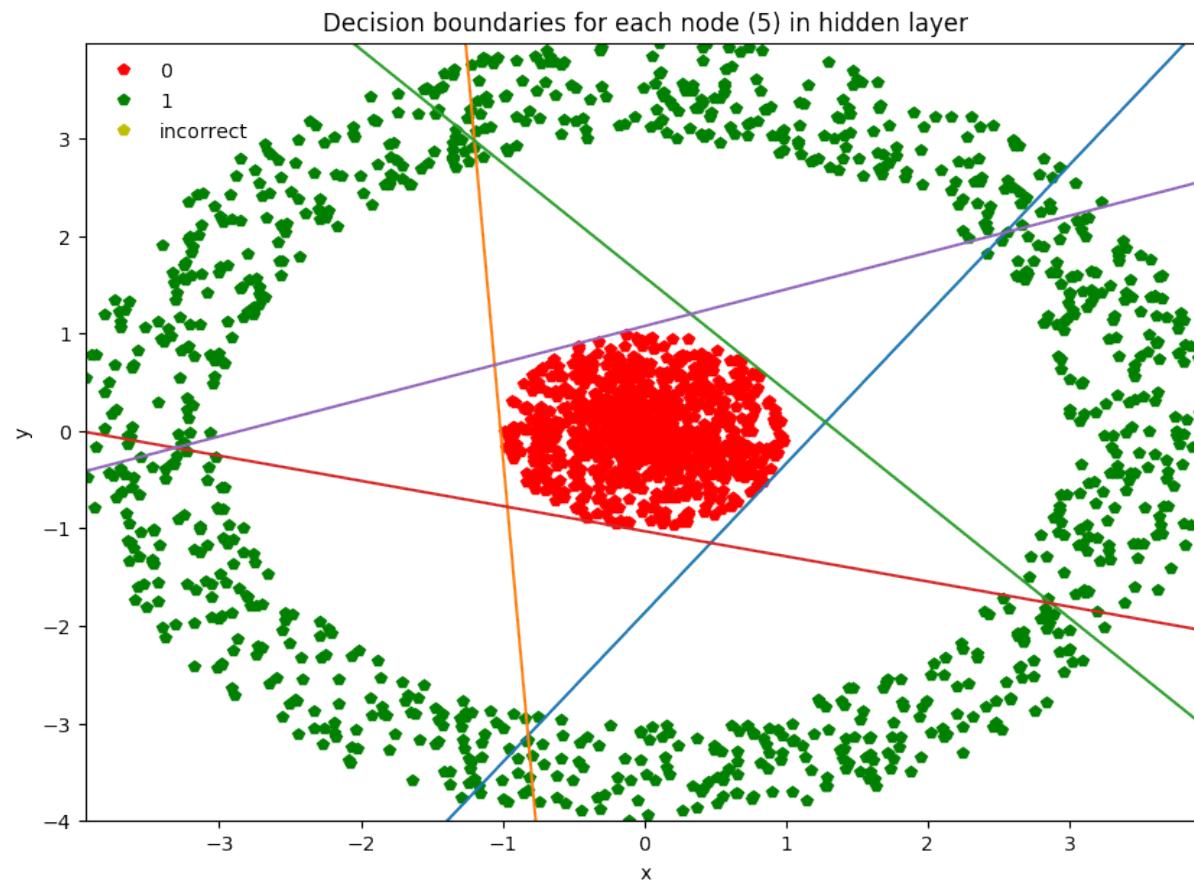


# Problem

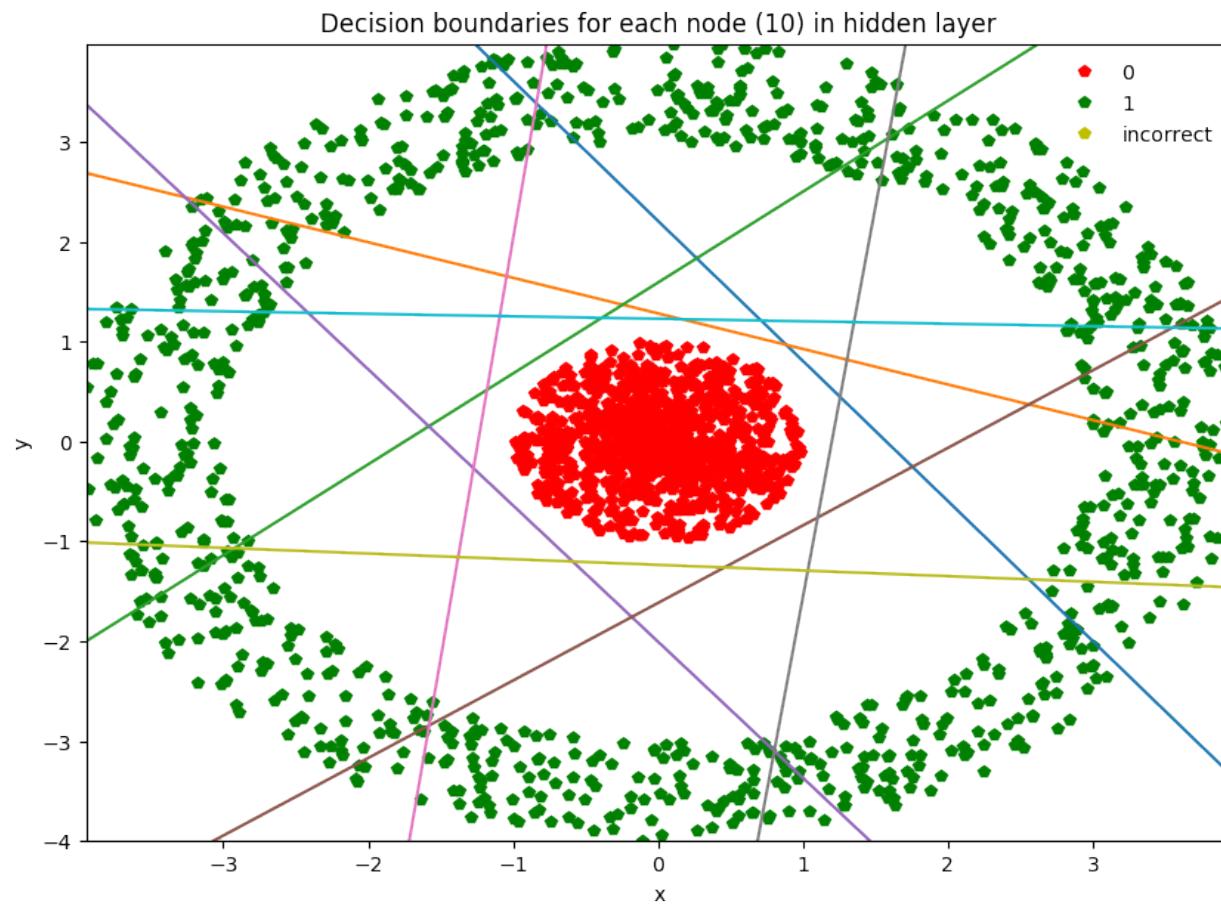


Binary Classification: Given  $(x, y)$ , predict class 0 or 1

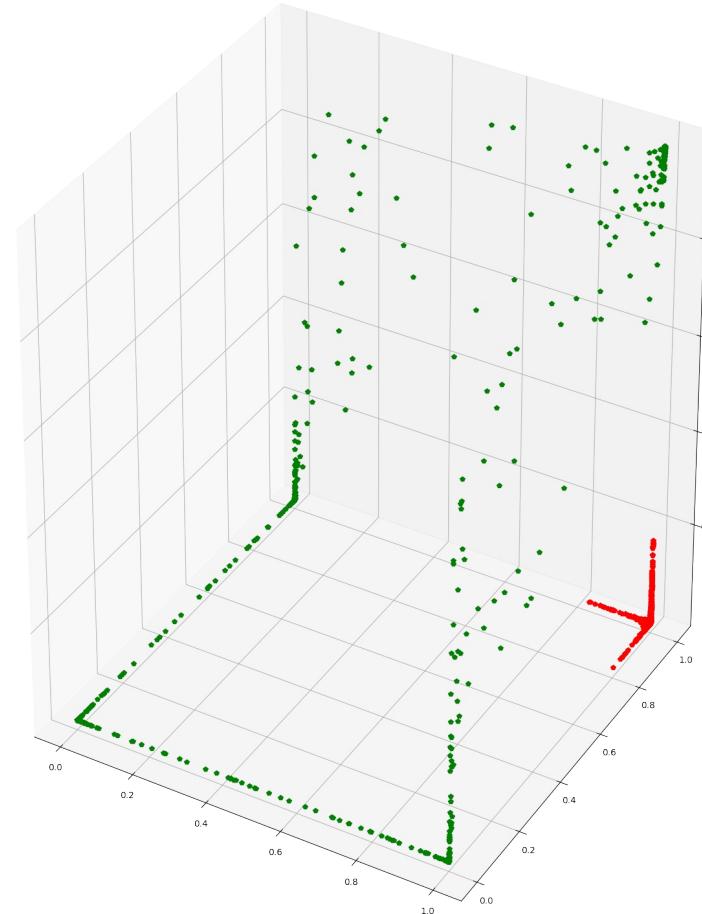
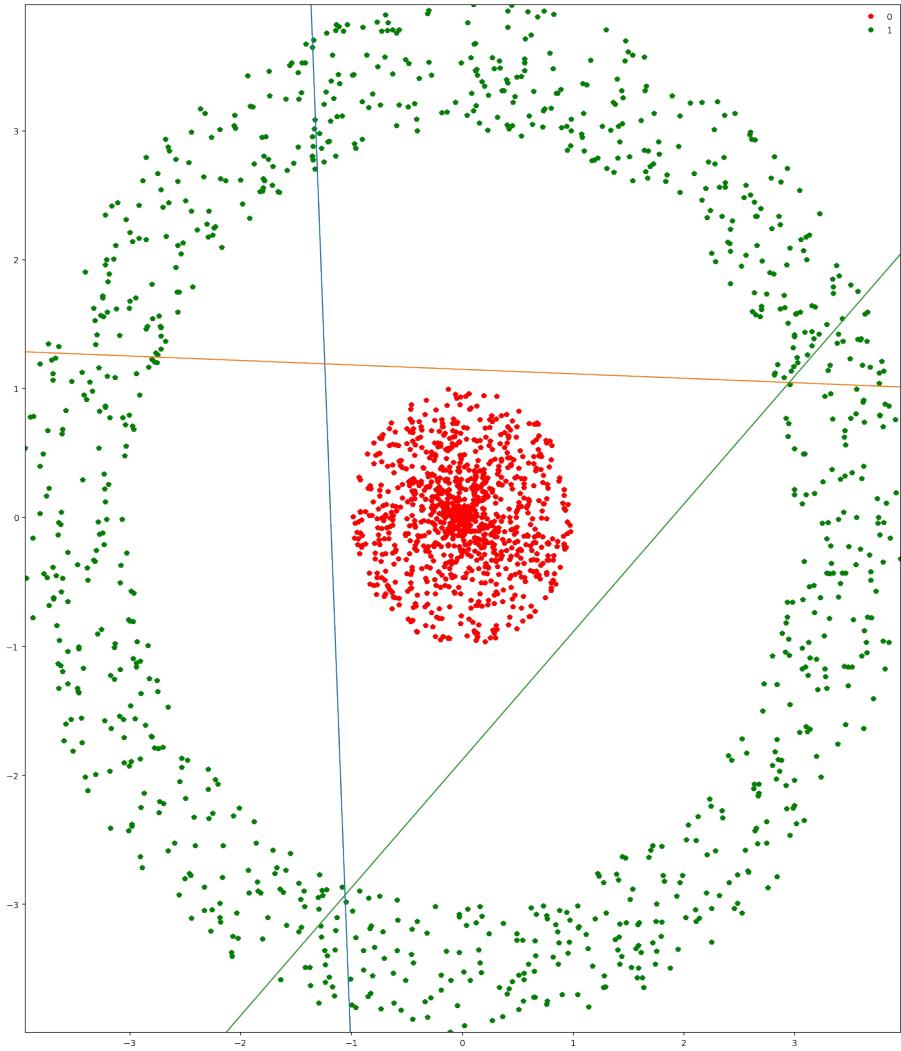
# Hidden Nodes = 5



# Hidden Nodes = 10



# Map 2-dim input to 3-dim space



Many more options:

Combine multiple (weak) trees → Bagging, Random Forests

Sequentially train trees → Boosting

Don't have to be trees → any set of models

Pick K-nearest neighbors and use their results

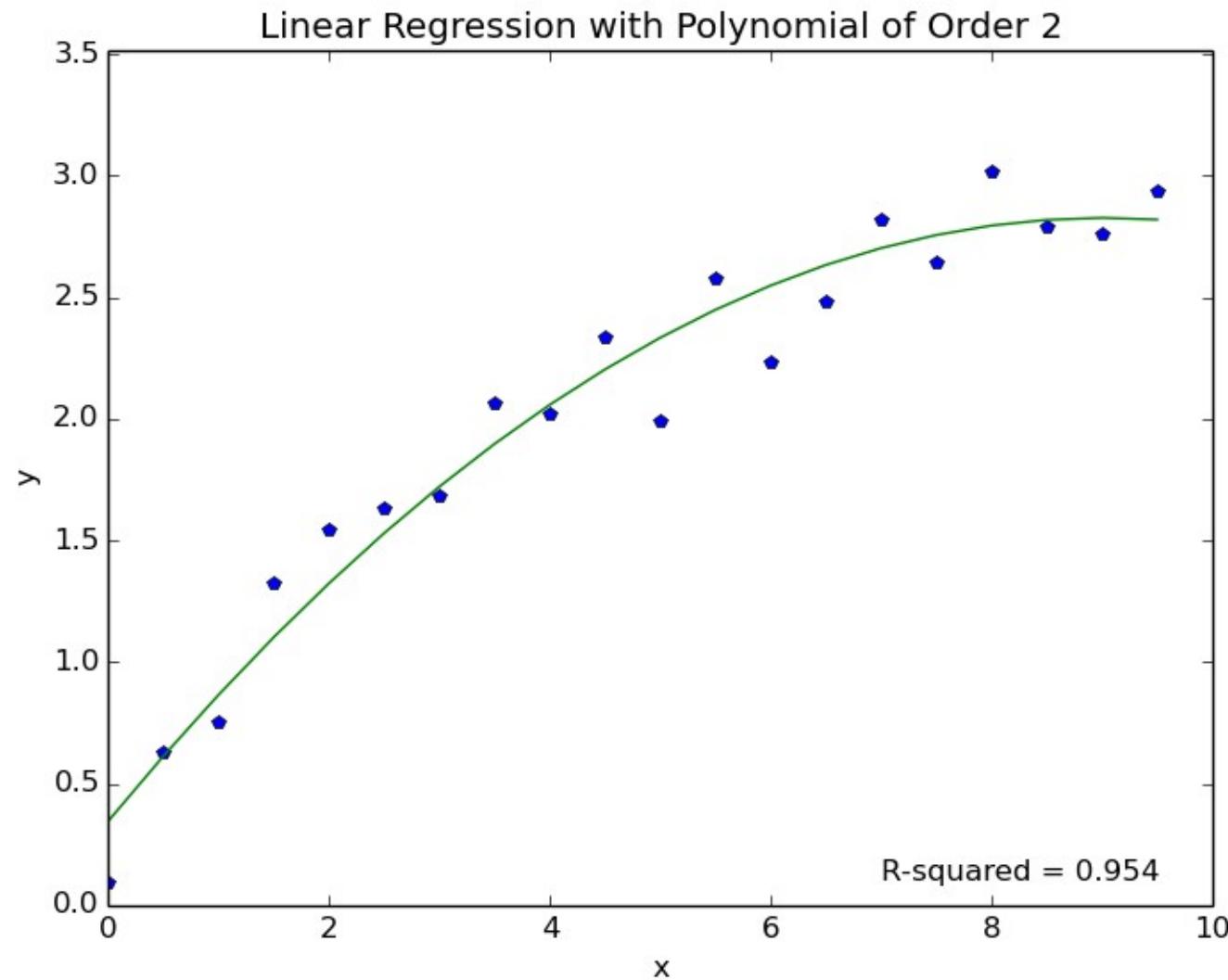
**Things get trickier in higher dimensions!**

# Problem Statement

Given features :  $x_1, x_2, \dots, x_n$

predict real number y

# Example



$$y = w_0 + w_1 x + w_2 x^2 = 0.344 + 0.551x - 0.031x^2$$

Can use all the techniques mentioned before with some modifications

All of this is called **supervised learning**

Supervised = have both inputs and outputs available

OR

Have **feedback** (output) for **every example** (input)

## **Supervised Learning:**

**Optimization:** Turn learning problem into optimization problem

$$\min_x f(x) \rightarrow \min_{\text{model parameters}} \text{error}(\text{pred}, \text{output})$$

**Generalization:** should be able to predict accurately on unseen inputs

What if don't have feedback for every input??

Maybe have only infrequent feedback

What if don't have feedback for every input??

Maybe have only infrequent feedback

College student taking a course

College student taking a course

Goal: maximize learning by end of semester

College student taking a course

Goal: maximize learning by end of semester

Question: how should the student organize studying?

College student taking a course

Goal: maximize learning by end of semester

Question: how should the student organize studying?

Occasional feedback: quiz grades, exam grade etc.

# Reinforcement Learning



state =  $s_t$

obs (partial state) =  $o_t$

# Reinforcement Learning



ML Model

Policy

$$\pi_{\theta}(a_t|o_t)$$

probability distribution over actions given obs

state =  $s_t$

obs (partial state) =  $o_t$

# Reinforcement Learning



state =  $s_t$   
obs (partial state) =  $o_t$

ML Model

Policy  
 $\pi_\theta(a_t|o_t)$   
probability distribution over actions given obs



a<sub>t</sub>(sample action)

# Reinforcement Learning



state =  $s_t$   
obs (partial state) =  $o_t$

ML Model

Policy

$$\pi_\theta(a_t|o_t)$$

probability distribution over actions given obs

$a_t$  (sample action)

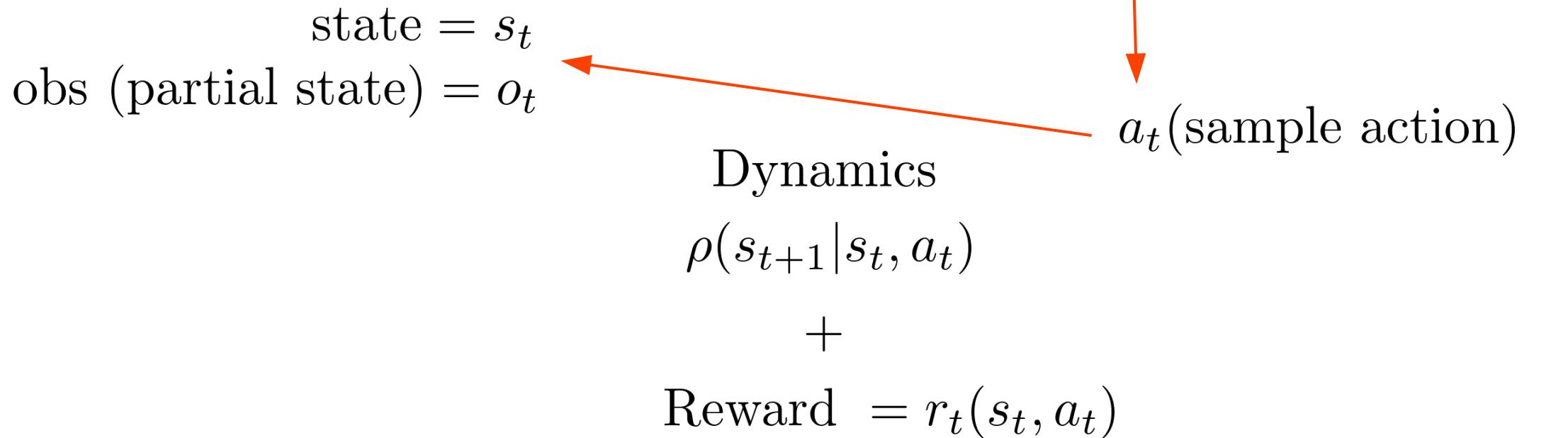
Dynamics

$$\rho(s_{t+1}|s_t, a_t)$$

+

$$\text{Reward} = r_t(s_t, a_t)$$

# Reinforcement Learning



# Reinforcement Learning

Dynamics are usually not known

Dynamics,  $\rho$  and  $r_t$  stochastic

$R = \sum_{t=1}^T r_t$  highly path dependent

More precisely, maximize:  $\mathbb{E}[R|s_0]$

## **Supervised Learning:**

**Optimization:** Turn learning problem into optimization problem

$$\min_x f(x) \rightarrow \min_{\text{model parameters}} \text{error}(\text{pred}, \text{output})$$

**Generalization:** should be able to predict accurately on unseen inputs

## Planning:

**Optimization:** Turn learning problem into optimization problem

$$\min_x f(x) \rightarrow \min_{\text{model parameters}} \text{error}(\text{pred}, \text{output})$$

**Generalization:** should be able to predict accurately on unseen inputs

**Delayed Feedback:** don't know if sequence of moves good/bad till experiment ends

## **Reinforcement Learning:**

**Optimization:** Turn learning problem into optimization problem

$$\min_x f(x) \rightarrow \min_{\text{model parameters}} \text{error}(\text{pred}, \text{output})$$

**Generalization:** should be able to predict accurately on unseen inputs

**Delayed Feedback:** don't know if sequence of moves good/bad till experiment ends

**Exploration:** Rules (dynamics) of system **unknown**. Have to discover them by trial-and-error.

Supervised = have **feedback** (output) for **every example** (input)

Reinforcement Learning = have **infrequent** feedback

What if no feedback?!

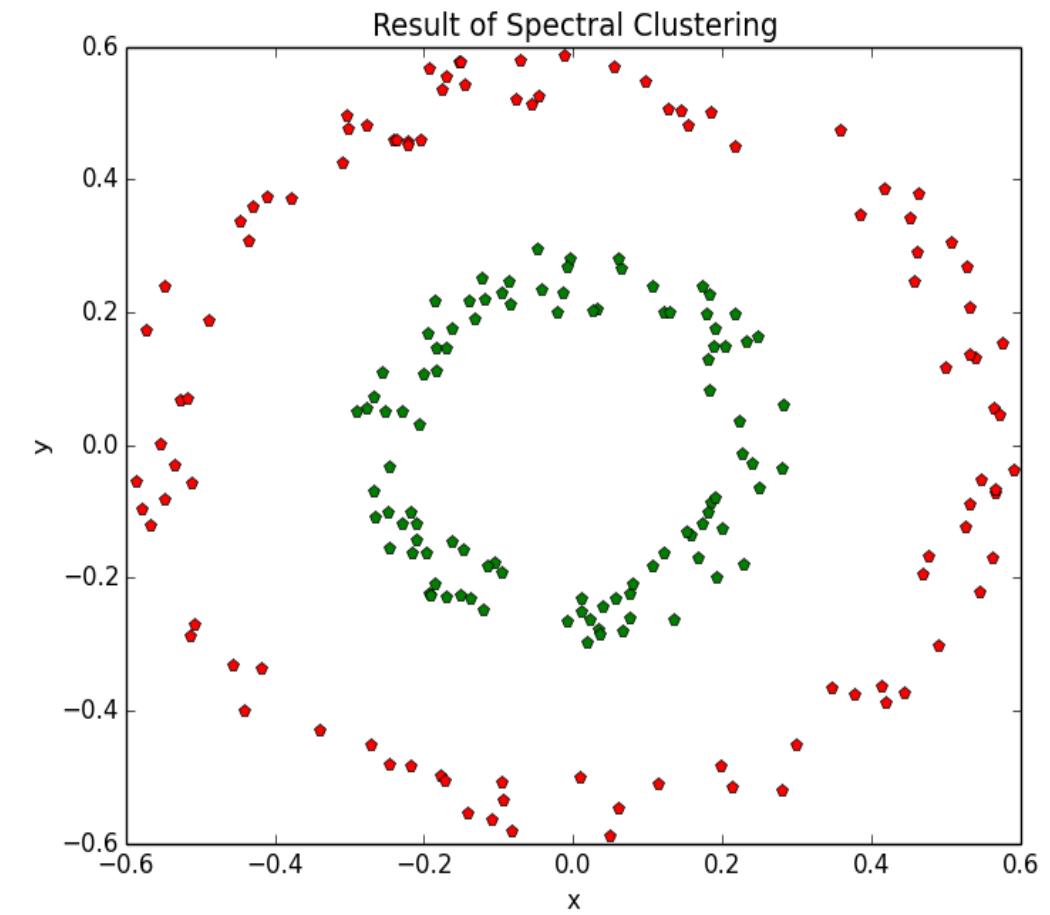
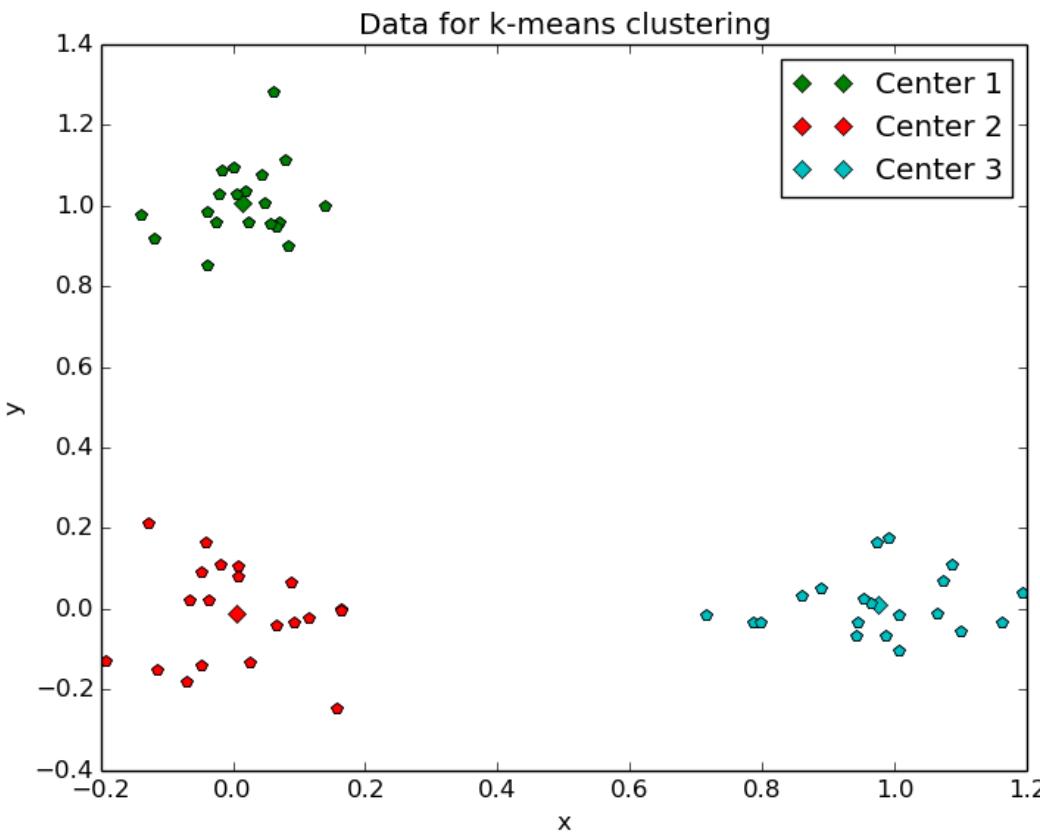
Supervised = have **feedback** (output) for **every example** (input)

Reinforcement Learning = have **infrequent** feedback

What if no feedback?!

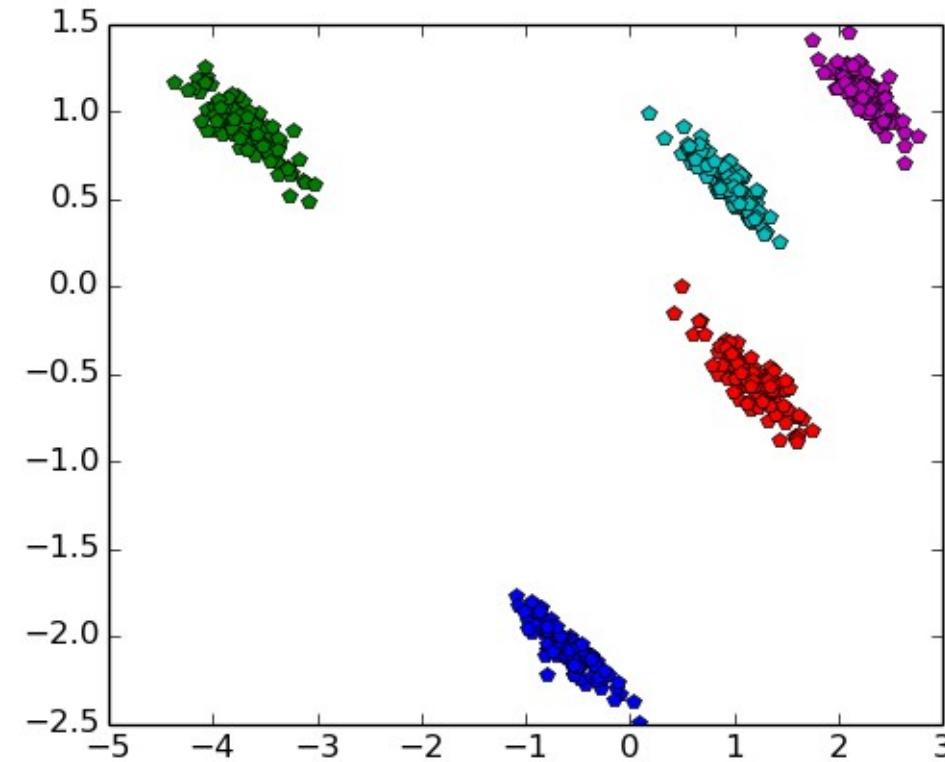
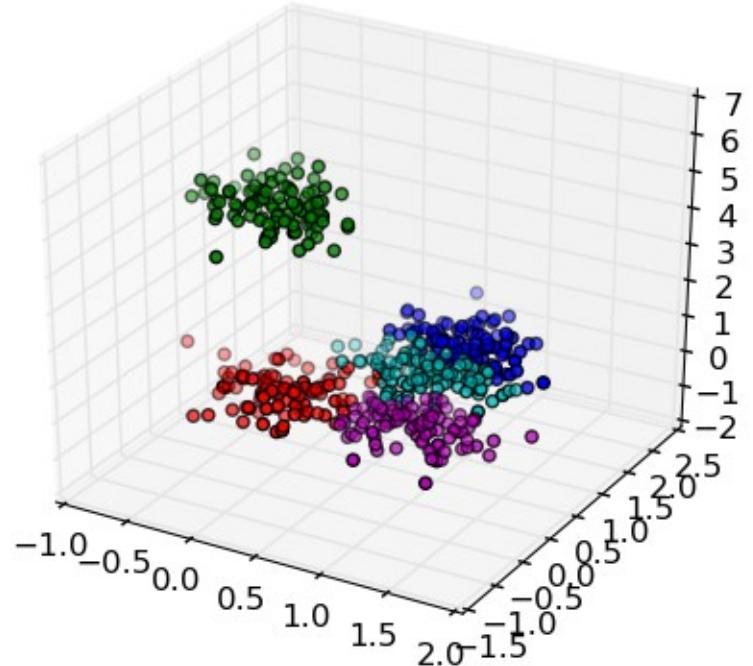
**Very common!** = just have some data

# Unsupervised Learning



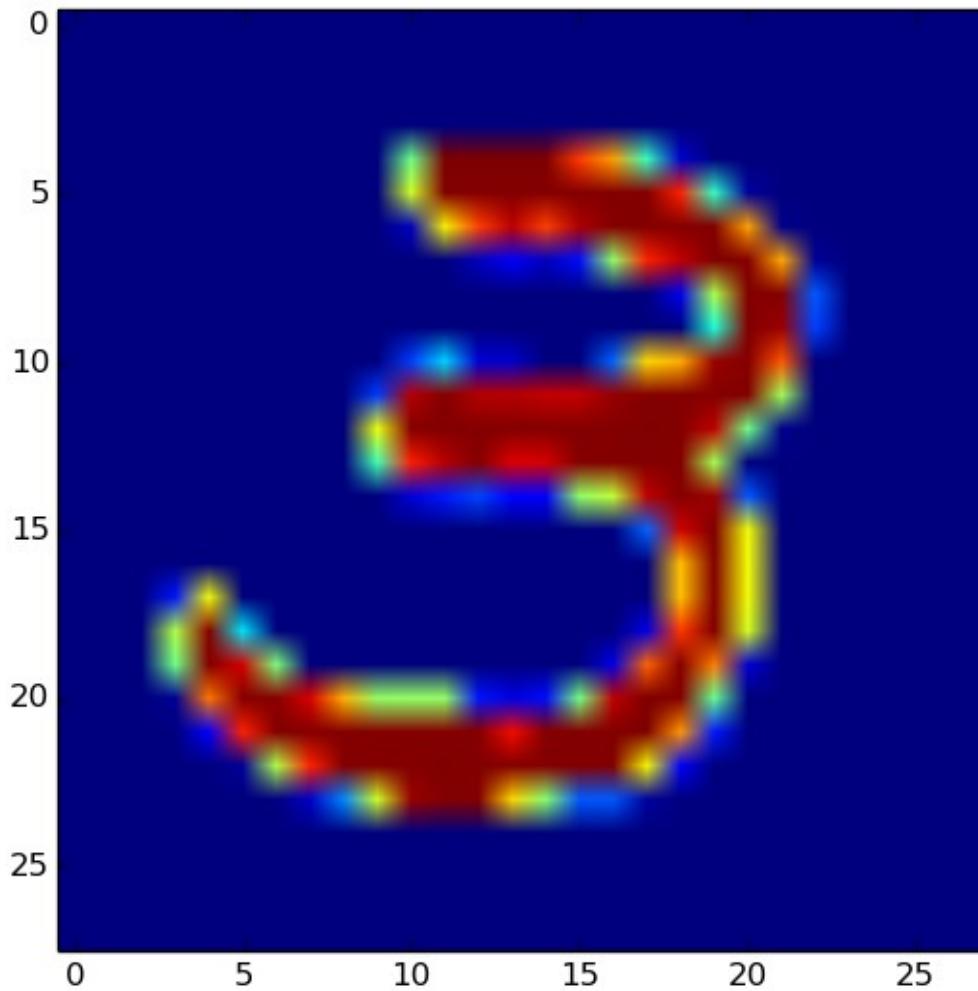
**Clustering → assign point/row to one of K groups**

# Unsupervised Learning



## Dimensionality Reduction

# Unsupervised Learning

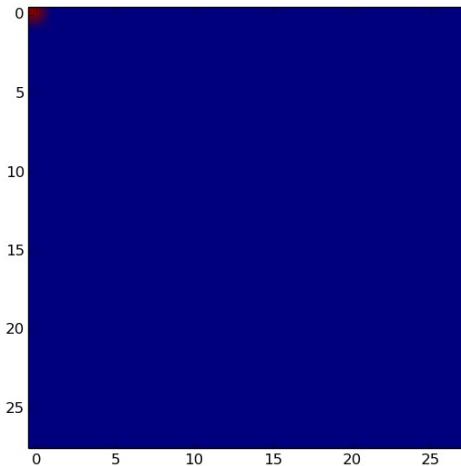


**Dimensionality Reduction**

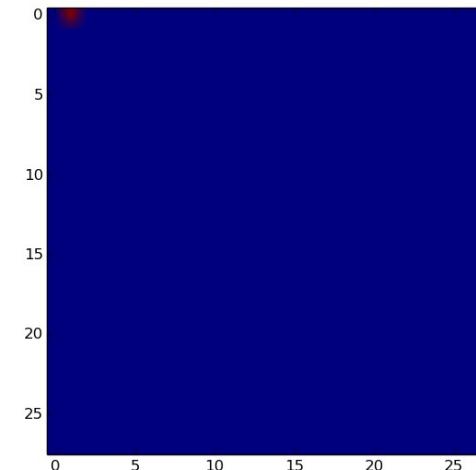
# Unsupervised Learning

Old Basis:

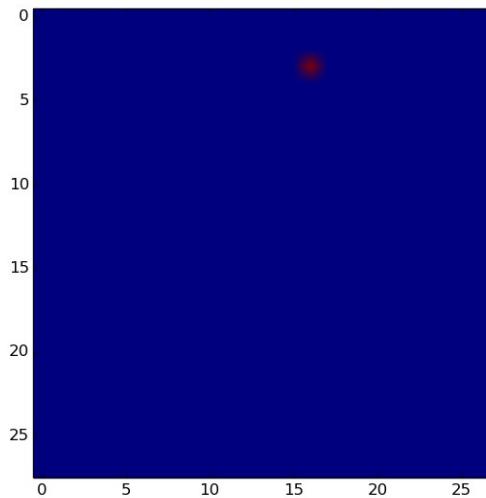
Each image =  $a_1$



$+ a_2$



$+ \dots + a_{100}$

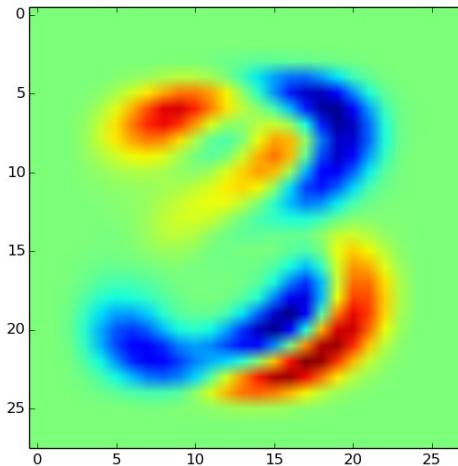


$+ \dots$

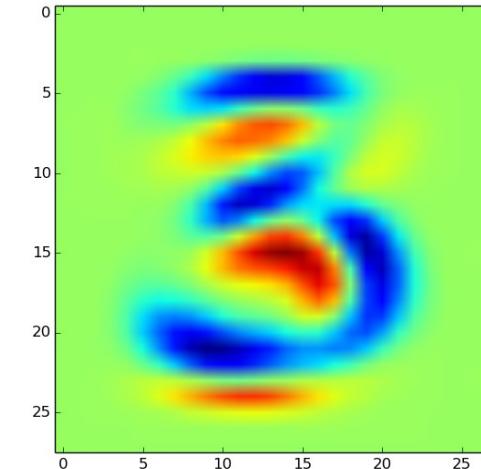
# Unsupervised Learning

New Basis:

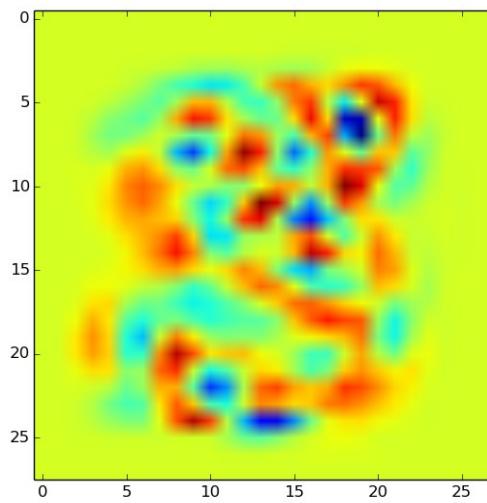
Each image =  $b_1$



$+ b_2$



$+ \dots + b_{100}$



$+ \dots$

# Unsupervised Learning

Many more

In some ways, this is the dominant problem

Generally:

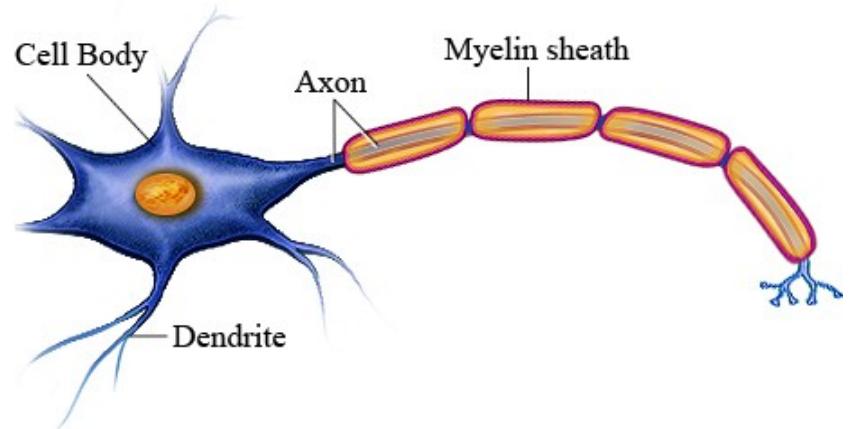
Supervised = straightforward if problem defined. Good implementations available

Unsupervised = Good implementations available often. Problem definition is much harder

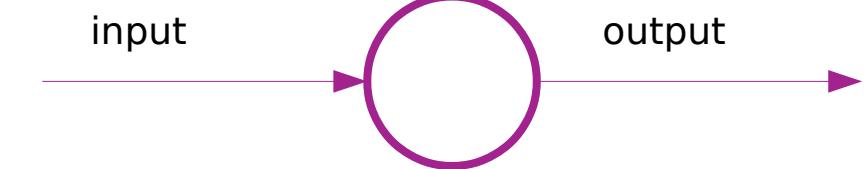
Reinforcement Learning = models can be unstable and very sensitive to details (rewards, dynamics, cardinality of state and action spaces ...)

# Neural Networks

# What is a neural network?

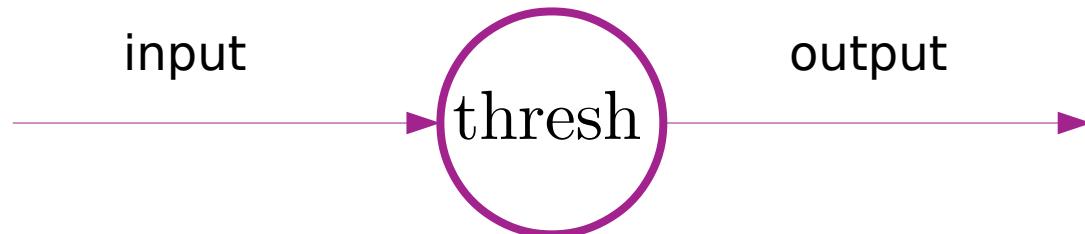


© Healthwise, Incorporated



**Loosely** inspired by neurons in the brain

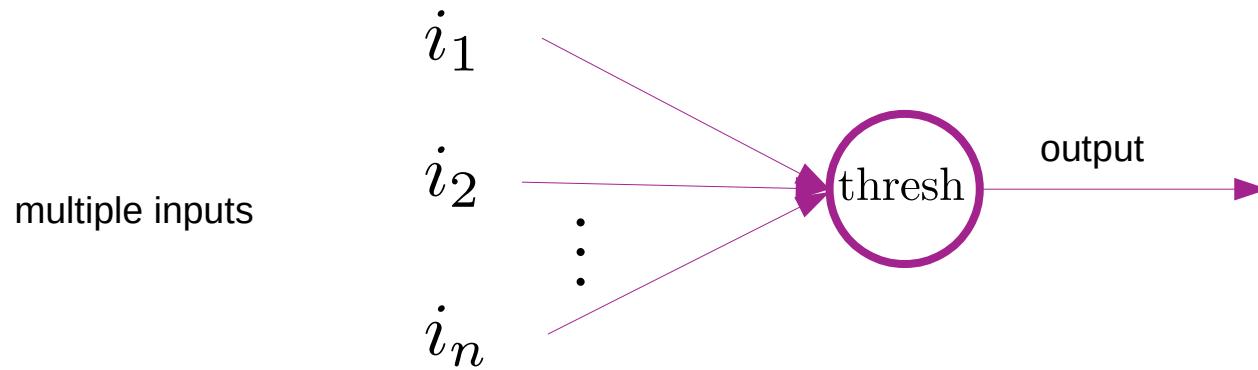
# What is a neural network?



$\text{input} > \text{threshold} \rightarrow \text{output} = 1$

$\text{input} \leq \text{threshold} \rightarrow \text{output} = 0$

# What is a neural network?



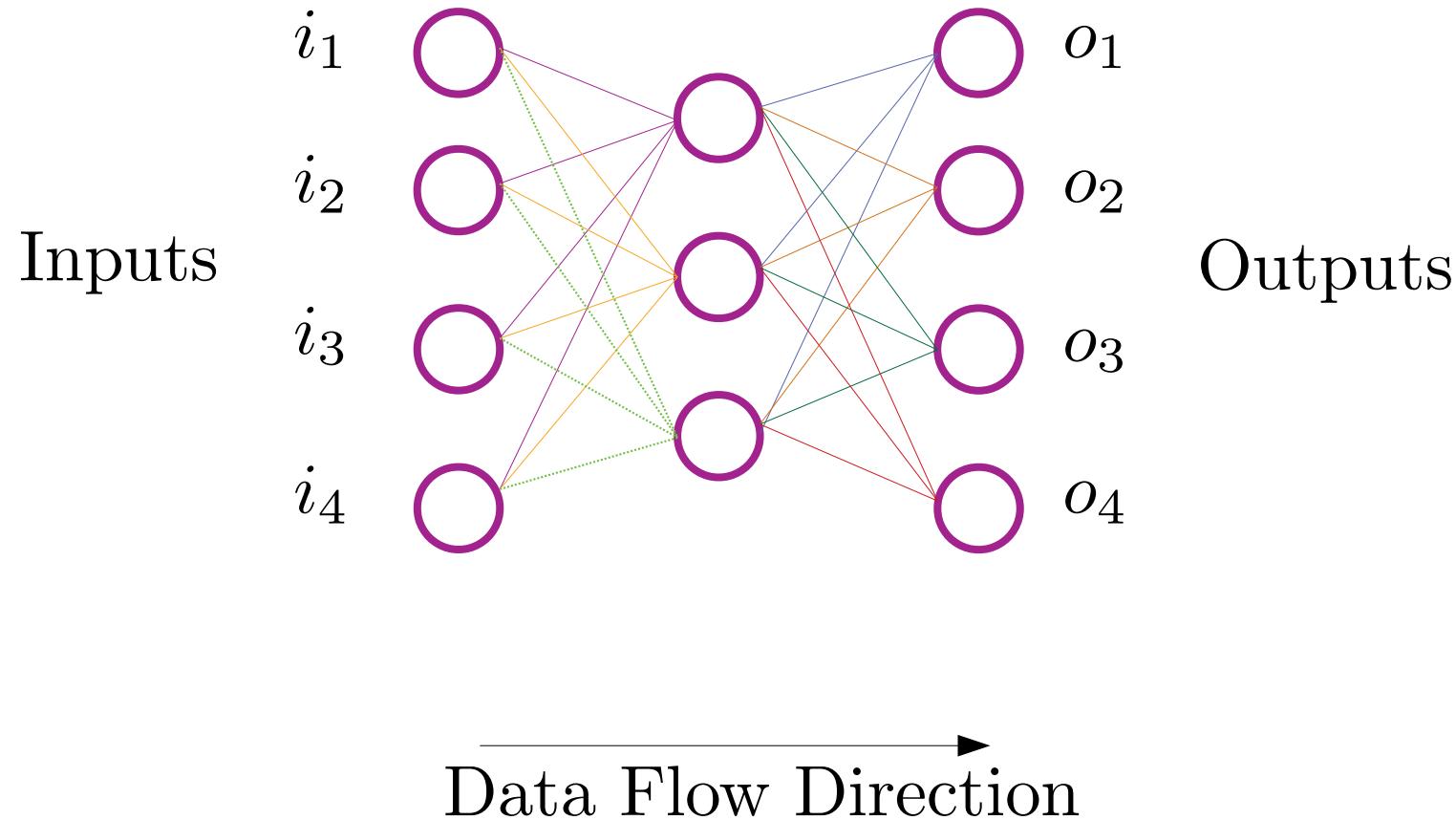
$$\text{input, } i = \underbrace{w_1}_{\text{weight}} i_1 + \underbrace{w_2}_{\text{weight}} i_2 + \dots + \underbrace{w_n}_{\text{weight}} i_n$$

$\text{input} > \text{threshold} \rightarrow \text{output} = 1$

$\text{input} \leq \text{threshold} \rightarrow \text{output} = 0$

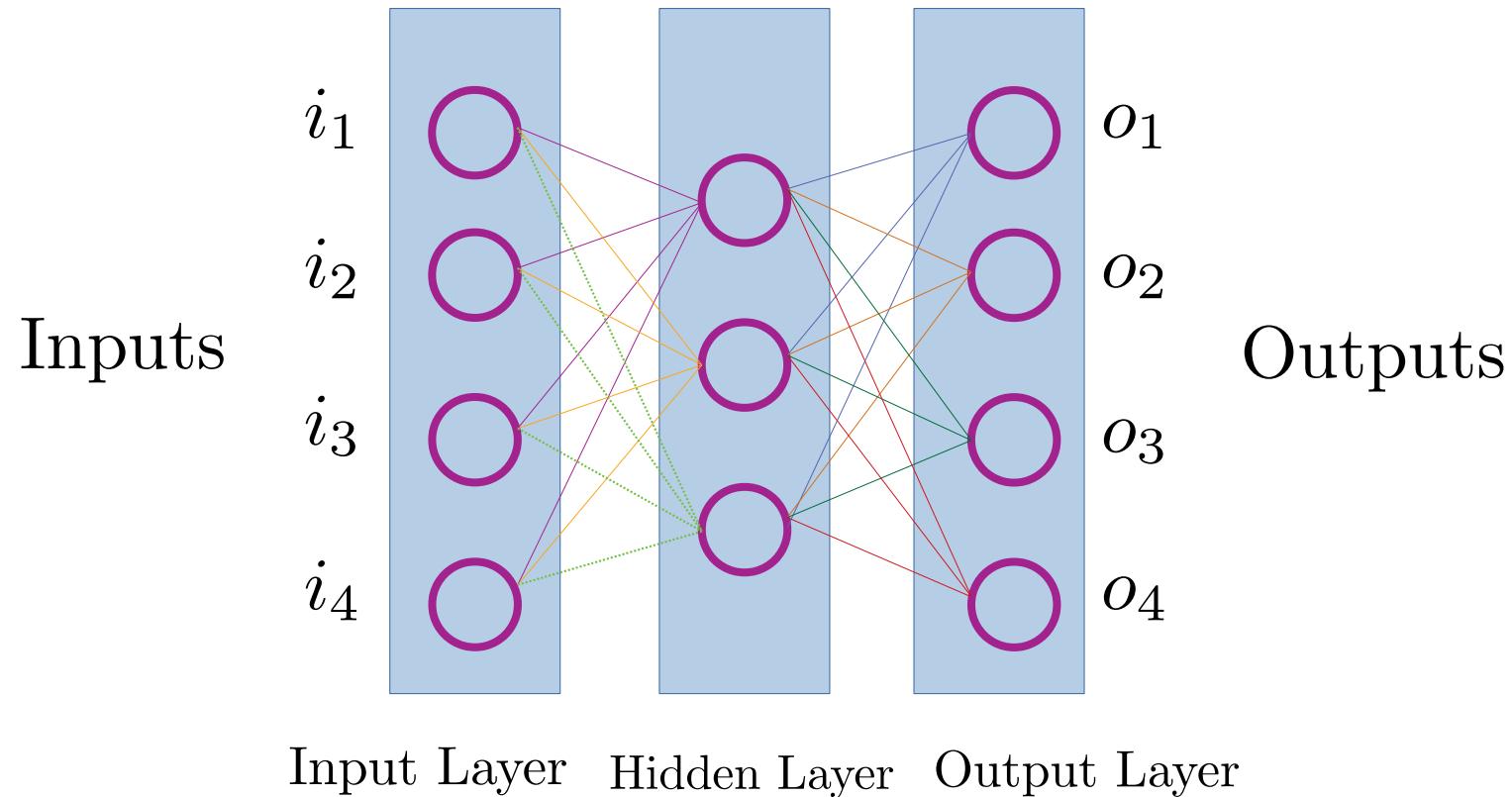
# What is a neural network?

Combine basic neurons into a network



# What is a neural network?

Feed-forward architecture



# Why go through all this trouble?

Why is any of this useful?

Can't I just use linear regression please?

Or random forests?

# **Why go through all this trouble: Universal Approximation Theorem**

## **Loose version**

A (feed-forward) neural network with one hidden layer can approximate any “reasonable” function,  $f$ , to arbitrary accuracy

# Why go through all this trouble: Universal Approximation Theorem

## Precise version

“Activation” function	$\sigma : \mathbb{R} \rightarrow \mathbb{R}$ $\sigma$ non-constant, bounded, continuous
Function to learn	$f : [0, 1]^n \rightarrow \mathbb{R}$ $f$ continuous on $[0, 1]^n$

# Why go through all this trouble: Universal Approximation Theorem

Precise version

$$\sigma : \mathbb{R} \rightarrow \mathbb{R} \qquad f : [0, 1]^n \rightarrow \mathbb{R}$$

can find  $w_i, u_i, b_i, N$

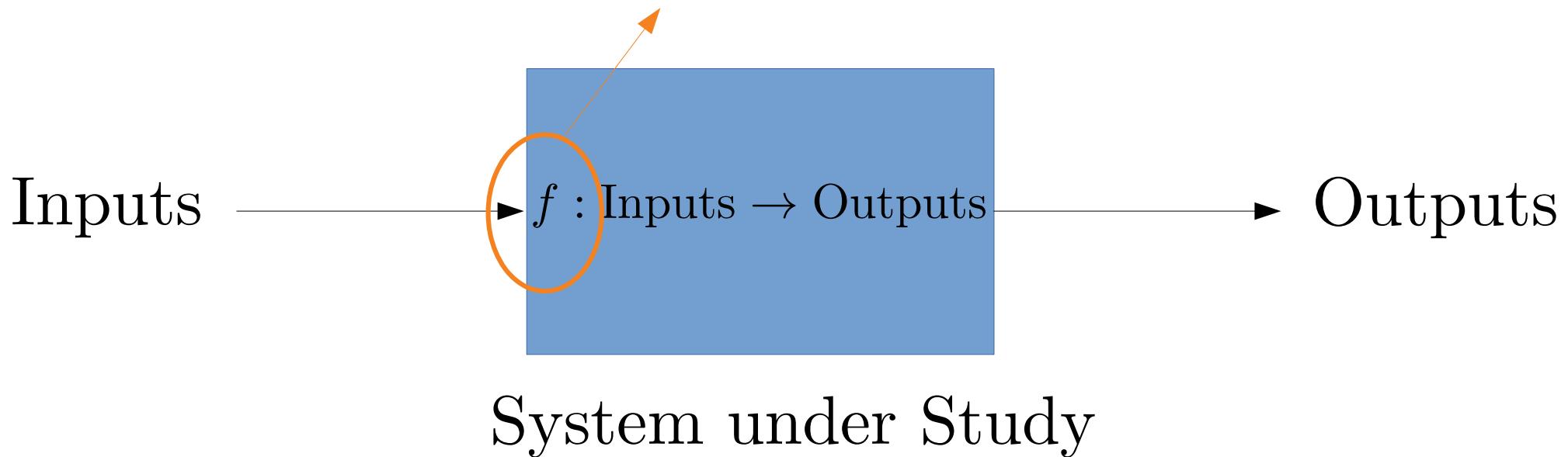
$$| \underbrace{f(x)}_{\text{function to learn}} - \underbrace{\sum_{i=1}^N w_i \sigma(u_i^T x + b_i)}_{\text{approximation to } f} | < \epsilon$$

Important Technicality: can extend to any compact subset for the domain

# Why go through all this trouble: Universal Approximation Theorem

**What does this mean for me?**

Can approximate to arbitrarily high accuracy with a neural network



# Why go through all this trouble: Universal Approximation Theorem

## Two caveats

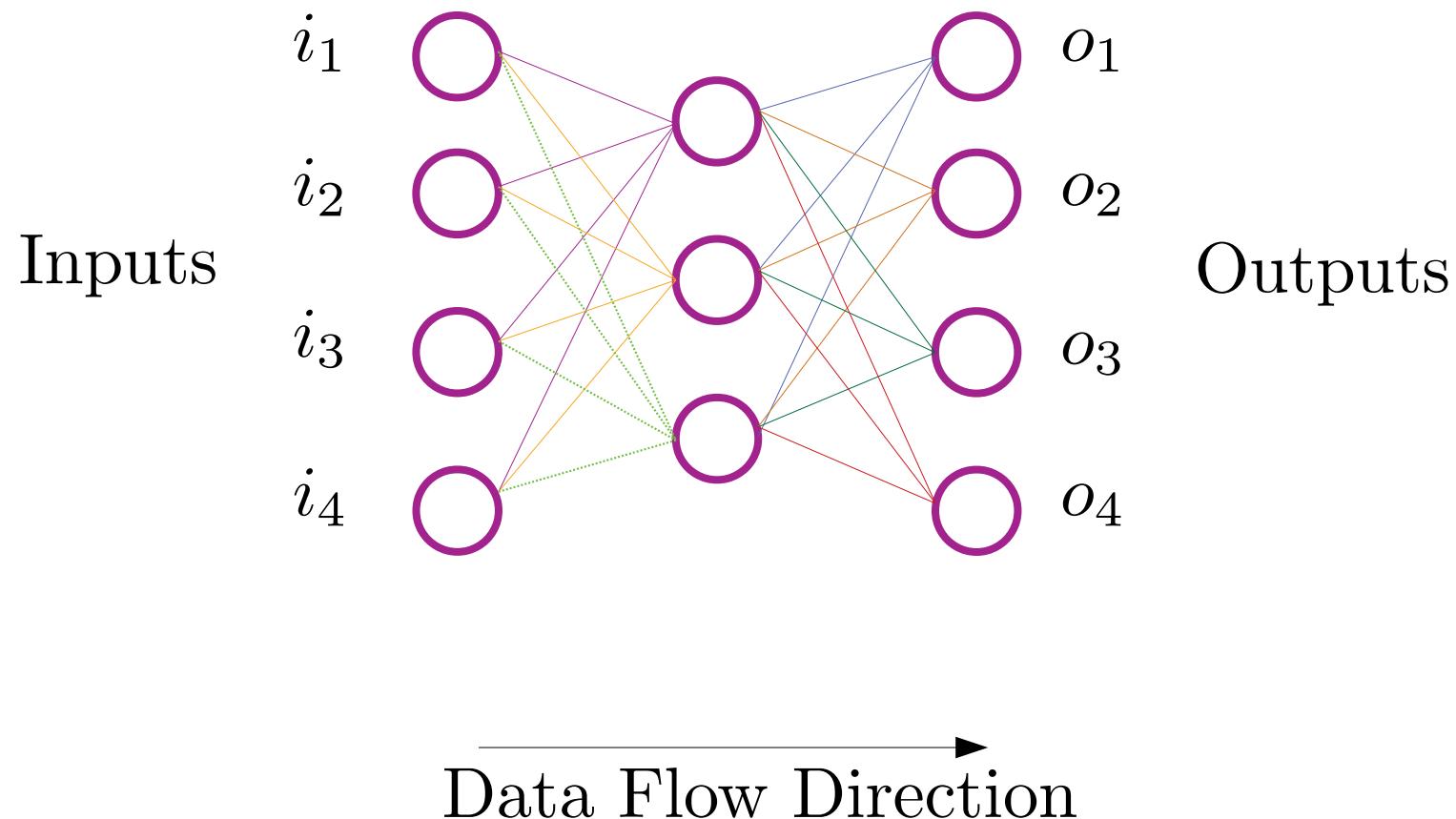
Need non-linear activation function,  $\phi$  in our neural network  
(we'll see what this means shortly)

$N$  refers to the number of neurons (nodes) in  
the hidden layer. This grows exponentially as  $\epsilon$  decreases

$$| \underbrace{f(x)}_{\text{function to learn}} - \underbrace{\sum_{i=1}^N w_i \sigma(u_i^T x + b_i)}_{\text{approximation to } f} | < \epsilon$$

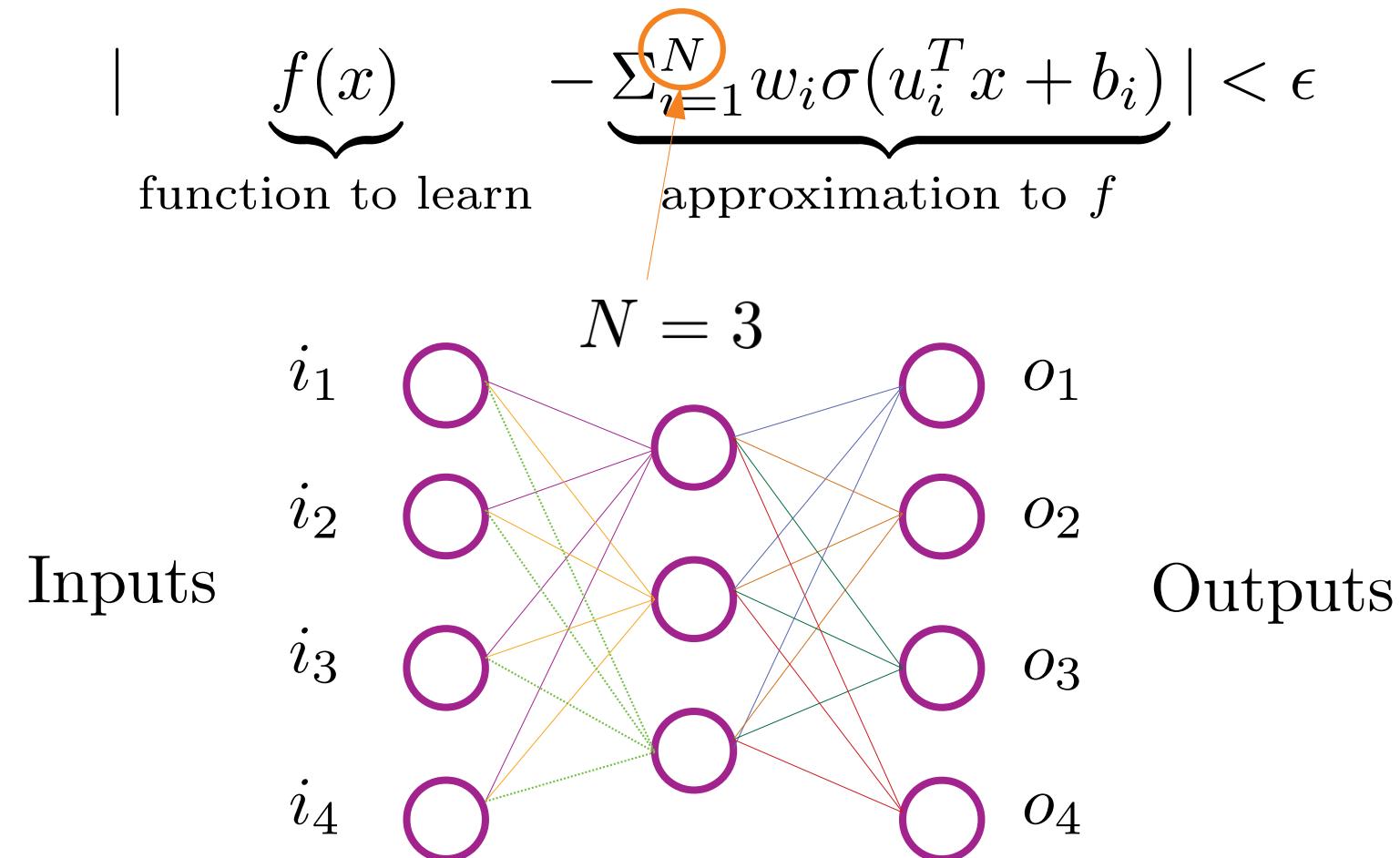
# Two Caveats

Need non-linear activation function,  $\sigma$  in our neural network  
(we'll see what this means shortly)

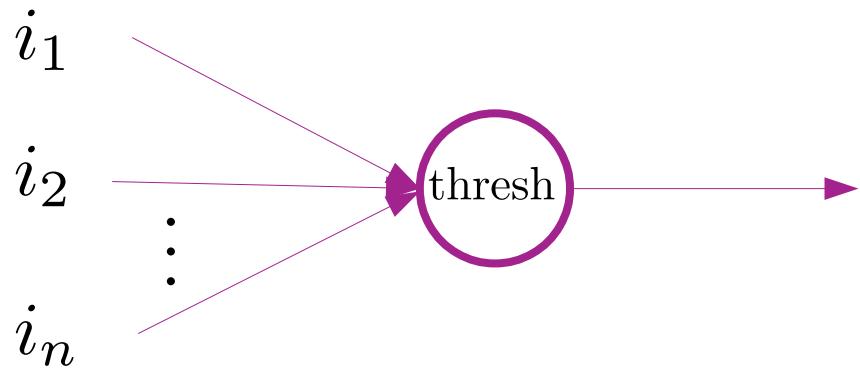


# Two Caveats

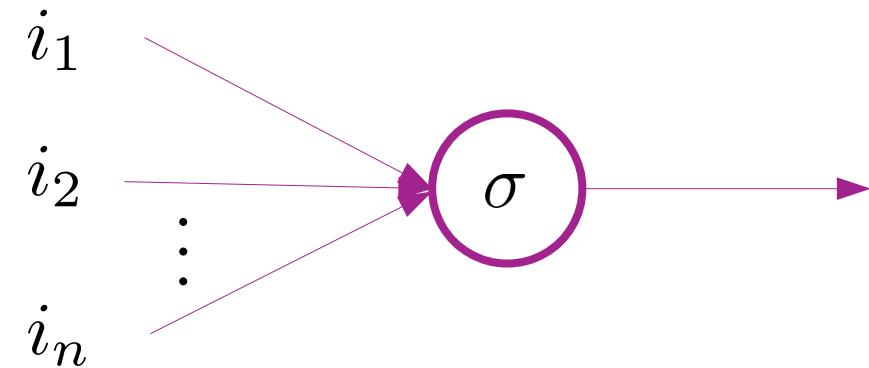
$N$  refers to the number of neurons (nodes) in the hidden layer. This grows exponentially as  $\epsilon$  decreases



# Structure of a node (neuron)



$$\text{thresh}(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$



$\sigma$  can be one of many  
non-linear  
activation functions

# Linear vs Non-linear

Linear:  $\sigma(x) = C_1x + C_2$

Applying one function after the other

Composition of two linear functions is linear

$$f(x) = C_1x + C_2 \quad g(x) = D_1x + D_2$$

$$g(f(x)) = D_1f(x) + D_2 = D_1(C_1x + C_2) + D_2$$

$$g(f(x)) = \underbrace{(D_1C_1)}_{E_1}x + \underbrace{(D_1C_2 + D_2)}_{E_2}$$

# Linear vs Non-linear

Non-linear

$$\log(x)$$

$$e^x$$

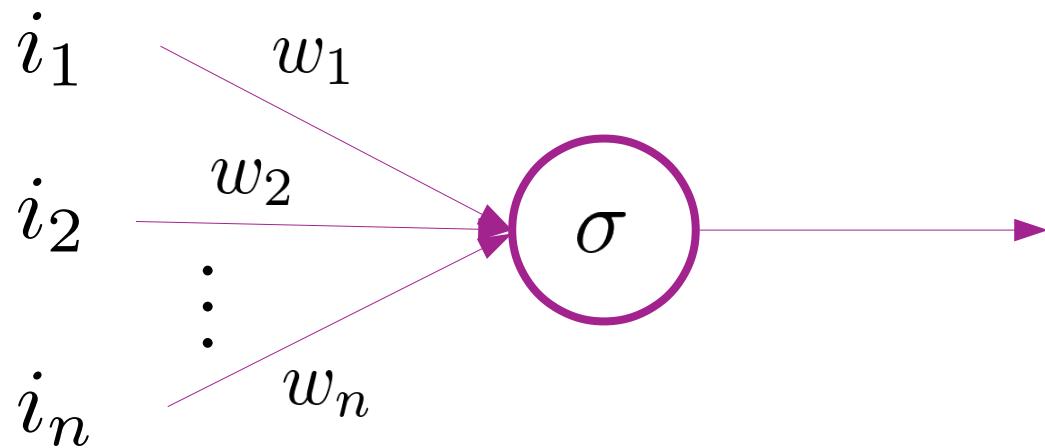
$$x^2$$

$$\sqrt{x}$$

infinitely more

Most complex systems behave non-linearly

# Structure of a node (neuron)



$$o = \sigma(w_1 i_1 + w_2 i_2 + \dots + w_n i_n + b)$$

Short-hand:  $o = \sigma(w^T i + b)$

$$i = \begin{bmatrix} i_1 \\ i_2 \\ \vdots \\ i_n \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad w^T = [w_1 \quad w_2 \quad \dots \quad w_n]$$

# Activation Functions: Requirements

Non-linear

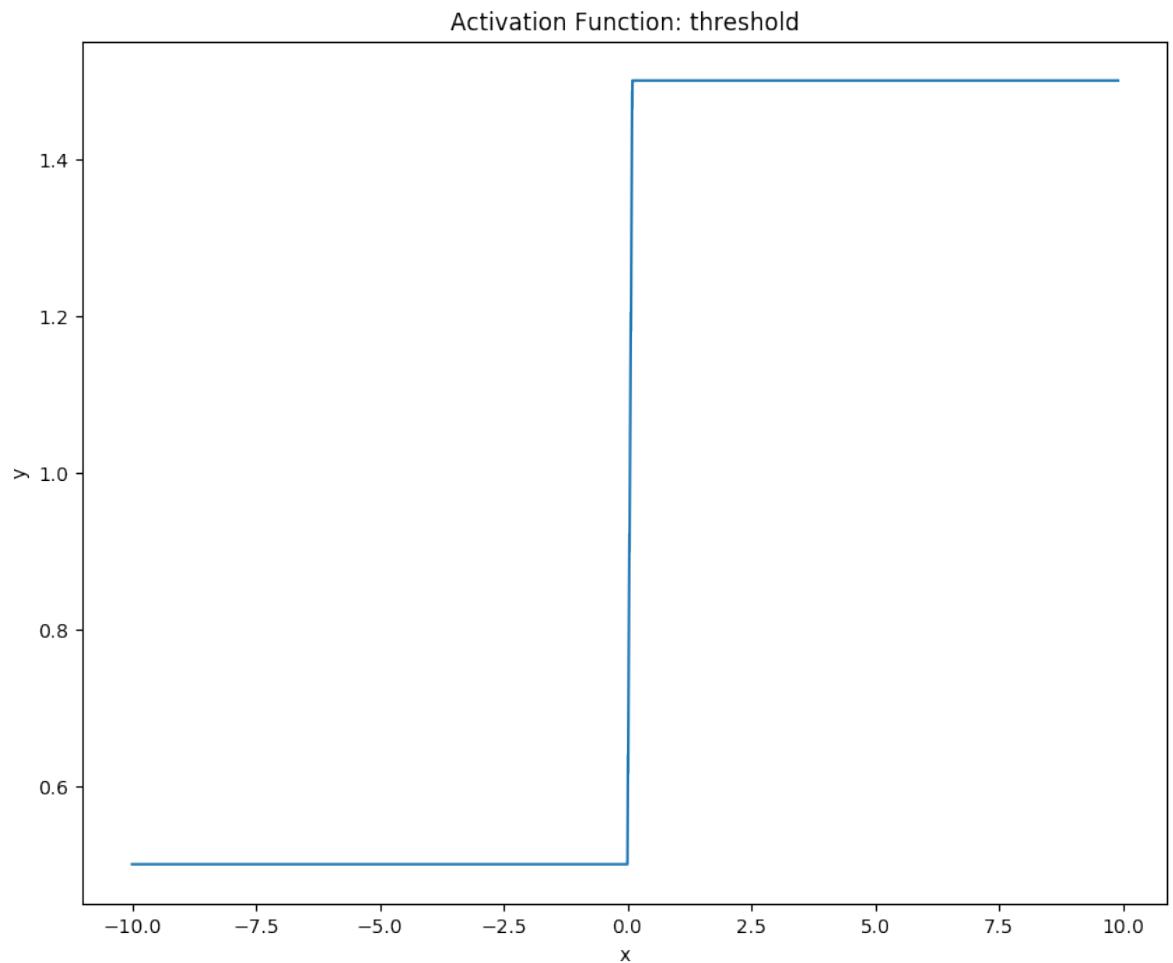
Simple to evaluate

Will need derivatives so hopefully derivatives easy to evaluate

# Binary switch

Derivative not defined at  $x = 0$   
(not a major issue)

Derivative not informative  
 $\sigma'(x) = 0$  everywhere except at  $x=0$

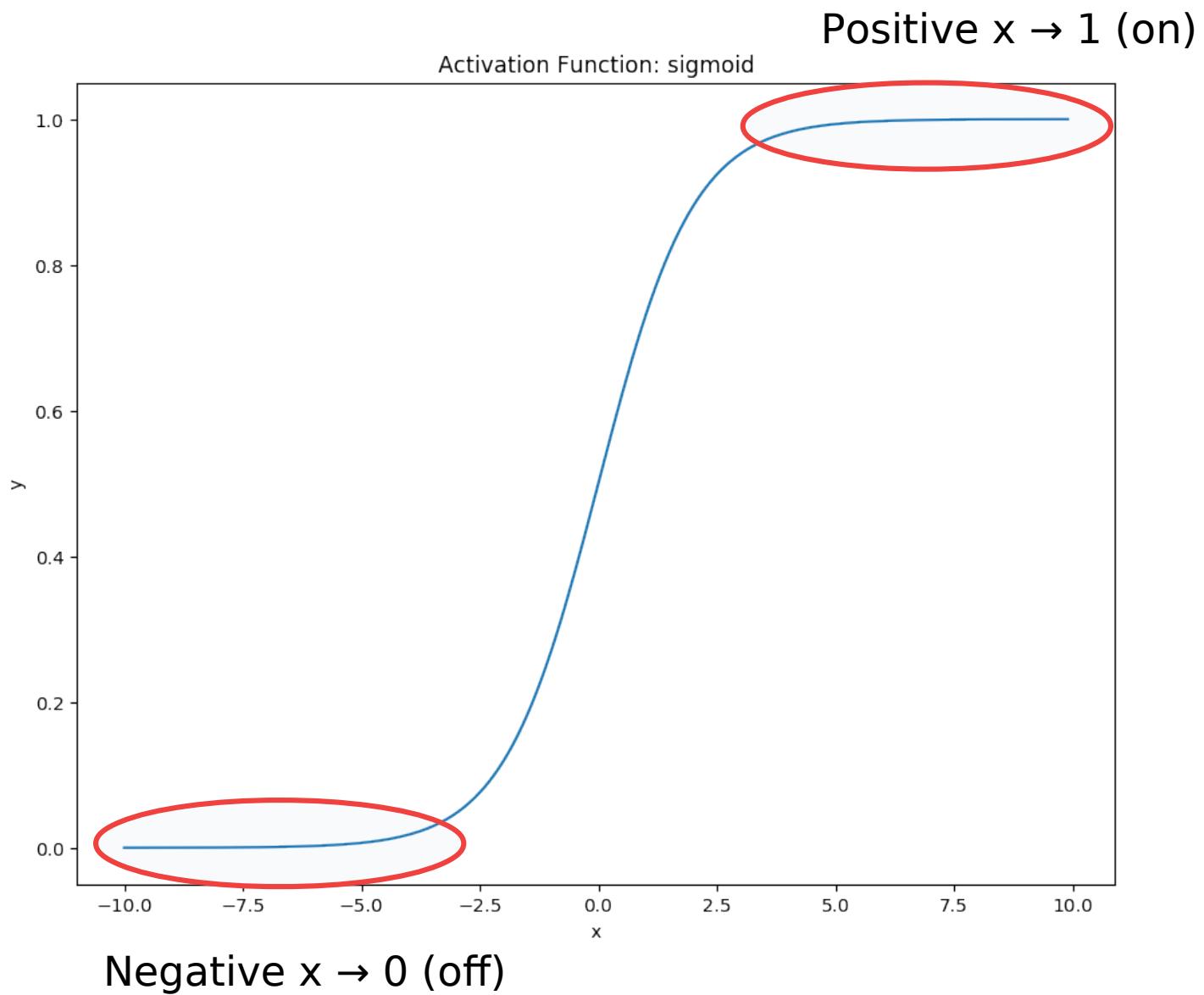


$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

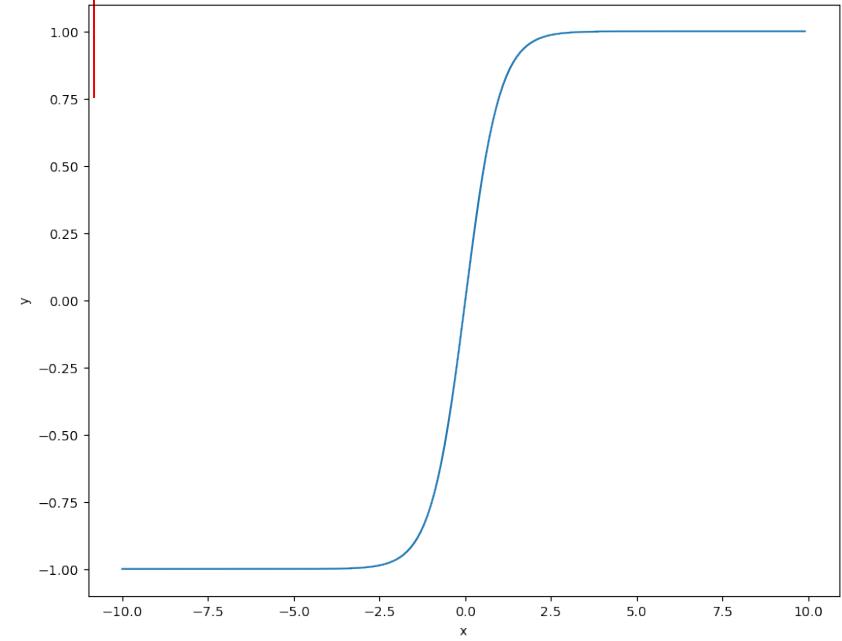
Continuous switch

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

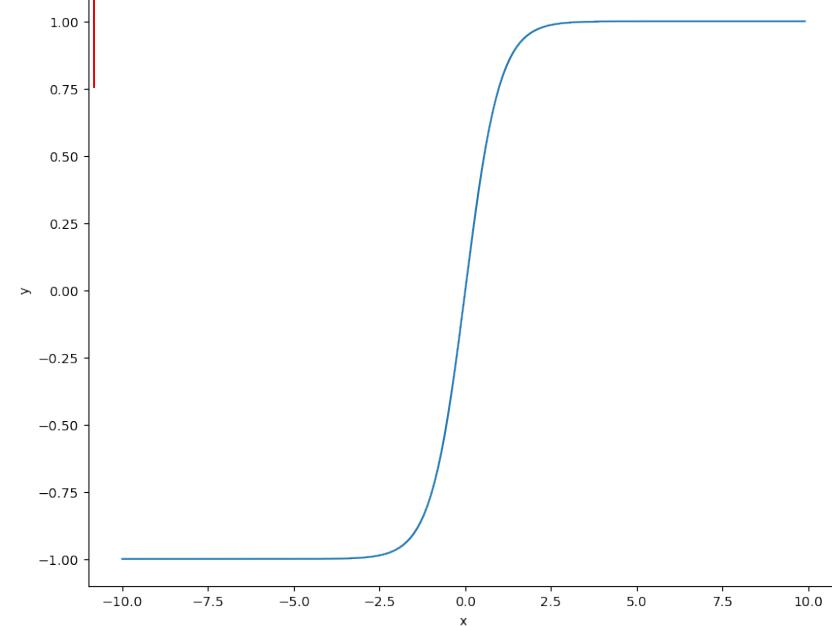
Derivatives in terms of activation itself



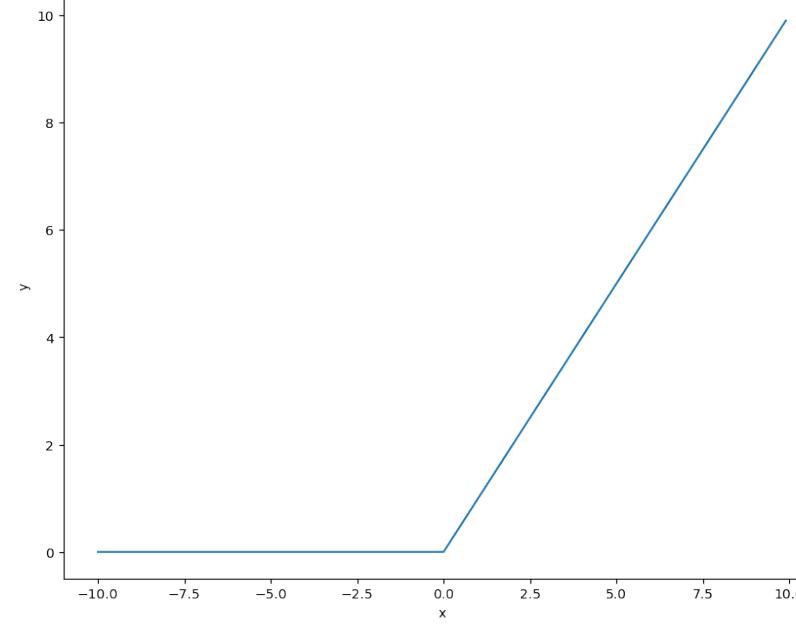
Activation Function: tanh



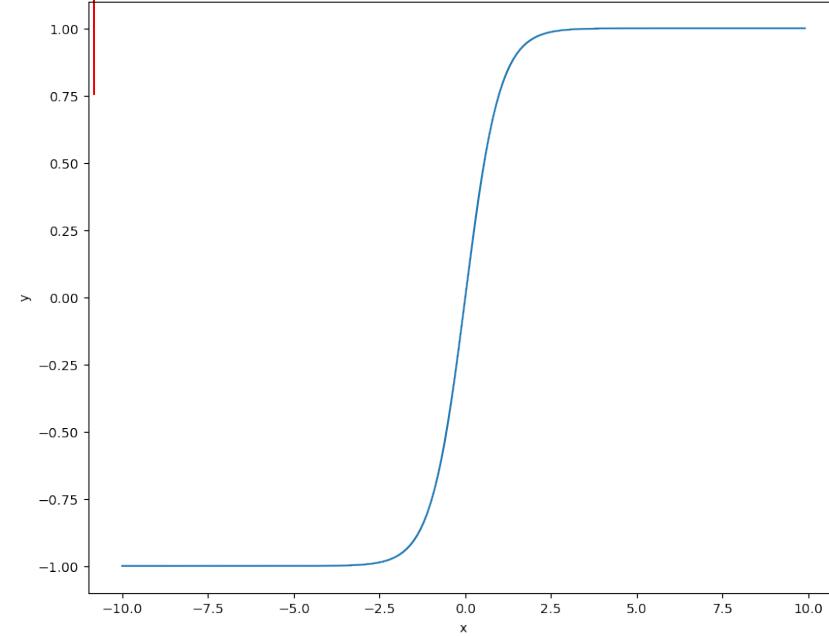
Activation Function: tanh



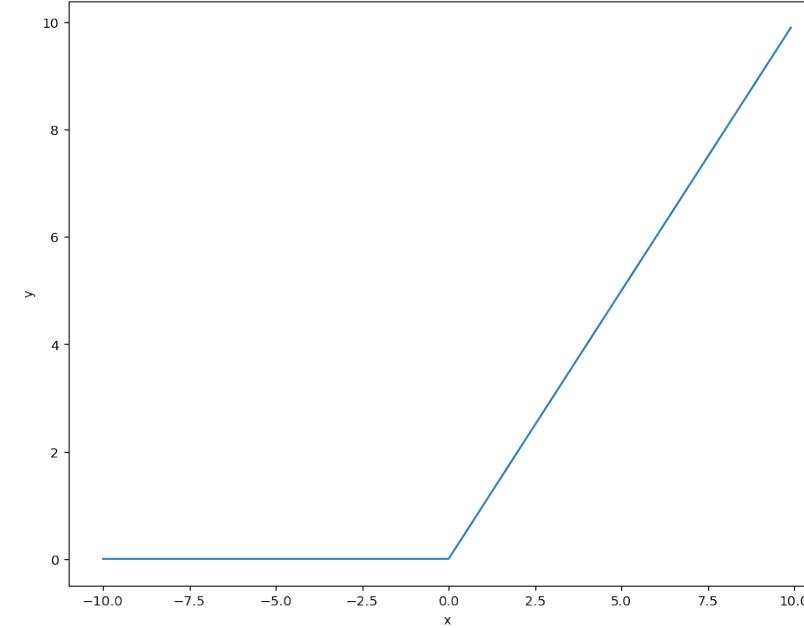
Activation Function: relu



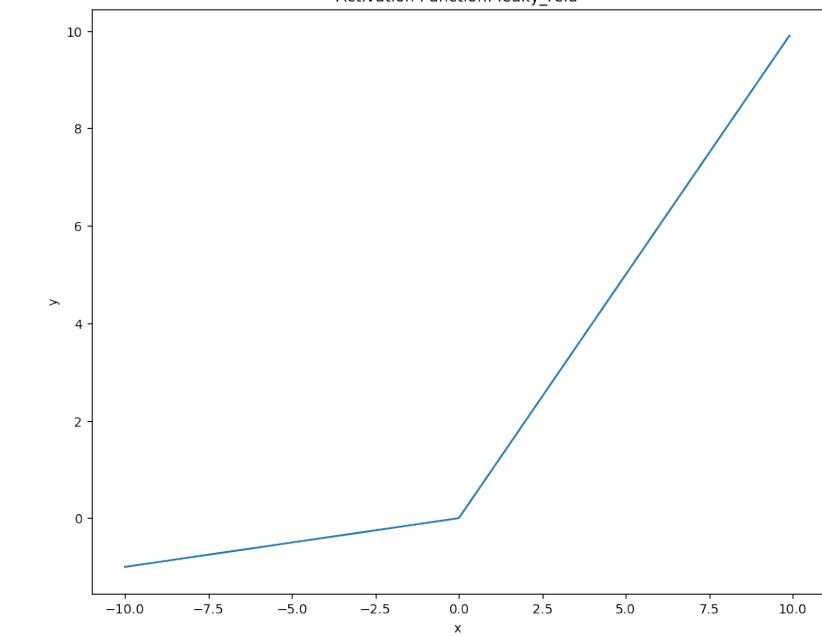
Activation Function: tanh



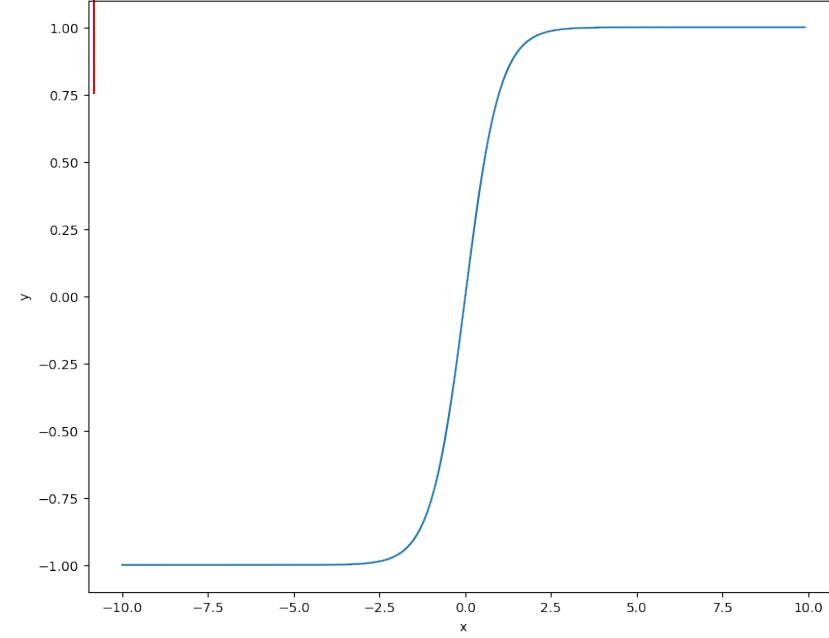
Activation Function: relu



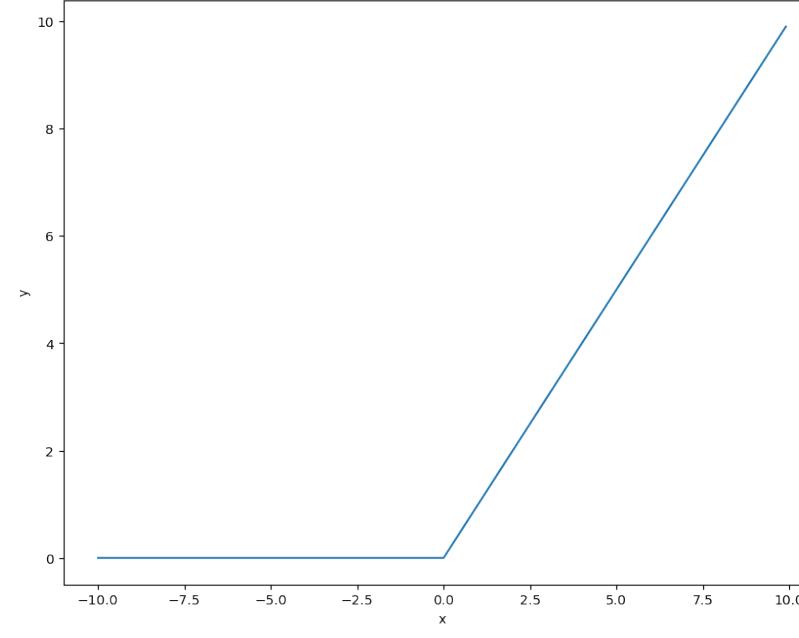
Activation Function: leaky\_relu



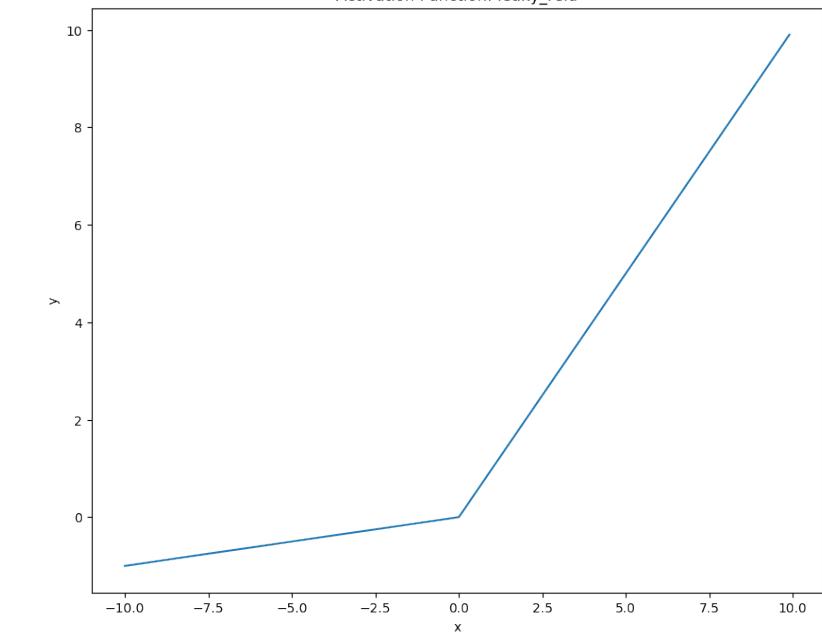
Activation Function: tanh



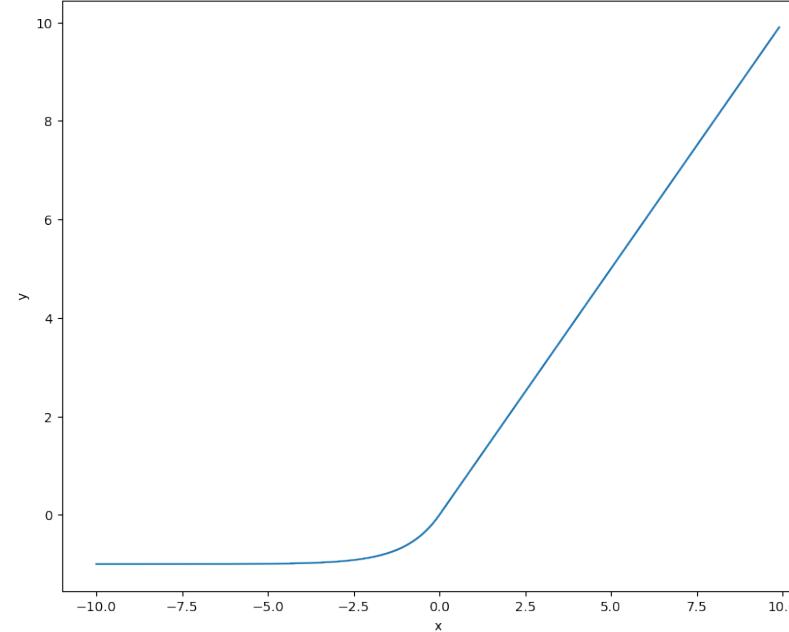
Activation Function: relu



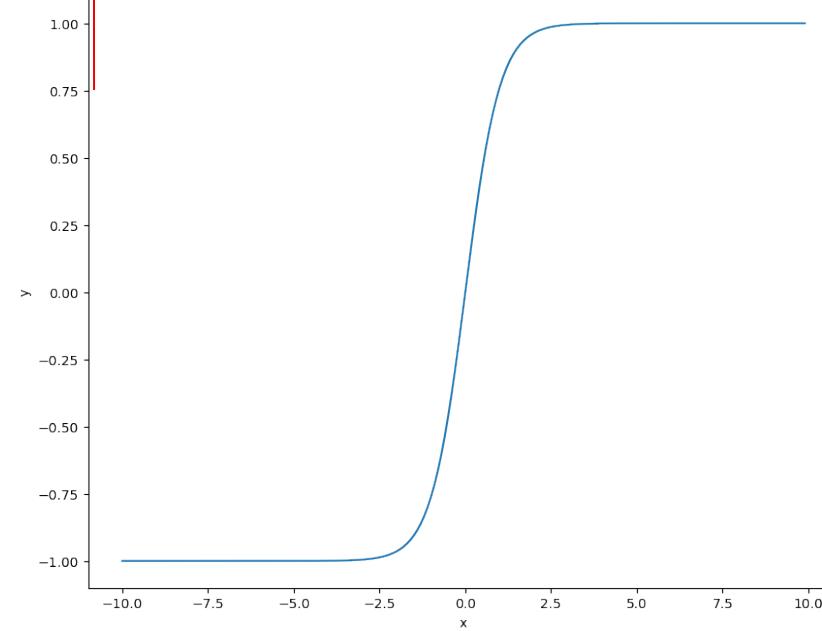
Activation Function: leaky\_relu



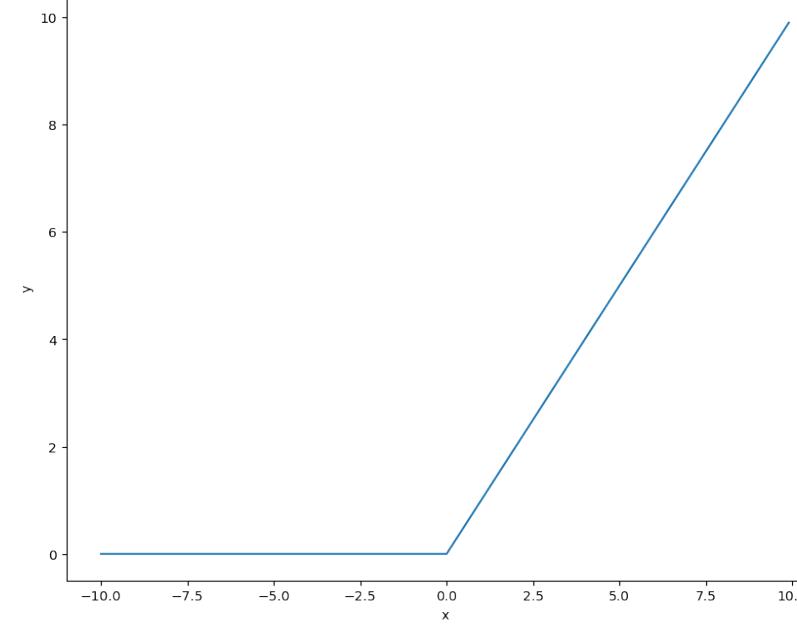
Activation Function: exponential\_leaky\_relu



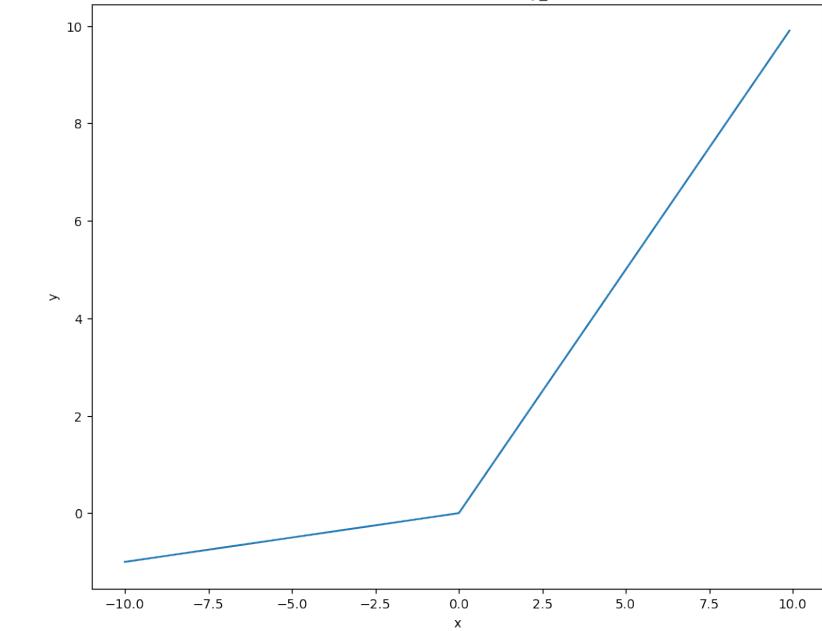
Activation Function: tanh



Activation Function: relu



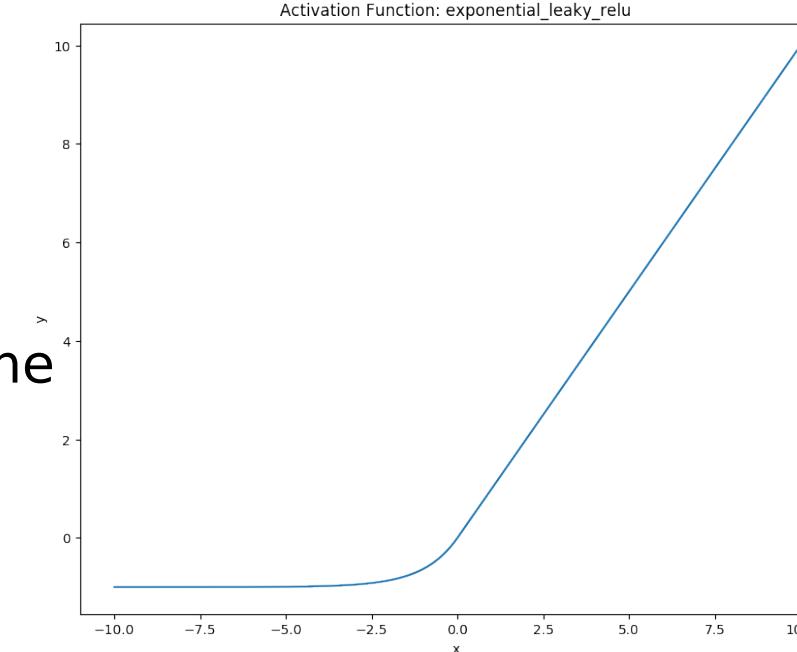
Activation Function: leaky\_relu



Many more choices

Choices dictated by:

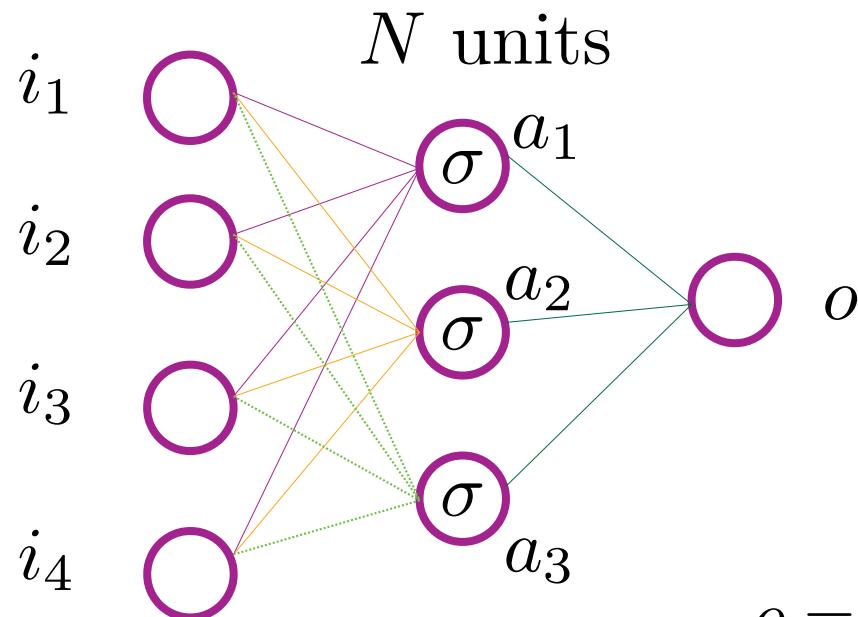
- Issues during training time (saturation of nodes)
- Type of predictions



# Recap

- Start with simple neuron
- Need to add non-linear activation to be able to learn/mimic any reasonable function
- Looked at a bunch of activations and added to our neuron
- Time to put neurons together into a network and see how prediction would work.

# Full Neural Network



Input:  $: (i_1, i_2, i_3, i_4)$

Hidden:

$$a_1 = \sigma(w_1^T i + b_1)$$

$$a_2 = \sigma(w_2^T i + b_2)$$

$$a_3 = \sigma(w_3^T i + b_3)$$

Output:

$$o = v_1 a_1 + v_2 a_2 + v_3 a_3 + c = v^T a + c$$

**Every edge has a “weight”**

# Full Neural Network

Input:  $i : (i_1, i_2, i_3, i_4)$

Hidden:

$$a_1 = \sigma(w_1^T i + b_1)$$

$$a_2 = \sigma(w_2^T i + b_2)$$

$$a_3 = \sigma(w_3^T i + b_3)$$

Output:

$$o = v_1 a_1 + v_2 a_2 + v_3 a_3 + c = v^T a + c$$

Concise matrix notation

$$o = V\sigma(Wi + b) + c$$

$W, V$  are matrices of weights

$b, c$  are vectors of biases

Imagine implementing this using a matrix library, say numpy in python.

# Full Training Cycle

Start with labeled data

elephant



dog



snow leopard



penguin

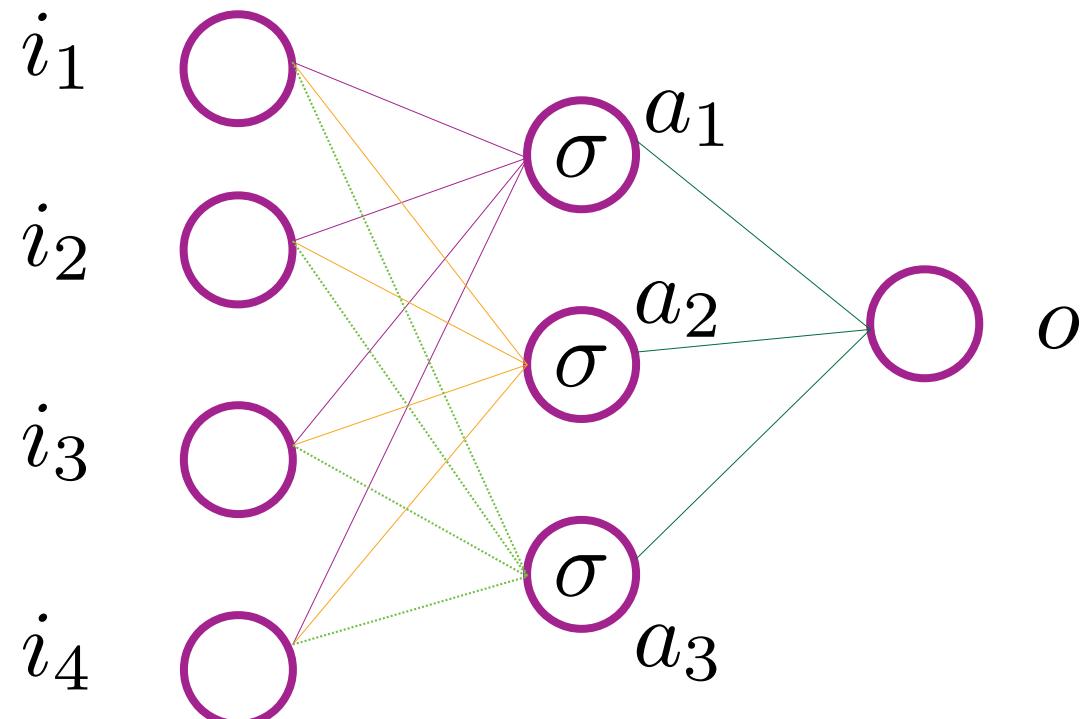


octopus



# Full Training Cycle

Pick network architecture (more details later)



# Full Training Cycle

## Make predictions (forward propagation)

$$o = V \sigma (W_i + b) + c$$

a number      weight activation    weight input      bias      bias

$V, W, b, c$  unknown

$V, W$  initialized randomly according to some distribution

$b, c$  initialized to vectors of zeros

# Full Training Cycle

**Make predictions (forward propagation)**

$$o = V \sigma (W_i + b) + c$$

a number      weight activation    weight input      bias      bias

**Compare output, o to actual value, v using a loss function  
(exactly like previous models)**

# Full Training Cycle

**Make predictions (forward propagation)**

$$o = V \sigma (W_i + b) + c$$

a number      weight activation      weight input      bias      bias

**Compare output, o to actual value, v using a loss function**

$$\text{Loss: } L = \frac{1}{2}(o - v)^2$$

$\rightarrow = 0$  if  $o = v$  (prediction matches value)  
 $\rightarrow$  larger the worse predictions get

**Loss function depends on your problem**

# Full Training Cycle

**Make predictions (forward propagation)**

$$o = V \sigma (W_i + b) + c$$

a number      weight activation    weight input      bias      bias

**Tweak  $V$ ,  $W$ ,  $b$  and  $c$  to adjust  $o$  and minimize  $L$**

Minimize:  $L = L(V, W, b, c)$

# Full Training Cycle



$v = \text{bunny}$

Predict (Forward propagation)

$$o = V\sigma(Wi + b) + c$$

$o = \text{polar bear}$

$$L(o, v)$$

Adjust  $V, W, b, c$

Minimize  $L$

# Loss Functions

# Loss Functions

Measuring deviation between predictions and targets/labels

Each example is treated independently:

$$\text{LOSS, } \mathcal{L} = \sum_{i=1}^n C[ \underbrace{\hat{y}_i}_{\text{prediction}}, \underbrace{y_i}_{\text{label}} ]$$

Cost for example i

# Loss Functions

$$\text{LOSS, } \mathcal{L} = \sum_{i=1}^n C[ \underbrace{\hat{y}_i}_{\text{prediction}}, \underbrace{y_i}_{\text{label}} ]$$

Cost for example i

Choice of  $C$  dictated by problem type

$C$  bounded by below i.e.  $C \geq 0$

# Loss Functions: Regression

Mean Squared Error

$$C[\hat{y}, y] = (\hat{y} - y)^2$$

Any deviation from  $y$  is punished

# Loss Functions: Binary Classification

Target: 0 or 1

Probability of belonging to class 1 :  $p \in [0, 1]$

How to measure deviation?

# Loss Functions: Binary Classification

Pick a threshold, say 0.5

Convert  $p$  to label:

$$p \leq 0.5 \rightarrow \text{prediction} = 0$$

$$p > 0.5 \rightarrow \text{prediction} = 1$$

Accuracy: % of predictions correct???

# Loss Functions: Binary Classification

Accuracy: % of predictions correct???

Problems?

# Loss Functions: Binary Classification

Accuracy: % of predictions correct???

Problems:

Dependent on threshold

Not a smooth function of predictions → indirect control of weights

# Loss Functions: Binary Classification

$$C = (\underbrace{y}_{\text{0 or 1}} - \underbrace{p}_{\in [0,1]})^2 \quad ???$$

Problems?

# Loss Functions: Binary Classification

$$C = (\underbrace{y}_{\text{0 or 1}} - \underbrace{p}_{\in [0,1]})^2 \quad ???$$

Unnatural to compare class label to probability

Penalty bounded above by 1 → too mild

# Loss Functions: Binary Classification

Given  $p, y$  : what's the **likelihood** of getting  $y$ ?

$$y = 0 : 1 - p$$

# Loss Functions: Binary Classification

Given  $p, y$  : what's the likelihood of getting  $y$ ?

$$y = 1 : p$$

# Loss Functions: Binary Classification

Given  $p, y$  : what's the **likelihood** of getting  $y$ ?

$$p^y(1 - p)^{1-y}$$

$$y = 0 : p^0(1 - p)^{1-0} = 1 - p$$

$$y = 1 : p^1(1 - p)^{1-1} = p$$

# Loss Functions: Binary Classification

All samples are independent

Total likelihood:  $\prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$

$p_i$  depends on the model:  $V, W, b, c$

Tune  $V, W, b, c \rightarrow$  Change  $p_i \rightarrow$  Max Likelihood

# Loss Functions: Binary Classification

$$\mathcal{L} = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

Maximizing  $\mathcal{L}$  equivalent to maximizing  $\log \mathcal{L}$   
( $\log$  is a monotonically increasing function)

$$\log \mathcal{L} = \sum_i^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

Nice: change things to a sum

# Loss Functions: Binary Classification

Usually have algorithms to minimize functions

Maximizing  $f$  equivalent to minimizing  $-f$

Minimize:

$$\mathcal{L}' = -\log \mathcal{L} = \boxed{-\sum_i^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i)}$$

# Loss Functions: Binary Classification

$$\mathcal{L} = -\sum_i^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

For one example, if  $y = 1$  :  $\mathcal{L} = -\log(p_i)$

$$y = 0 : \mathcal{L} = -\log(1 - p_i)$$

# Loss Functions: Binary Classification

$$y = 1 : \mathcal{L} = -\log(p_i)$$

$$p_i = 1(\text{correct}) \mathcal{L} = -\log(1) = 0$$

$$p_i = 0(\text{incorrect}) \mathcal{L} = -\log(0) \rightarrow \infty$$

# Loss Functions: Binary Classification

$$y = 0 : \mathcal{L} = -\log(1 - p_i)$$

$$p_i = 0(\text{correct}) \mathcal{L} = -\log(1 - 0) = 0$$

$$p_i = 1(\text{incorrect}) \mathcal{L} = -\log(1 - 1) \rightarrow \infty$$

# Loss Functions: Multi-class Classification

$$y \in 1, 2, \dots, n$$

# Loss Functions: Multi-class Classification

$$y \in 1, 2, \dots, n$$

$$\mathcal{L} = -\sum_{i=1}^{n_{examples}} y_{i1} \log(p_{i1}) + y_{i2} \log(p_{i2}) + \dots + y_{in} \log(p_{in})$$

# Loss Functions: Multi-class Classification

$$y \in 1, 2, \dots, n$$

$$\mathcal{L} = -\sum_{i=1}^{n_{examples}} y_{i1} \log(p_{i1}) + y_{i2} \log(p_{i2}) + \dots + y_{in} \log(p_{in})$$

$$\mathcal{L} = -\sum_{i=1}^{n_{examples}} \sum_{j=1}^n y_{ij} \log(p_{ij})$$

# Loss Functions: Multi-class Classification

$$y \in 1, 2, \dots, n$$

$$\mathcal{L} = -\sum_{i=1}^{n_{examples}} y_{i1} \log(p_{i1}) + y_{i2} \log(p_{i2}) + \dots + y_{in} \log(p_{in})$$

$$\mathcal{L} = -\sum_{i=1}^{n_{examples}} \sum_{j=1}^n y_{ij} \log(p_{ij})$$

Note:  $p_{i1} + \dots + p_{in} = 1, y_{i1} + \dots + y_{in} = 1, y_{ij} \in 0, 1$

# Loss Functions: Custom

You can design one based on your problem

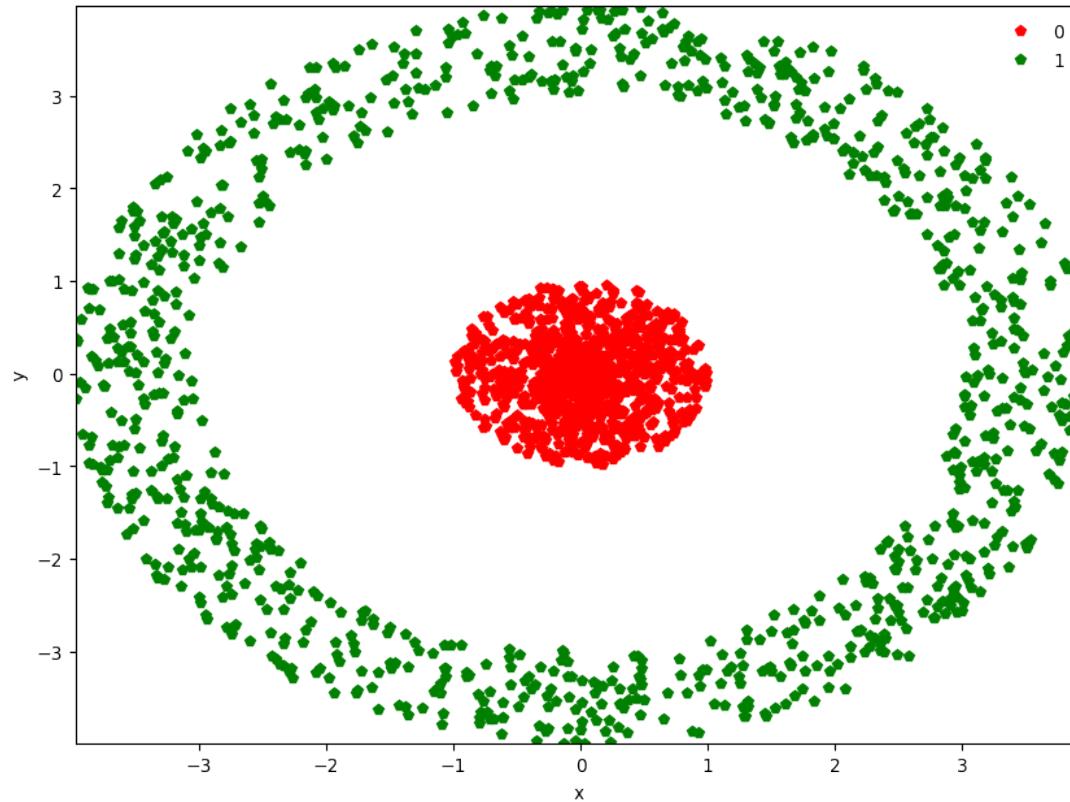
Well-designed loss greatly aids learning/training

Example: Kullback-Leibler divergence for comparing two distributions

Exercise: Go to kaggle and look at evaluation metrics for 10 competitions

Intuition: Why do  
neural networks  
work?

# Problem



Binary Classification: Given  $(x, y)$ , predict class 0 or 1

# Problem

Logistic Regression

$$\text{Compute: } \sigma(x, y) = \frac{1}{1 + e^{-(ax+by+c)}}$$

Probability of belonging to class 1

$a, b, c$  found by minimizing cross-entropy loss

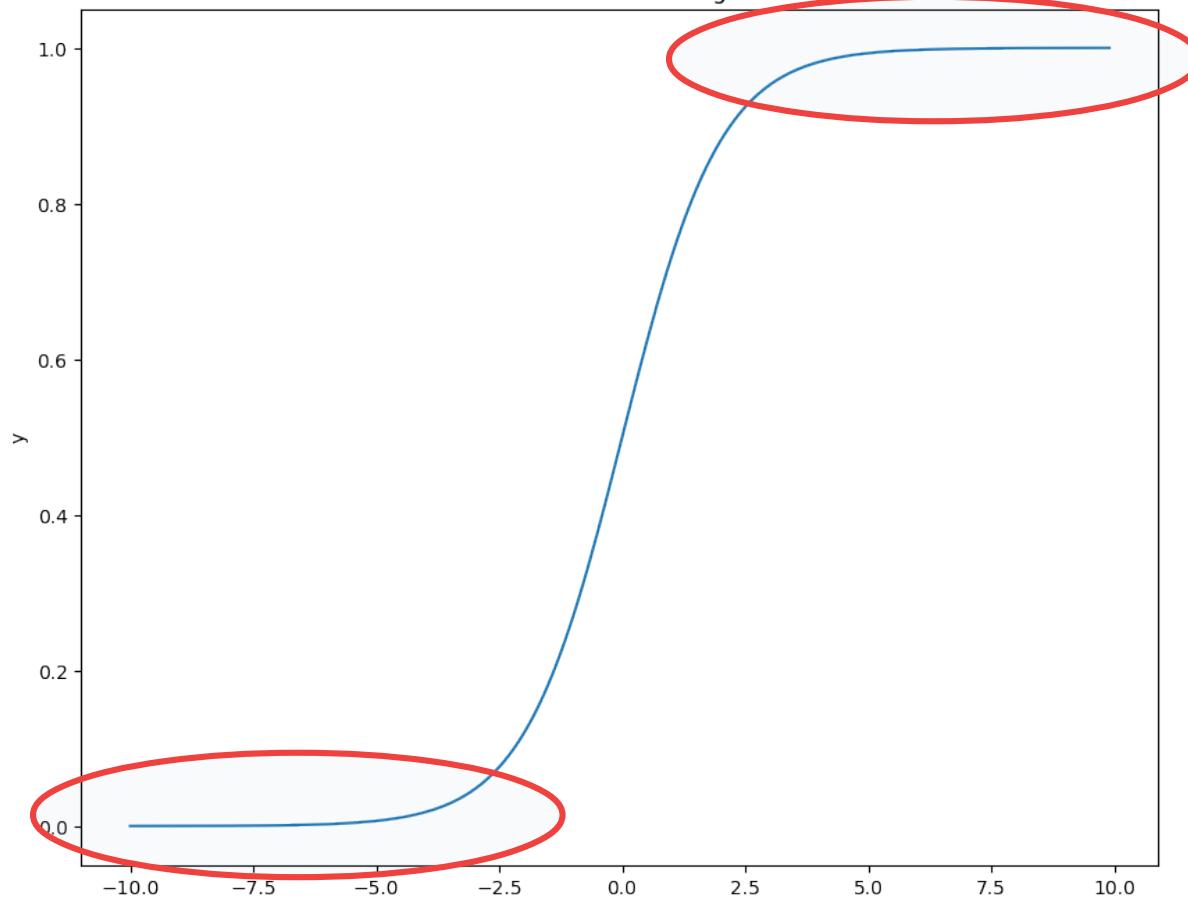
$$\mathcal{L} = -\sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i))$$

Label  $y_i = 0$  or  $y_i = 1$

Prediction:  $p_i = p_i(a, b, c) \in [0, 1]$

Positive  $x \rightarrow 1$  (on)

Activation Function: sigmoid



Negative  $x \rightarrow 0$  (off)

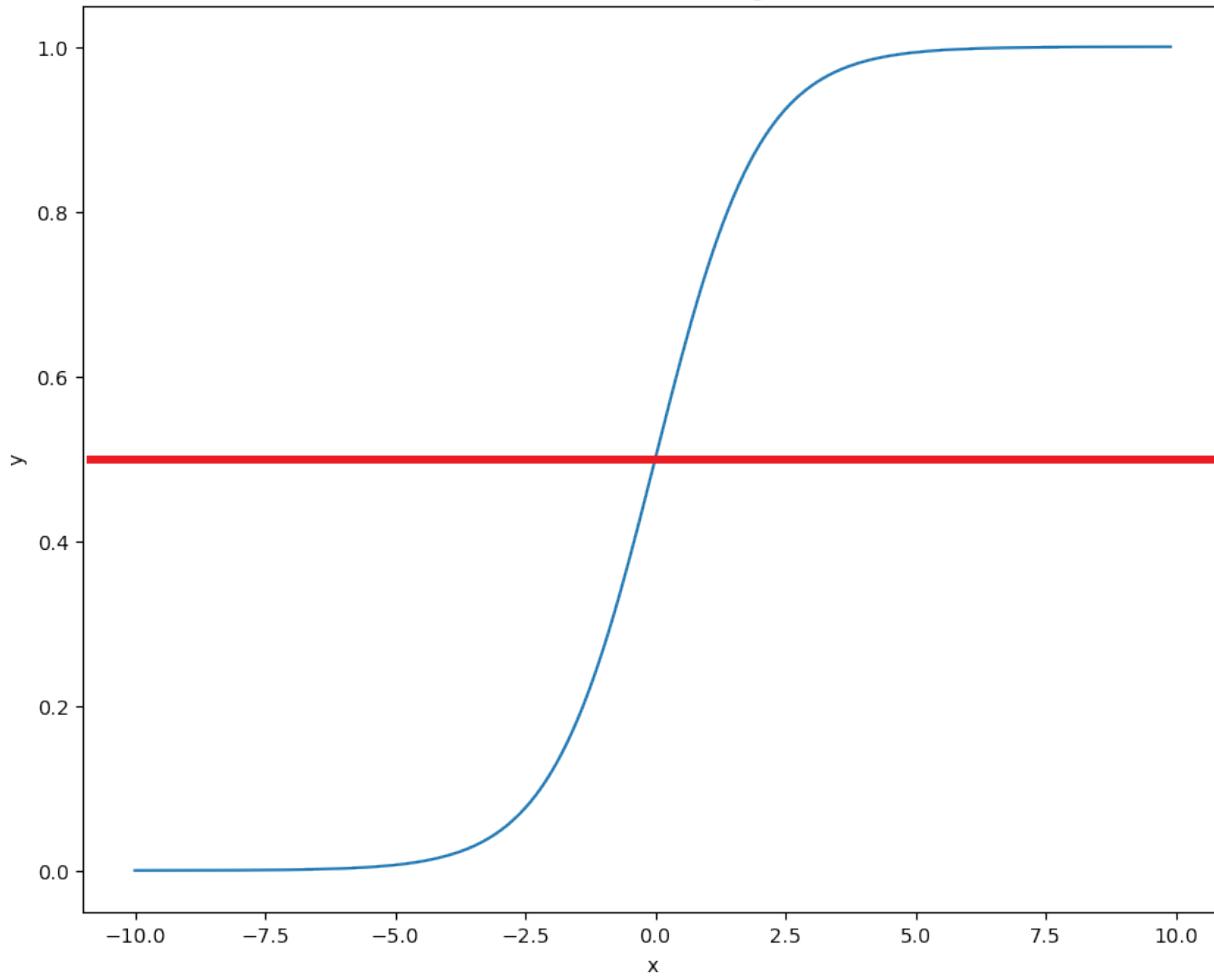
$$\sigma(ax + by + c) = \frac{1}{1 + e^{-(ax+by+c)}}$$

Class 1 (100%)

Probability

Class 0 (0%)

Activation Function: sigmoid



120

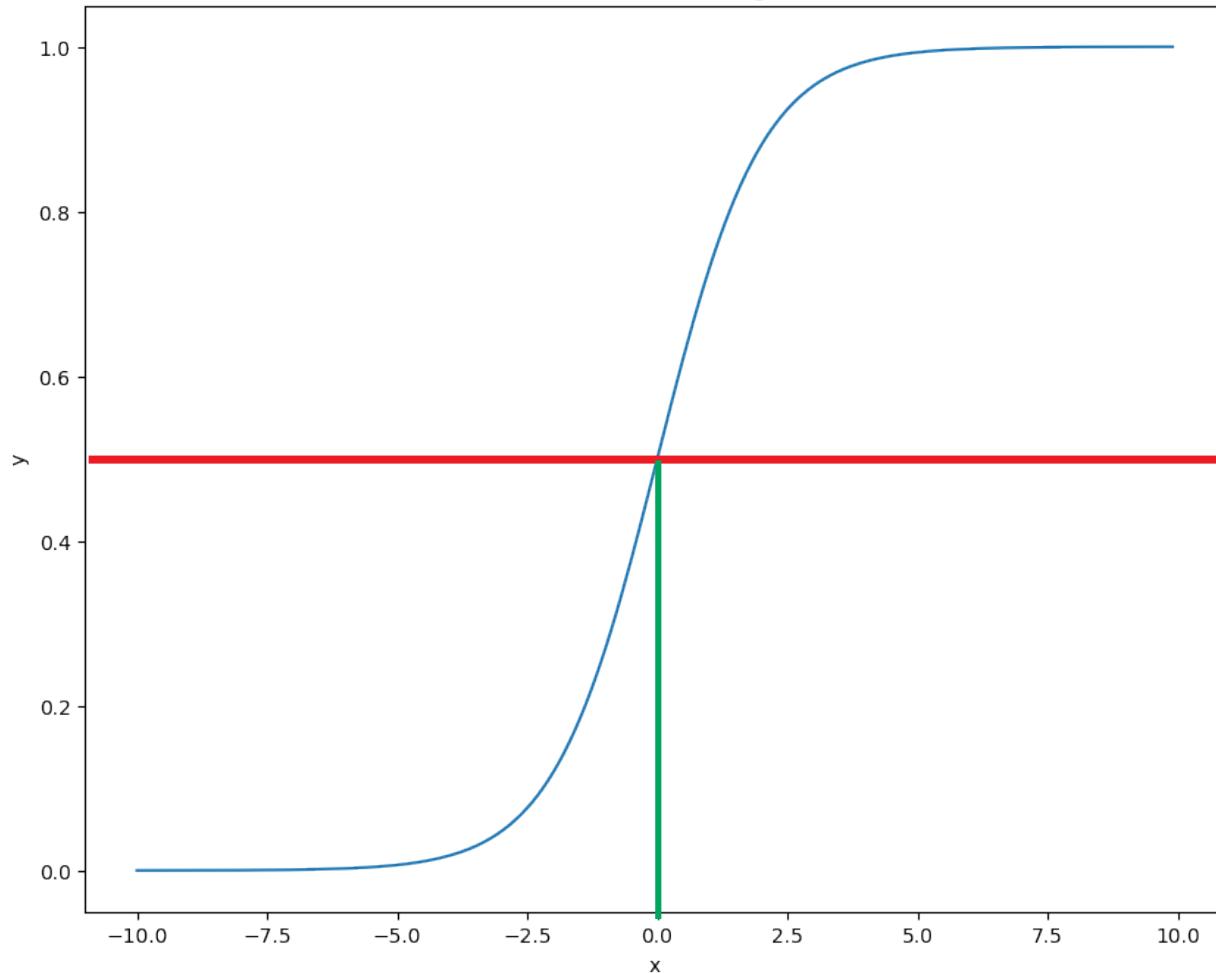
$y\text{-axis} = \sigma(x) = \text{probability of being in class 1}$

Class 1 (100%)

Probability

Class 0 (0%)

Activation Function: sigmoid



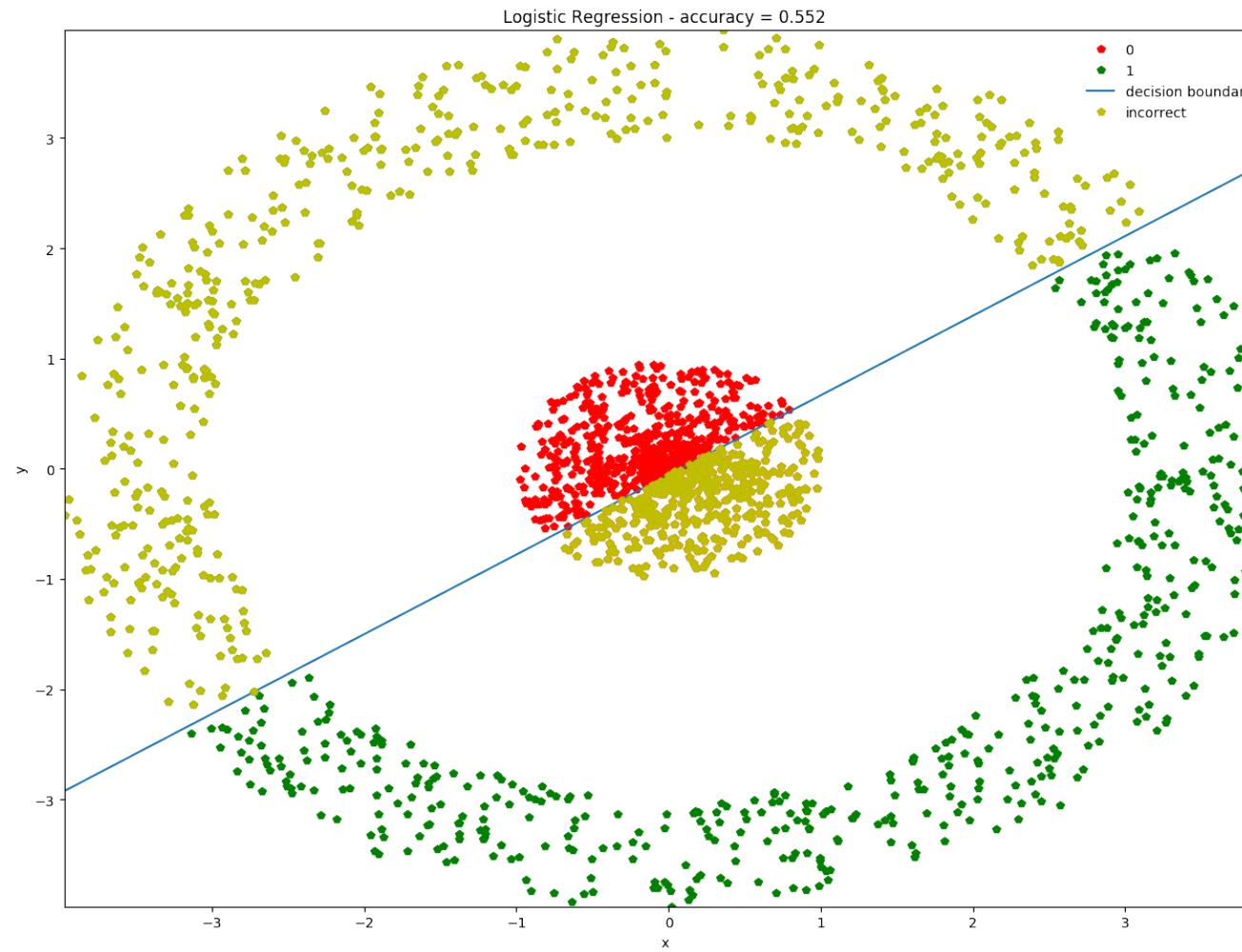
50% →  
Decision boundary

$y\text{-axis} = \sigma(x) = \text{probability of being in class 1}$

$$\sigma(x) = 0.5 \implies x(\text{input to sigmoid}) = 0$$

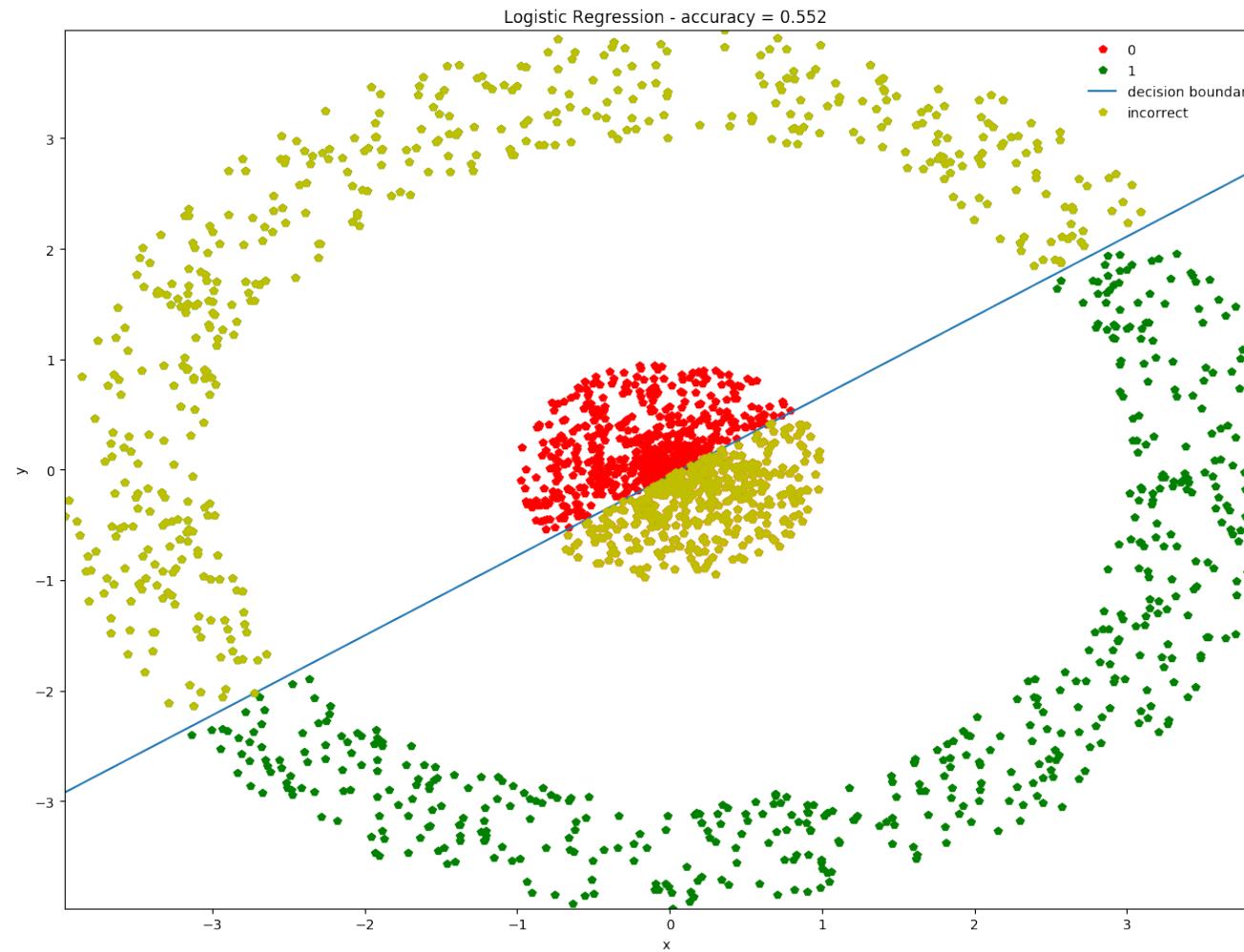
$$\sigma(ax + by + c) = 50\% \rightarrow \underbrace{ax + by + c = 0}_{\text{Straight Line}}$$

$$ax + by + c = 0 \implies y = -\frac{c}{b} - \frac{a}{b}x$$



No straight line can separate  
class 0 from class 1

Not **linearly separable**



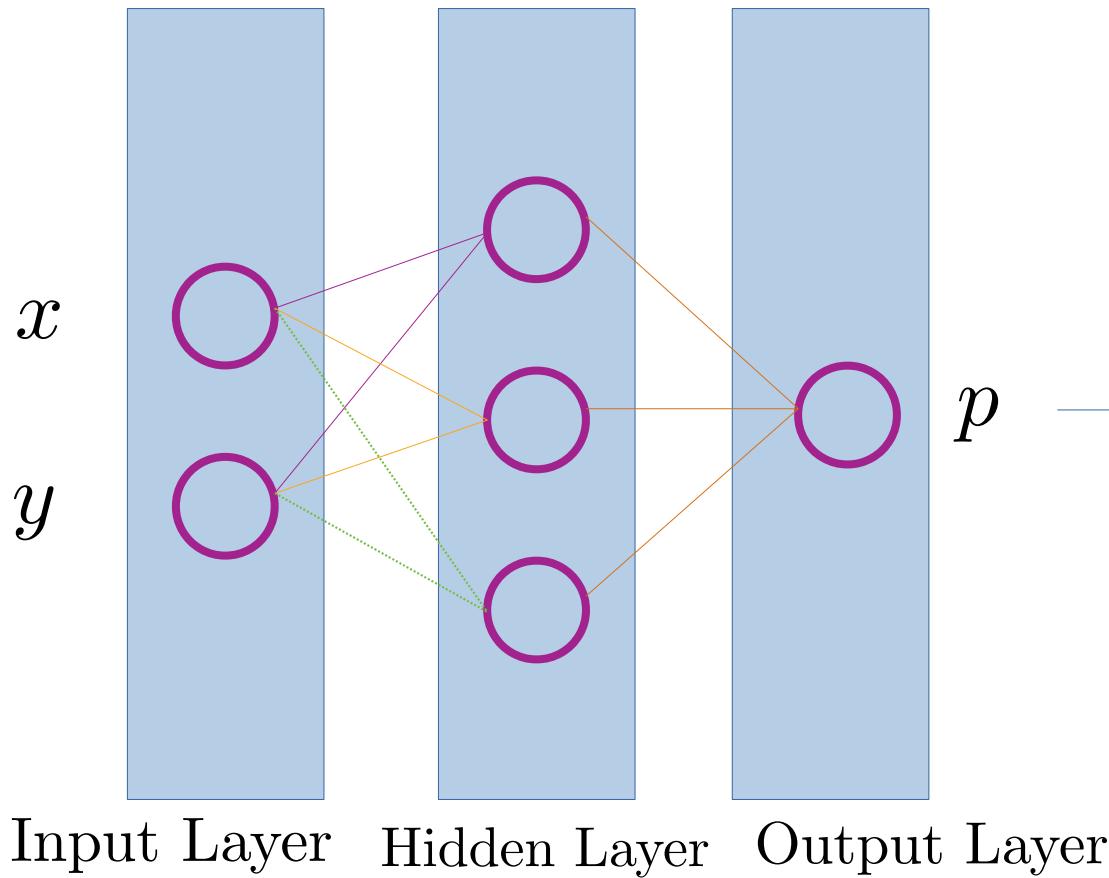
No straight line can separate  
class 0 from class 1

**One solution:**

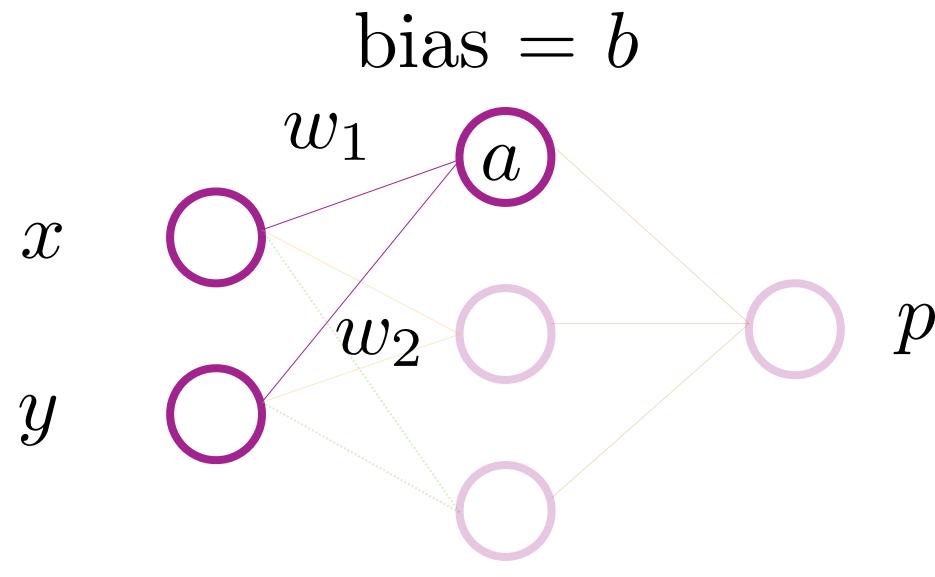
Convert to polar coordinates

$$r^2 = x^2 + y^2$$

**What about high  
dimensional data  
when we can't visualize?**



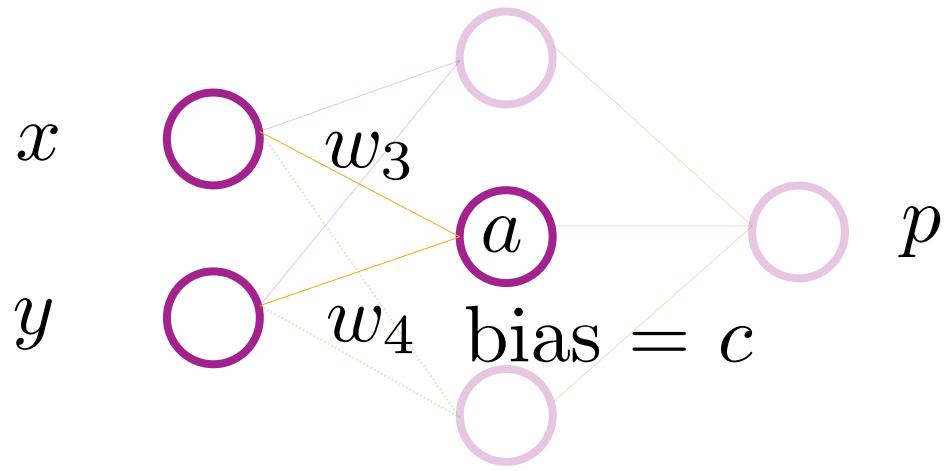
Minimize:  
Cross-entropy  
Activation = Sigmoid



$$a = \sigma(w_1x + w_2y + b)$$

Decision boundary:  
 $w_1x + w_2y + b = 0$

Input Layer   Hidden Layer   Output Layer

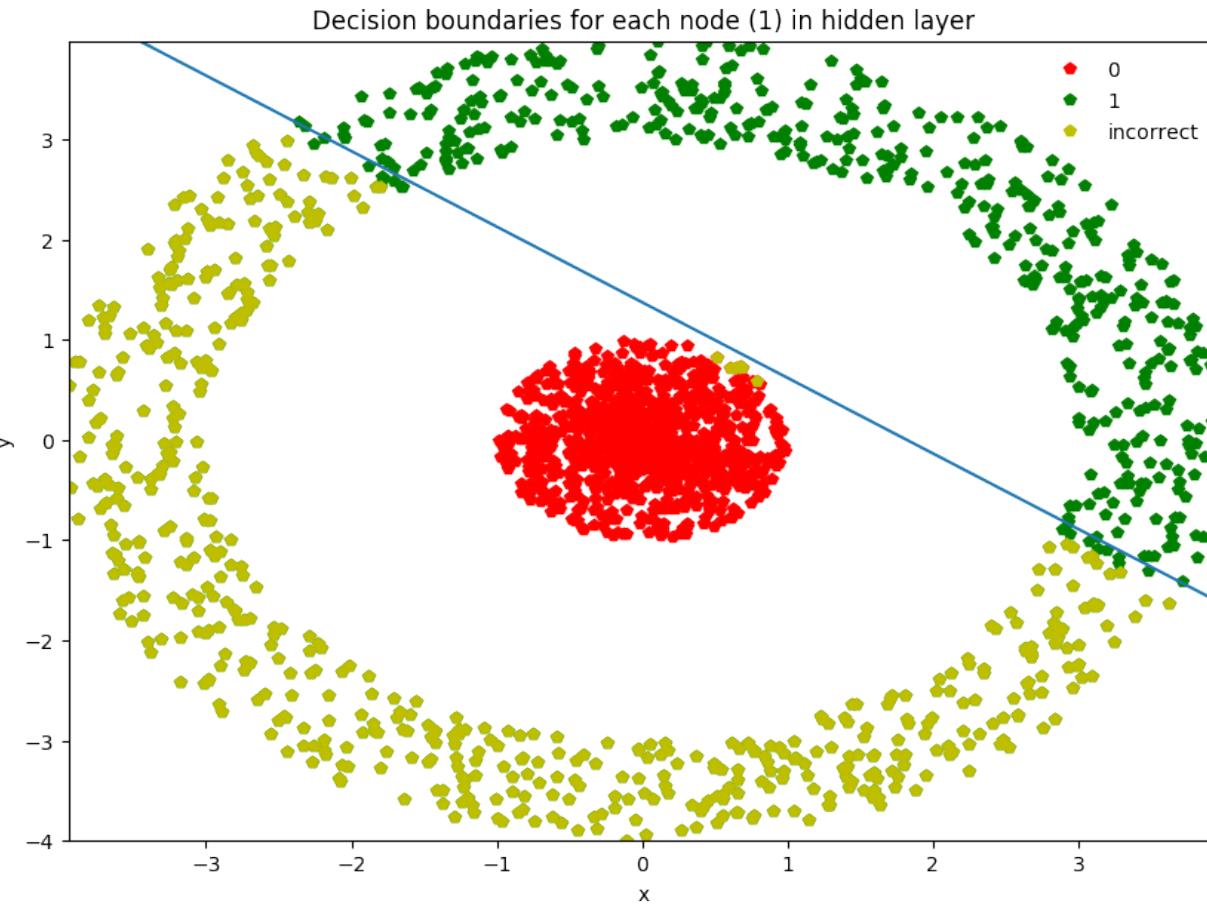


$$a = \sigma(w_3x + w_4y + c)$$

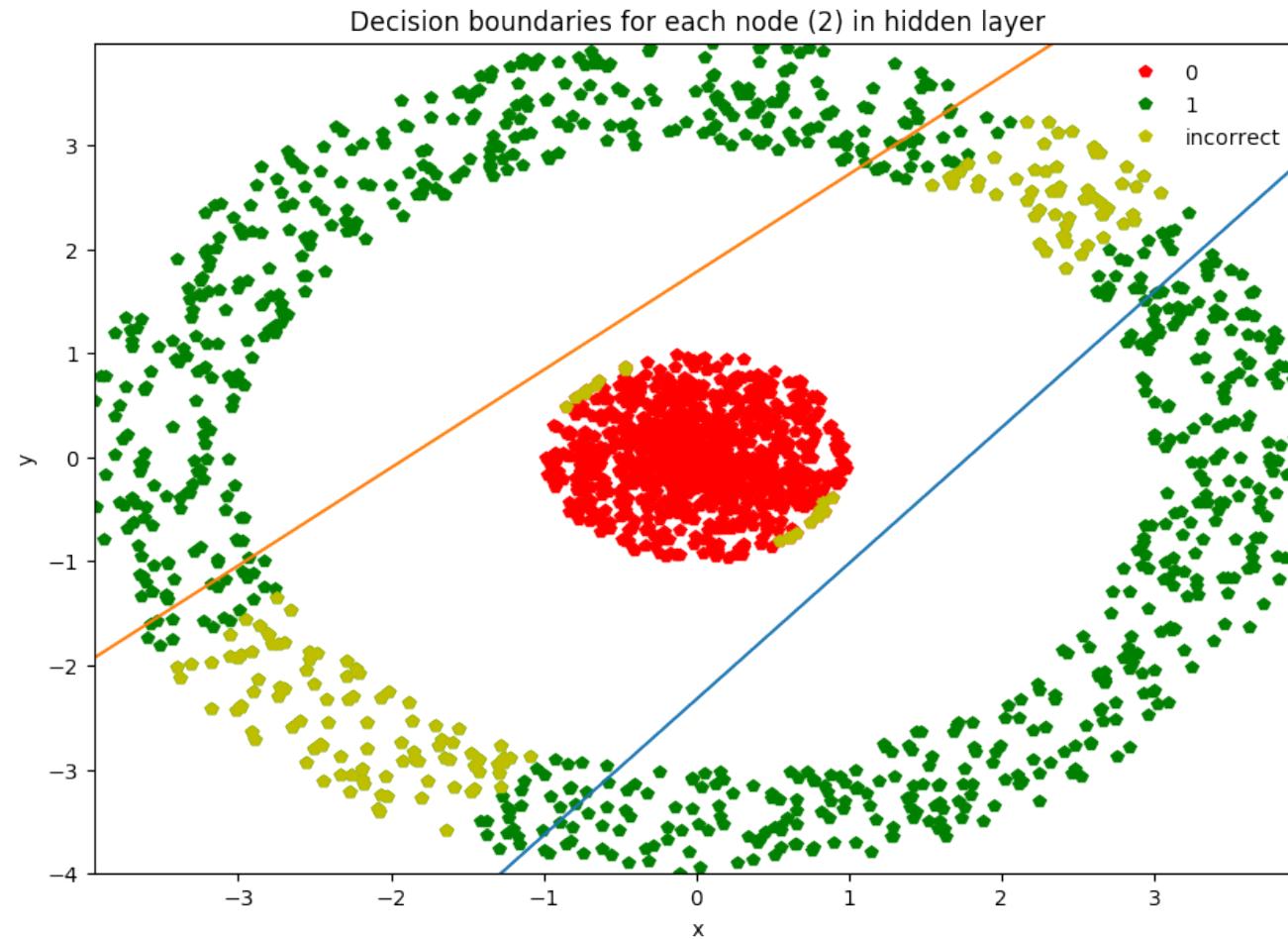
Decision boundary:  
 $w_3x + w_4y + c = 0$

Input Layer   Hidden Layer   Output Layer

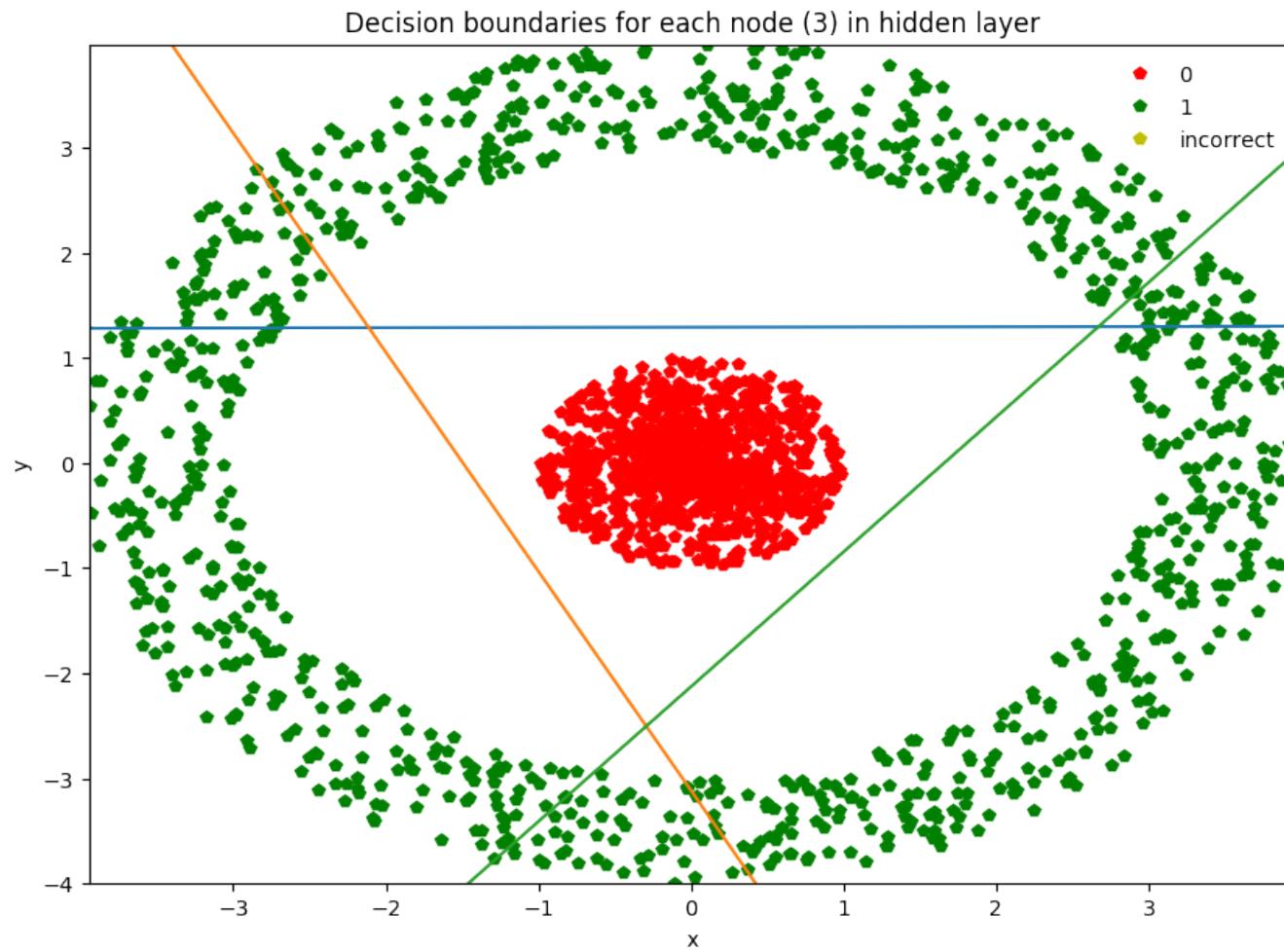
# Hidden Nodes = 1



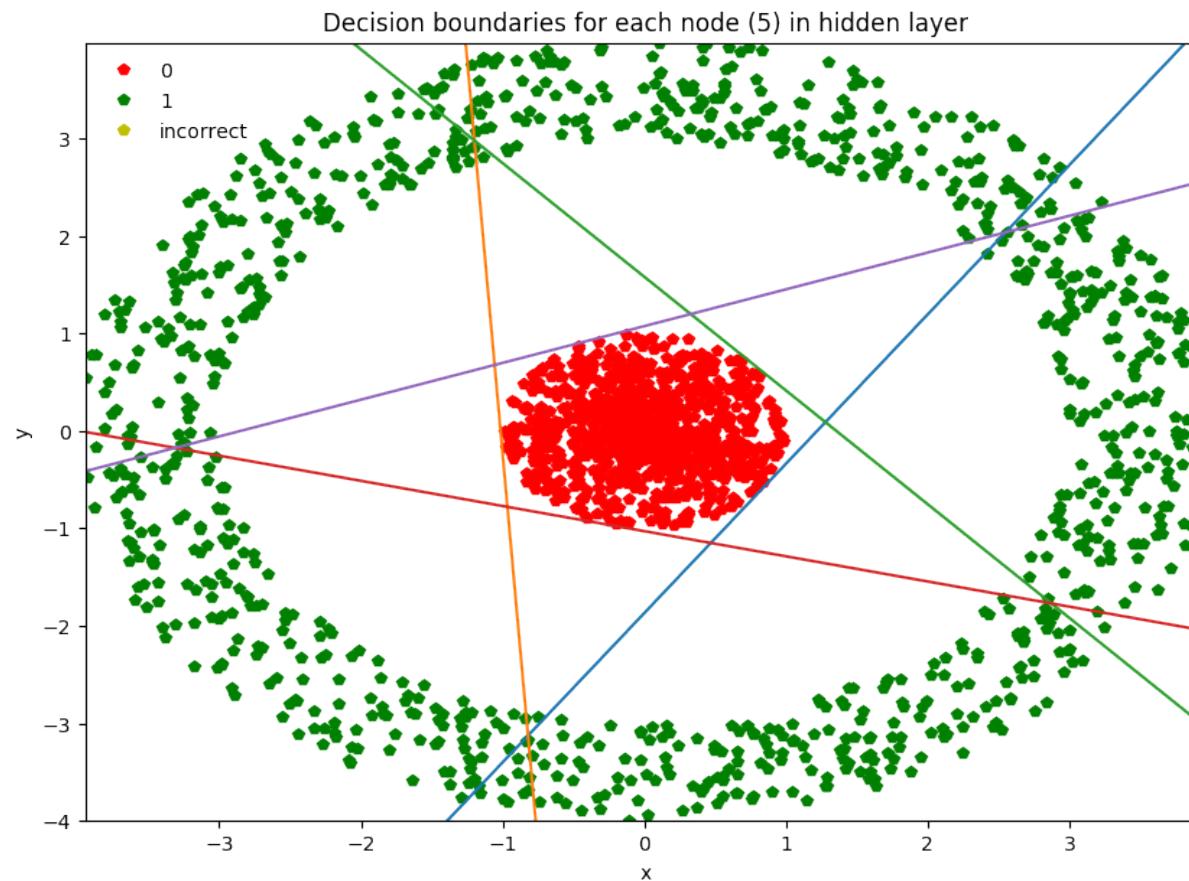
# Hidden Nodes = 2



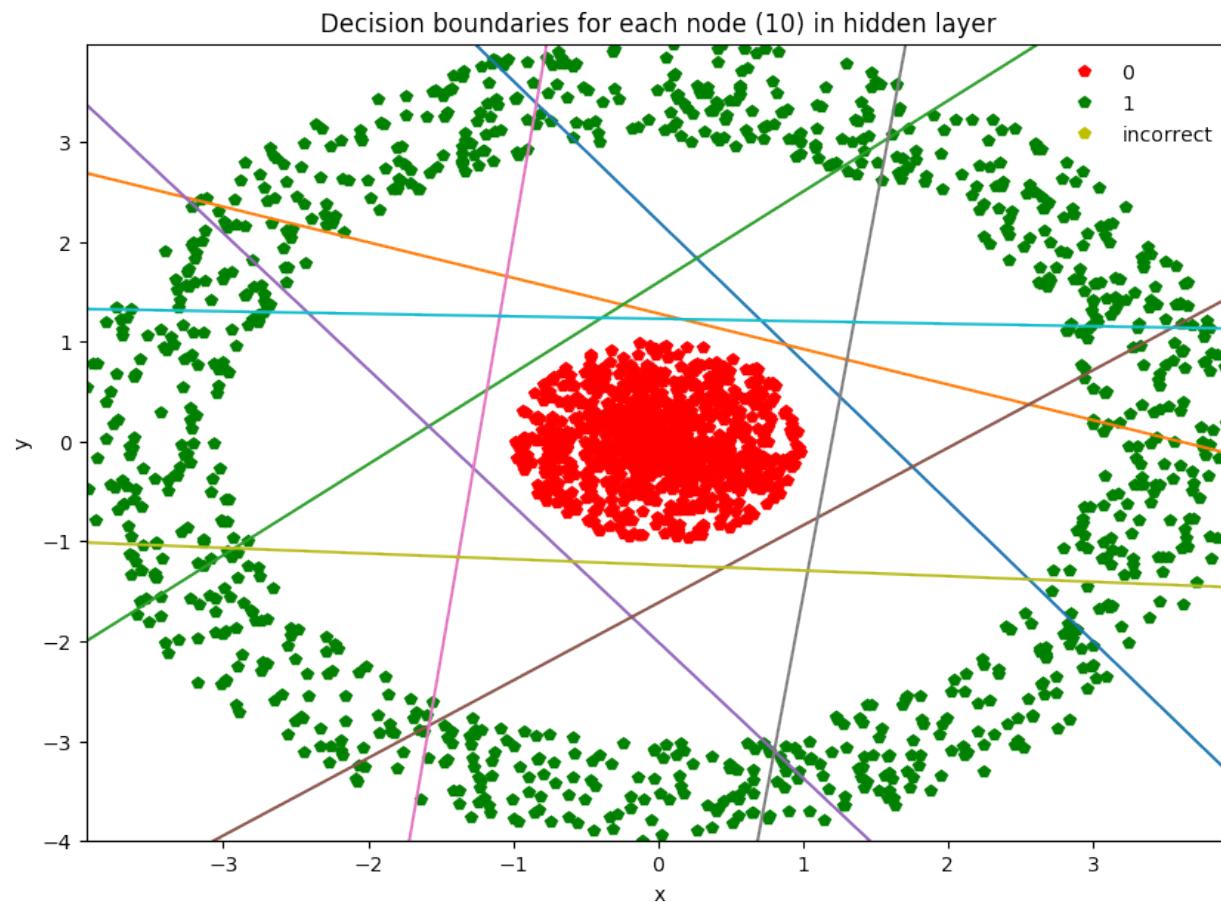
# Hidden Nodes = 3



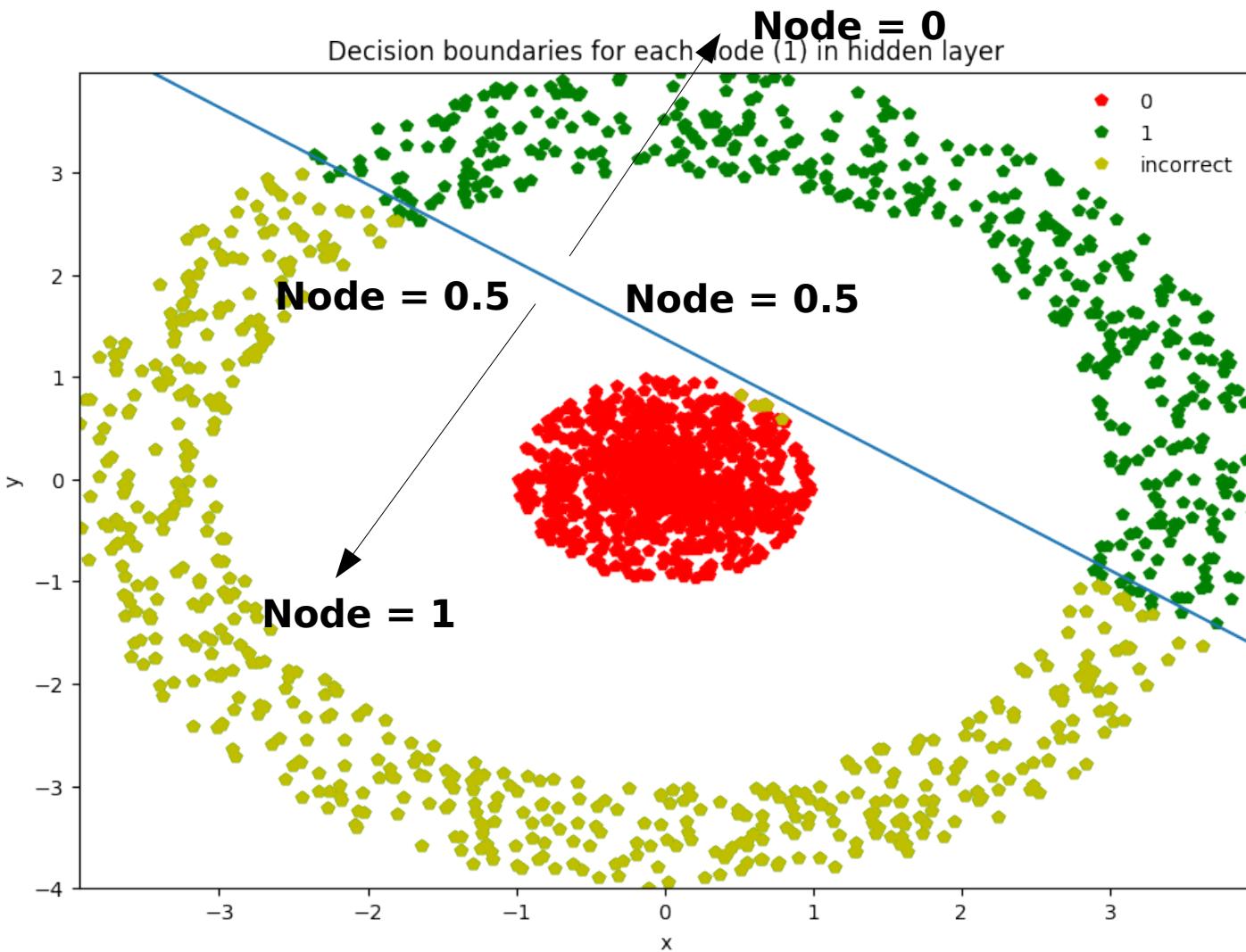
# Hidden Nodes =  
5



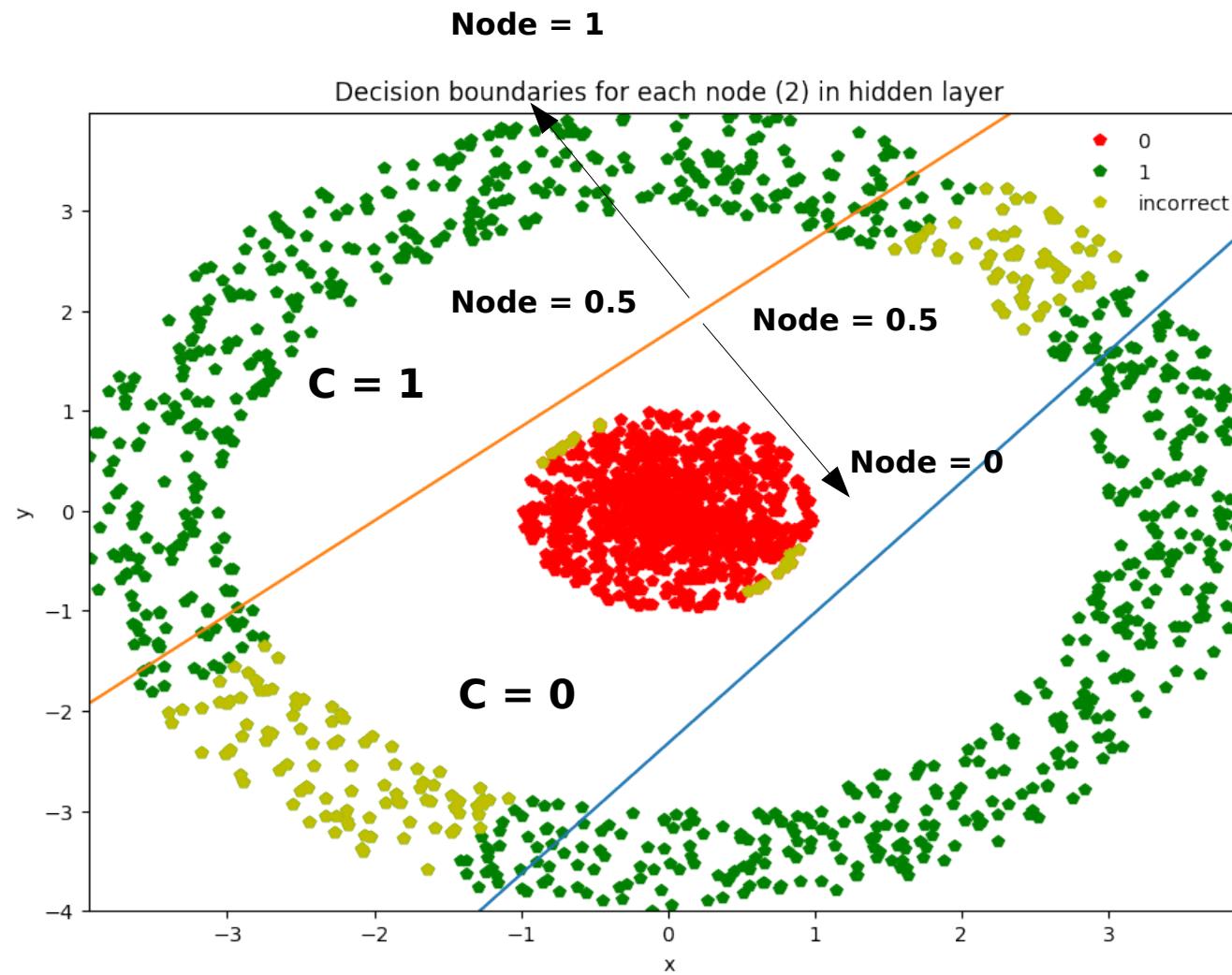
# Hidden Nodes = 10



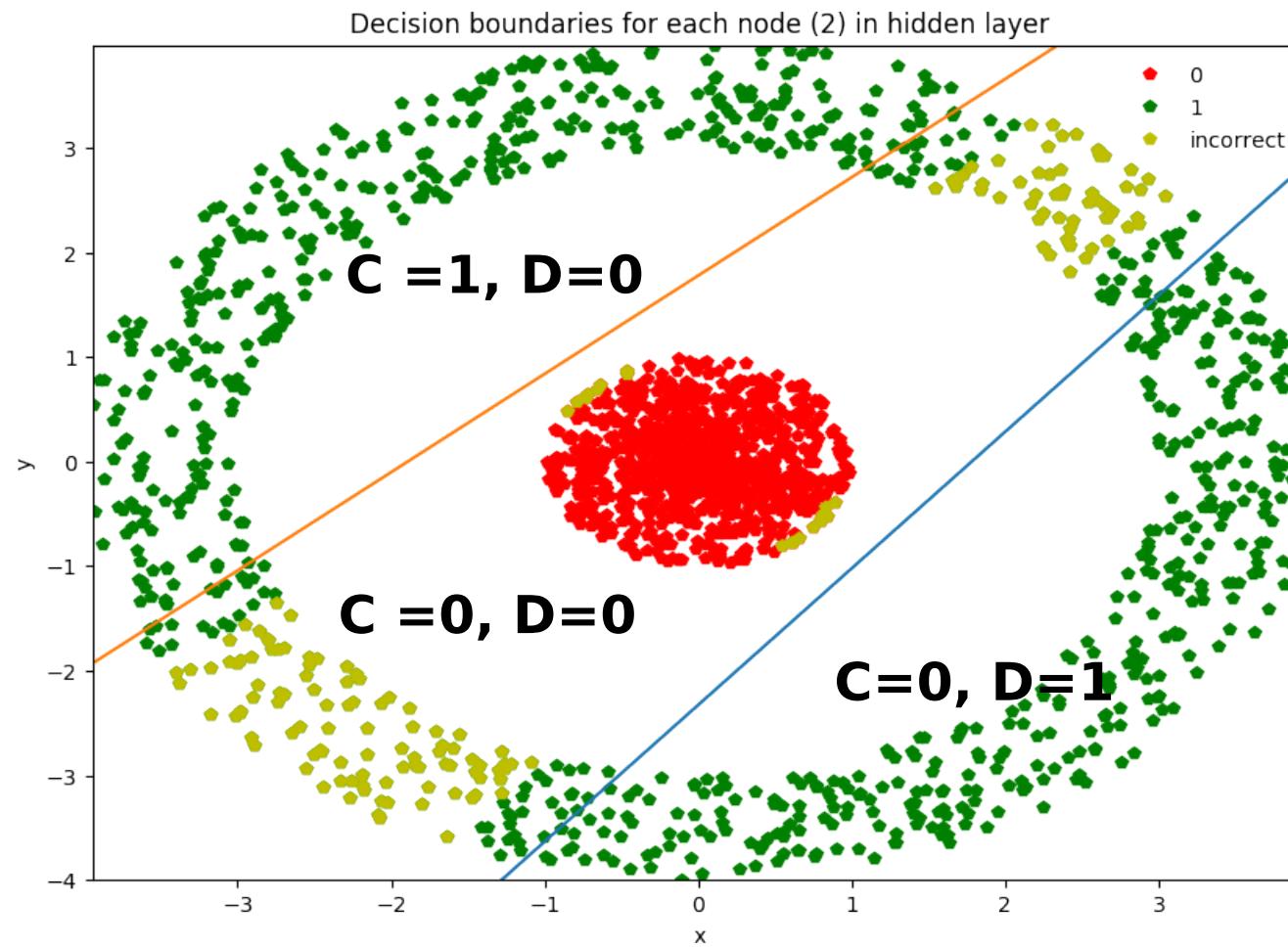
# Hidden Nodes = 1



# Hidden Nodes = 2



# Hidden Nodes = 2



Split input space into 3 regions

$C = 0, D = 0 \longrightarrow$  Mixed - both class 0 and class 1

$C = 1, D = 0 \longrightarrow$  Pure class 1

$C = 0, D = 1 \longrightarrow$  Pure class 1

$N$  hidden binary nodes can split input into  $2^N$  regions

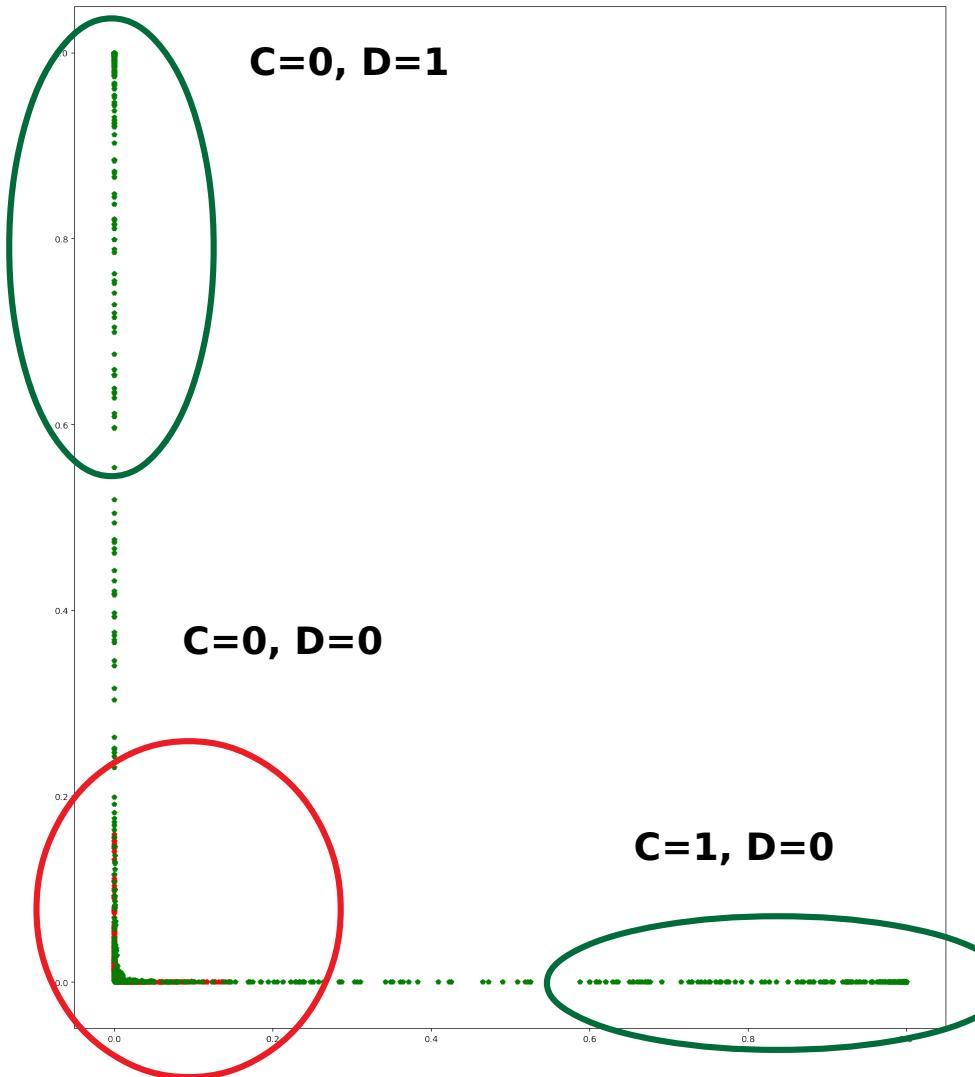
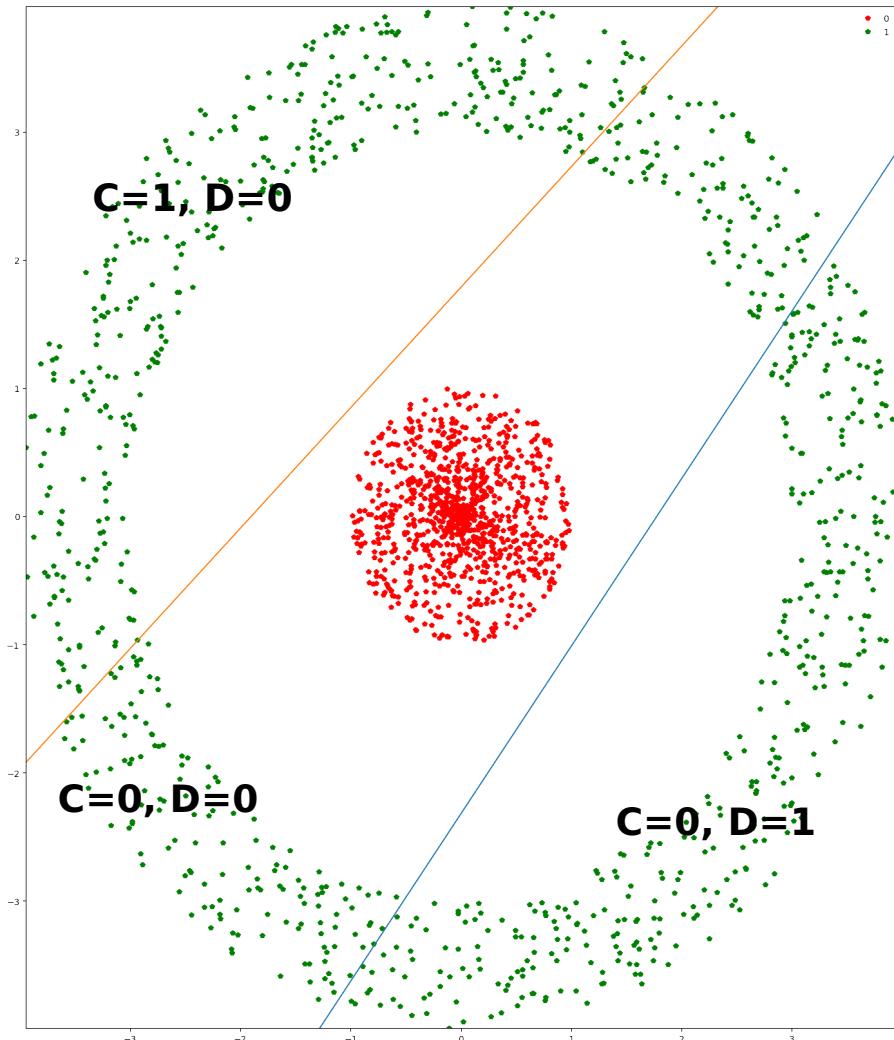
Keep increasing  $N$  (without overfitting) till regions "pure"

Map 2-dim input to N-dim space

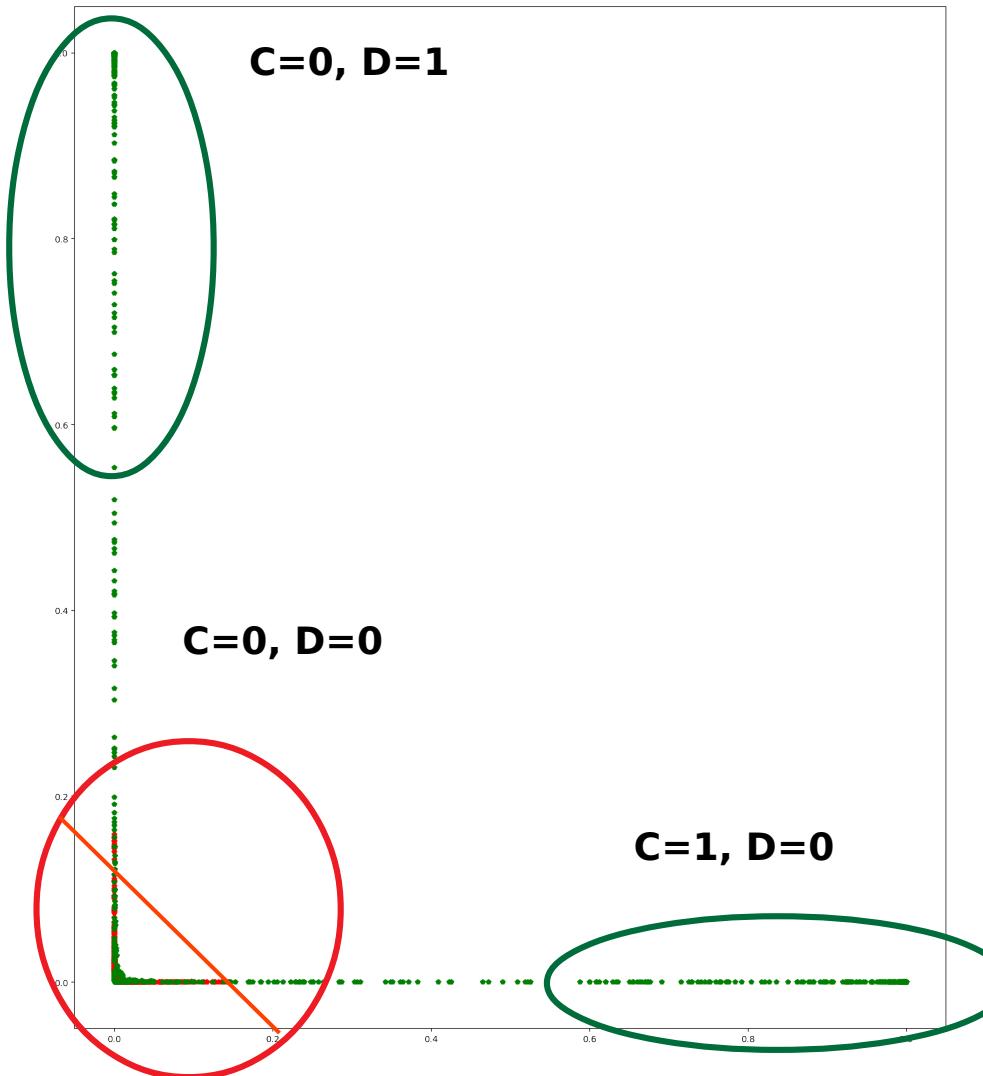
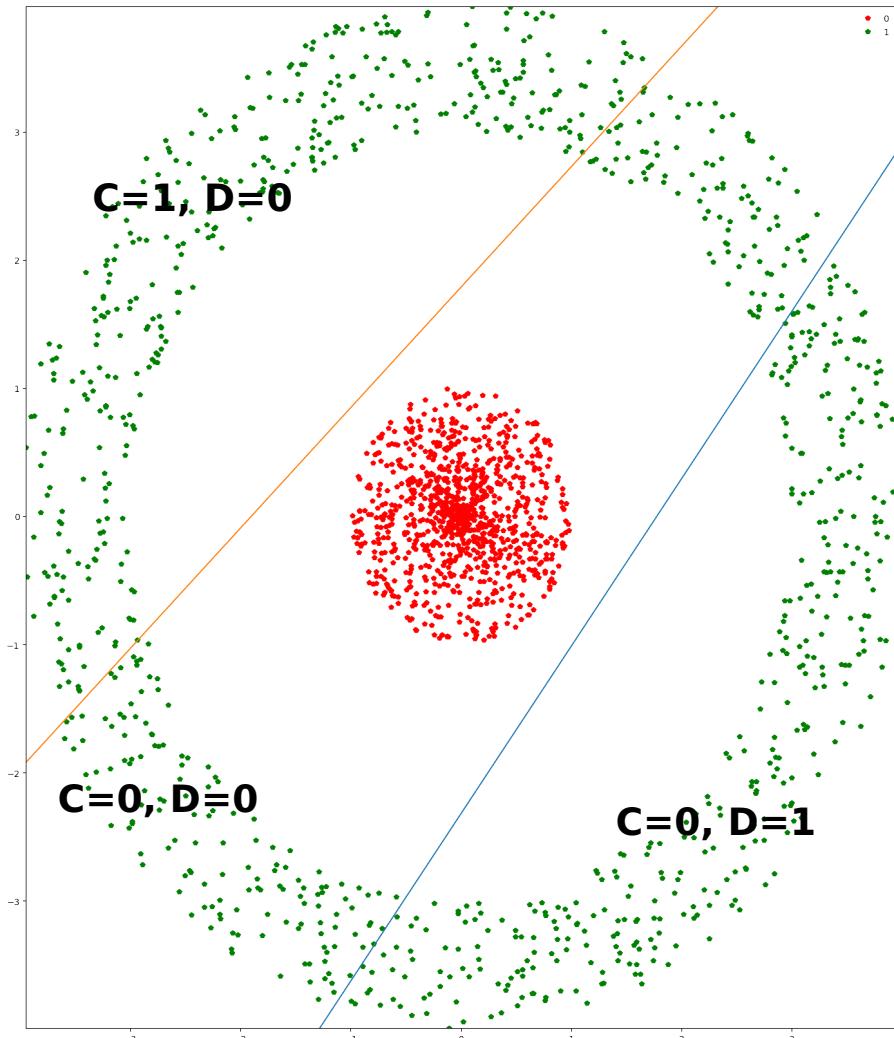
$$(x, y) \rightarrow (a_1, a_2, \dots, a_N) = \underbrace{(1, 0, \dots, 1)}_{\text{corners of N-dim hypercube}}$$

in the hope of linearly separating classes

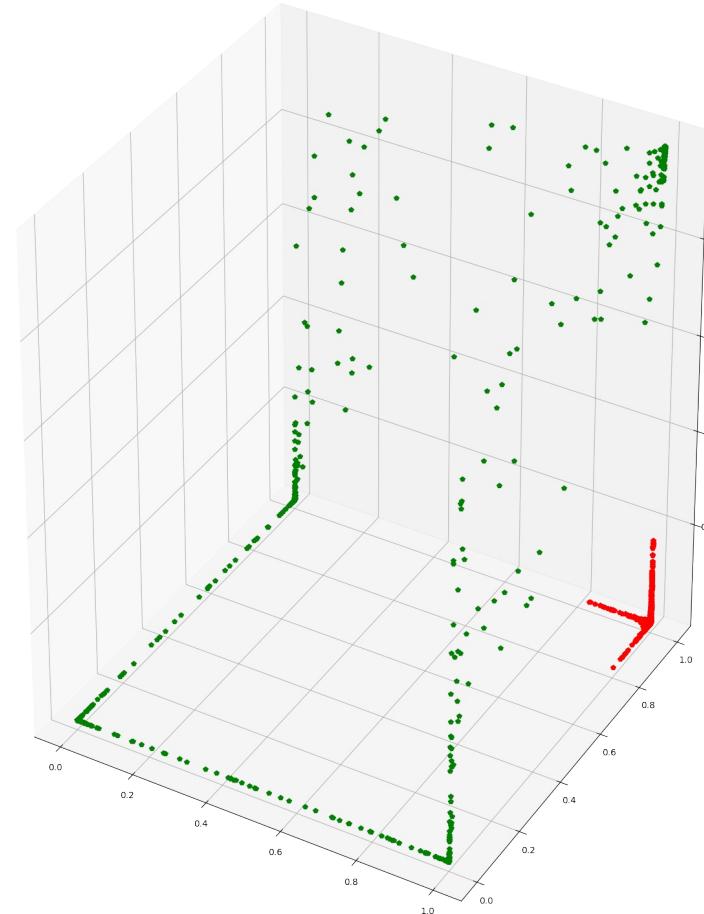
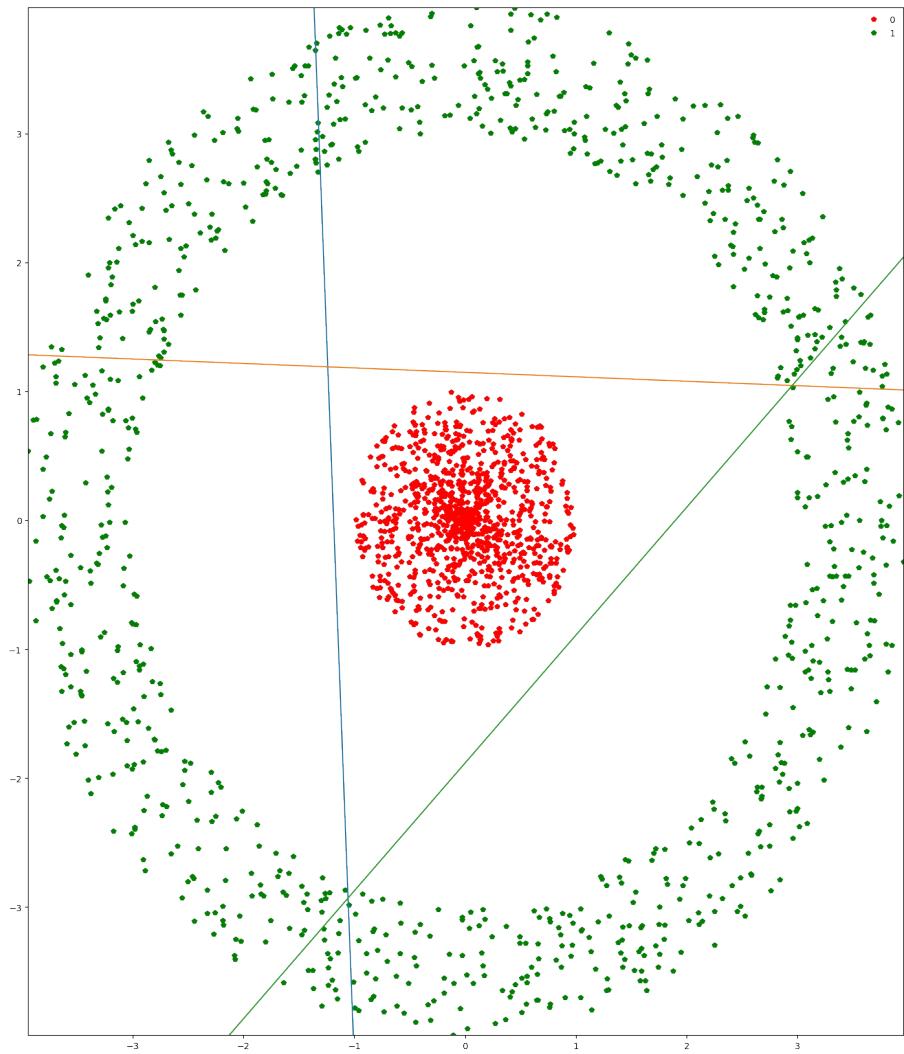
# Map 2-dim input to 2-dim space



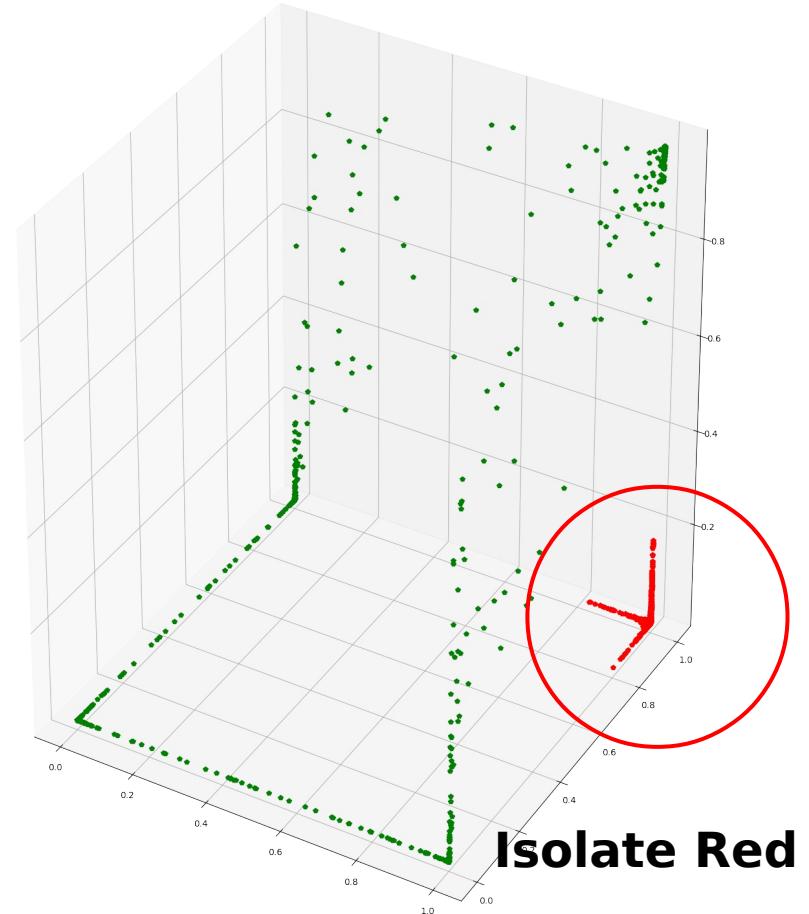
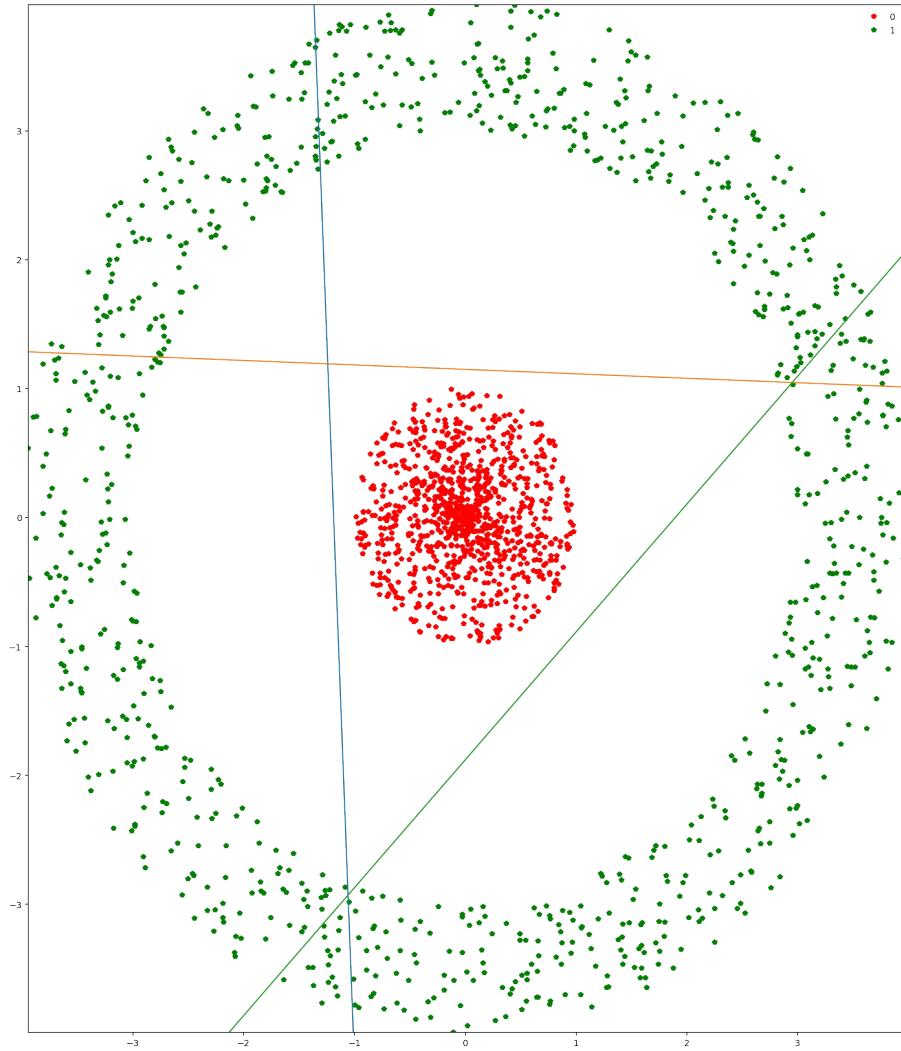
# Map 2-dim input to 2-dim space



# Map 2-dim input to 3-dim space

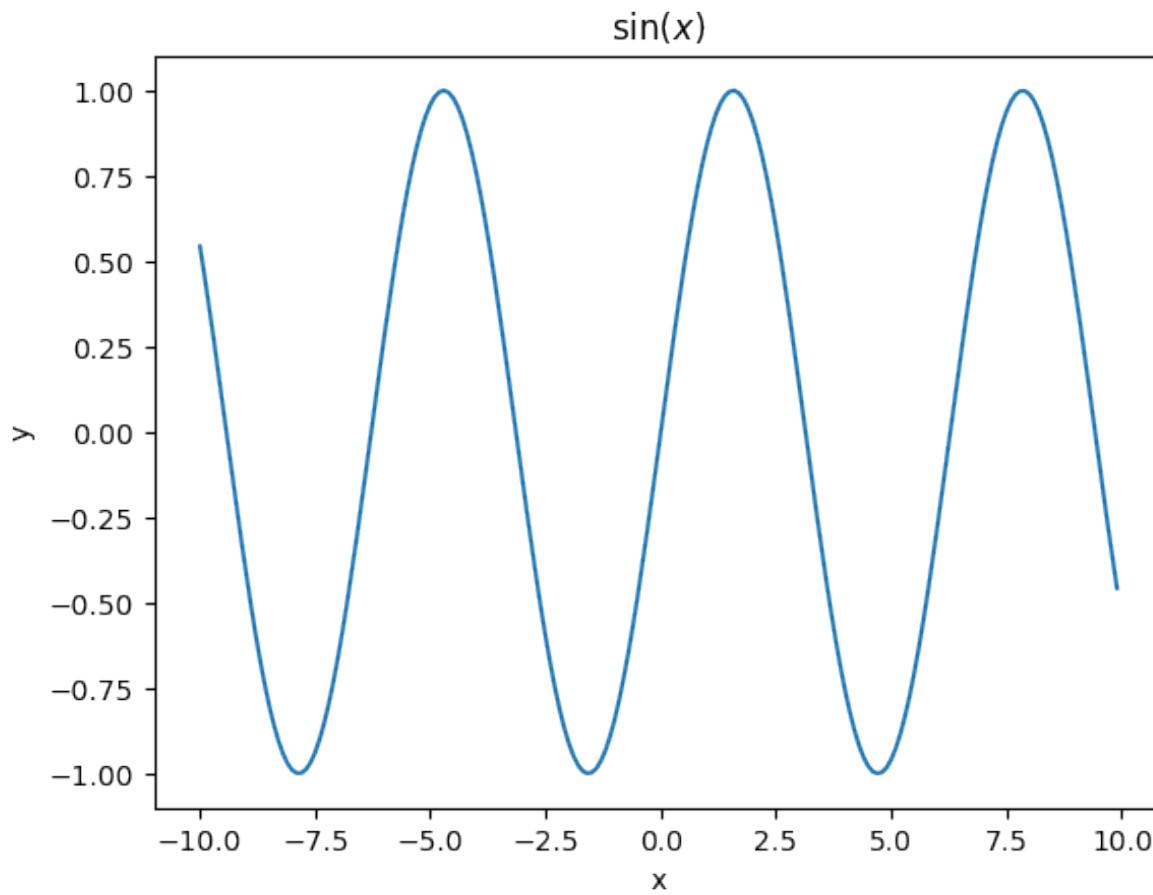


# Map 2-dim input to 3-dim space

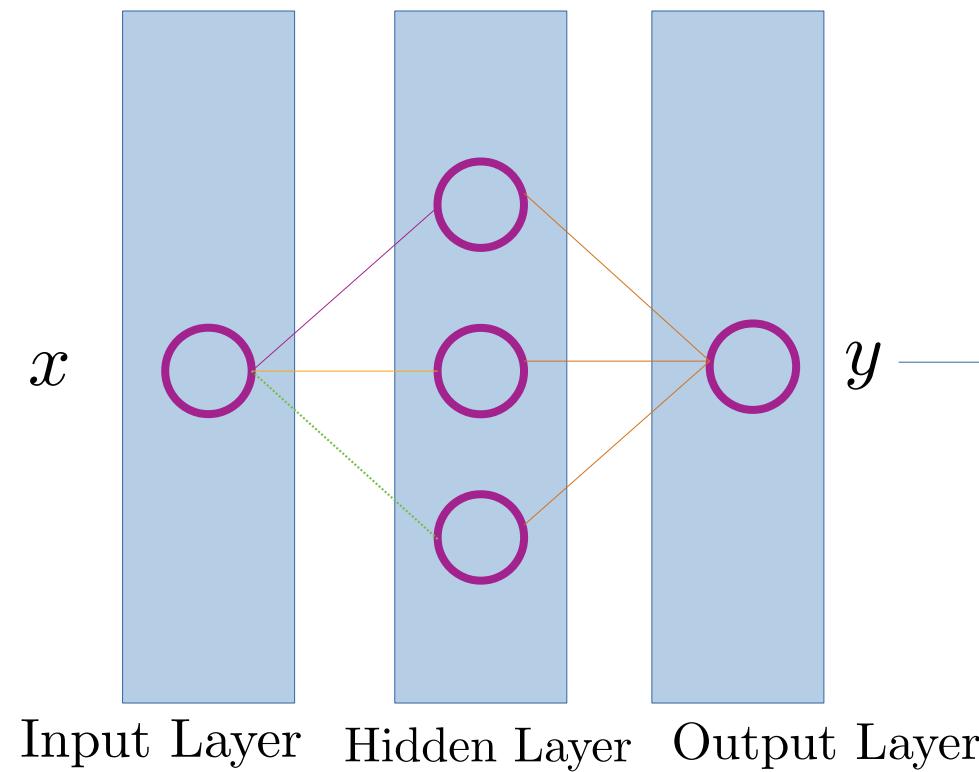


**Isolate Red in corner**

# A Second Example: Regression

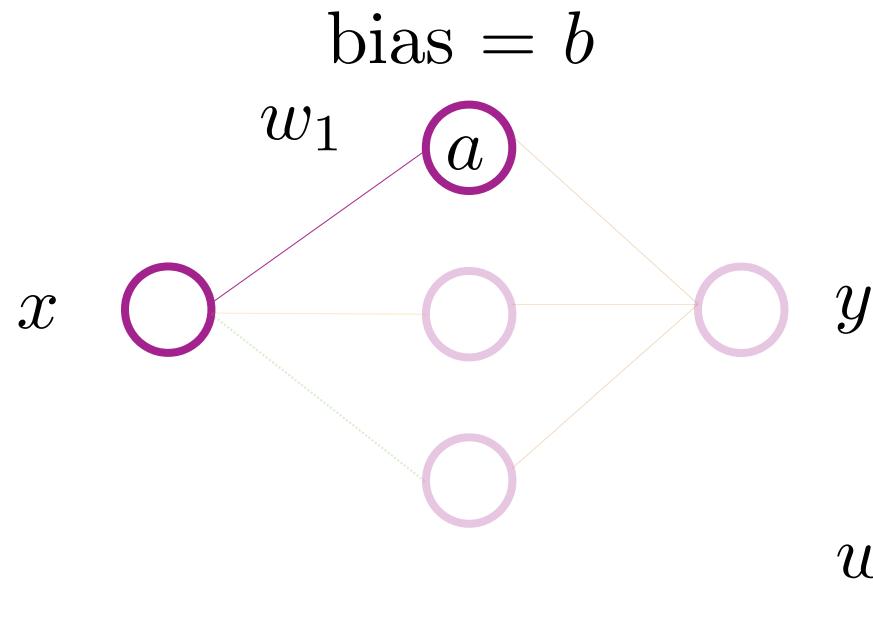


Given  $x$ , predict  $y = \sin(x)$



Minimize:  
Mean Squared Error

Activation = Sigmoid  
Output = Linear

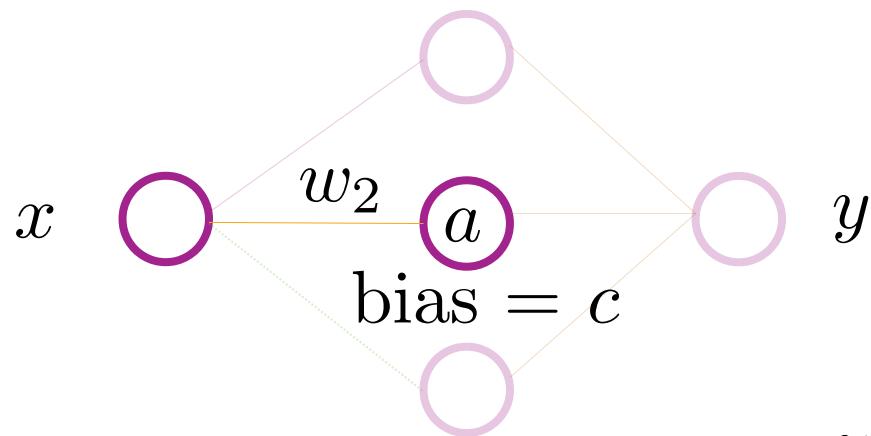


$$a = \sigma(w_1x + b)$$



Decision boundary:

$$w_1x + b = 0 \implies x = -\frac{b}{w_1}$$



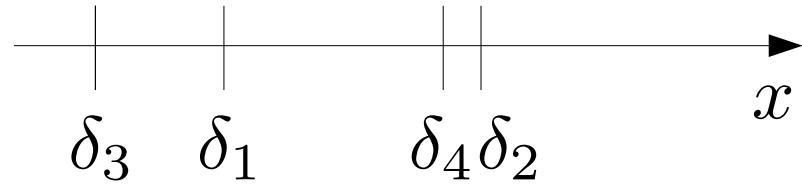
$$a = \sigma(w_2x + c)$$



Decision boundary:

$$w_2x + c = 0 \implies x = -\frac{c}{w_2}$$

## 1 decision boundary per hidden node



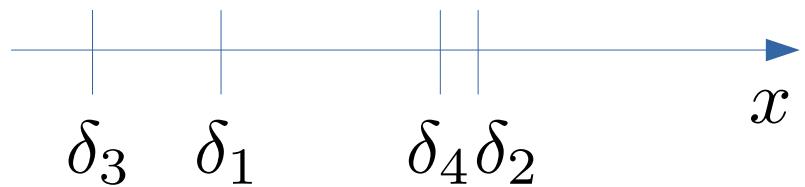
$$\delta_1 = -\frac{b_1}{w_1}$$

$$\delta_2 = -\frac{b_2}{w_2}$$

⋮

$$\delta_n = -\frac{b_n}{w_n}$$

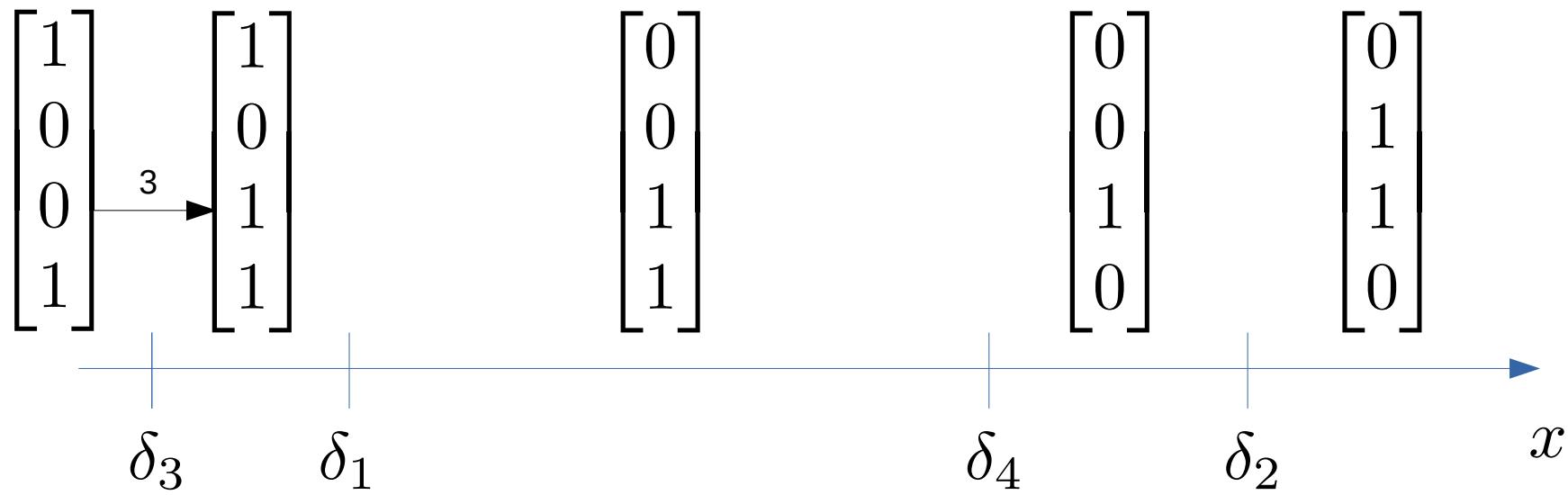
**Walk from left to right**

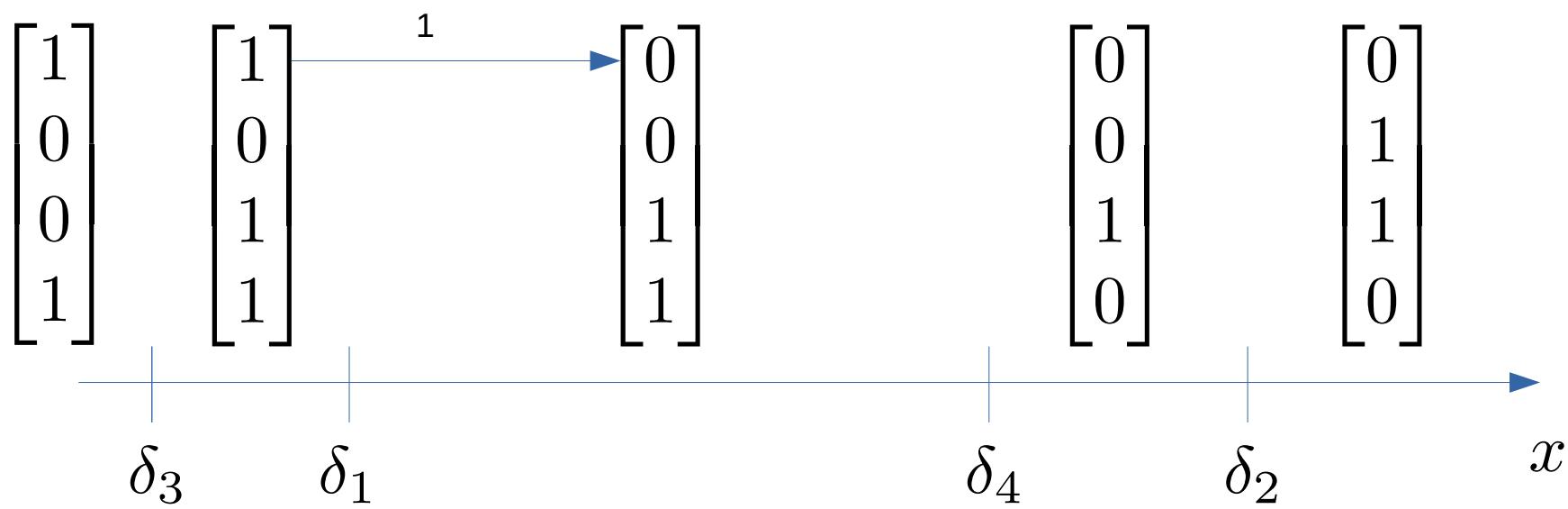


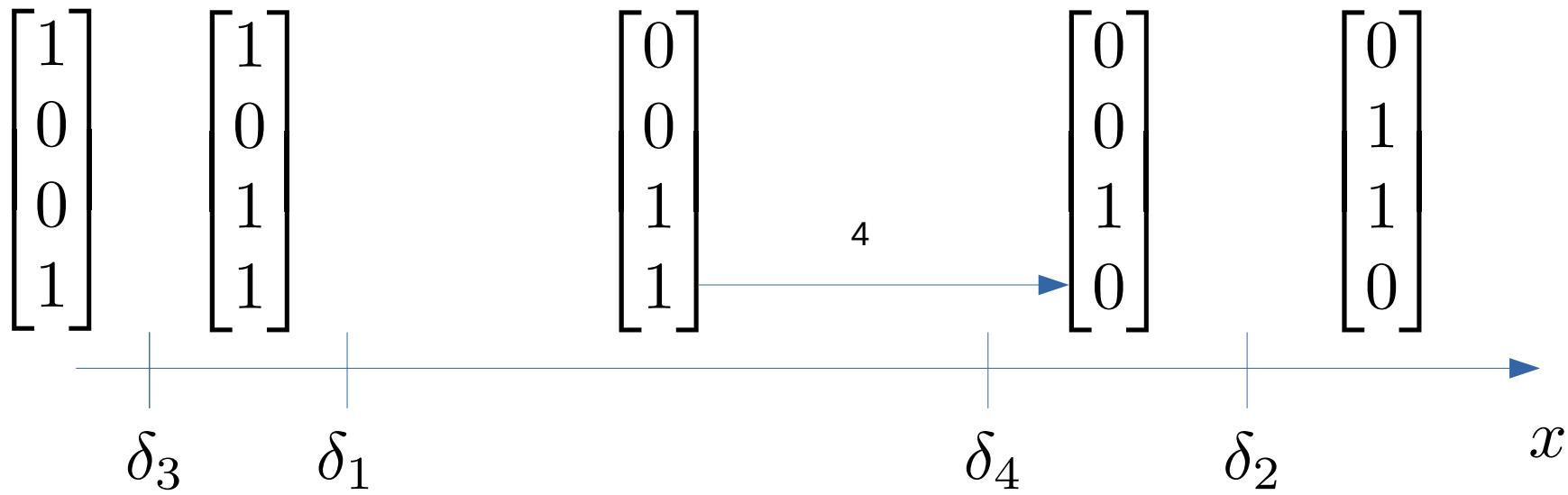
Each time you cross a boundary,  
a hidden node either:

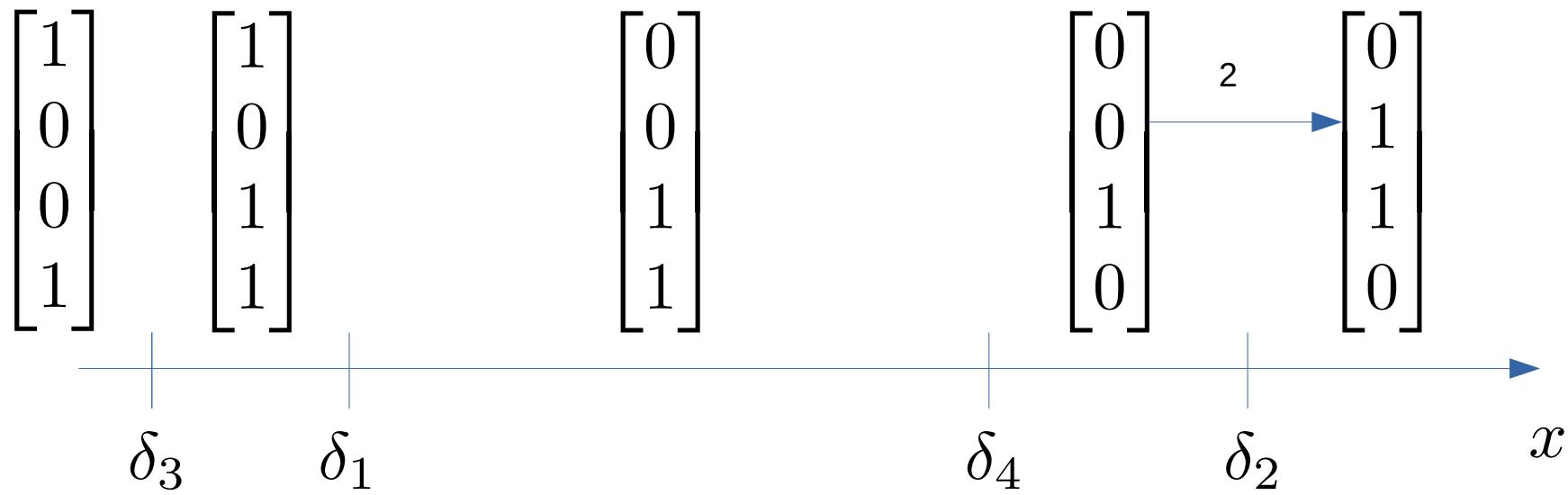
Turns on i.e. goes from  $1 \rightarrow 0$

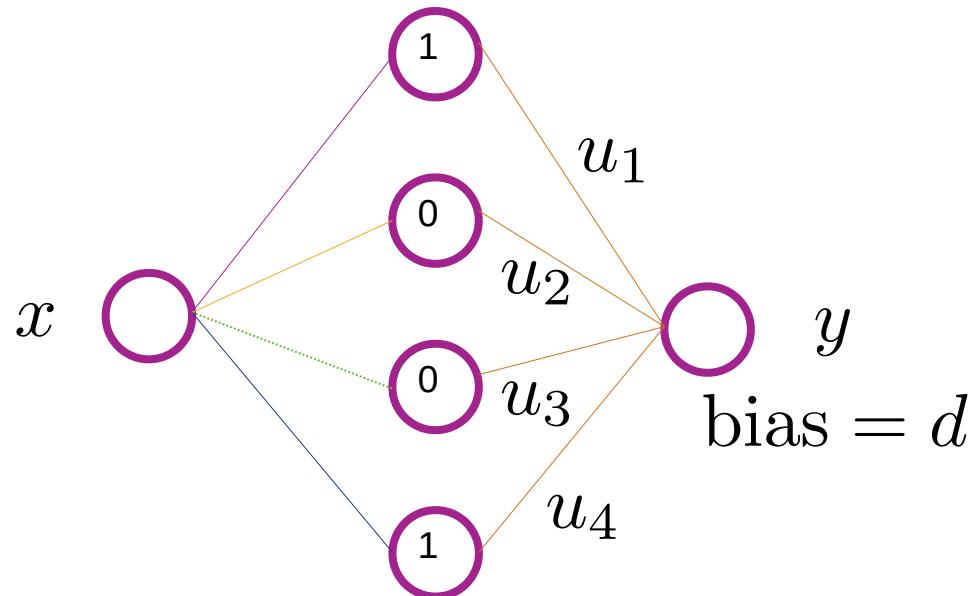
Turns off i.e. goes from  $0 \rightarrow 1$



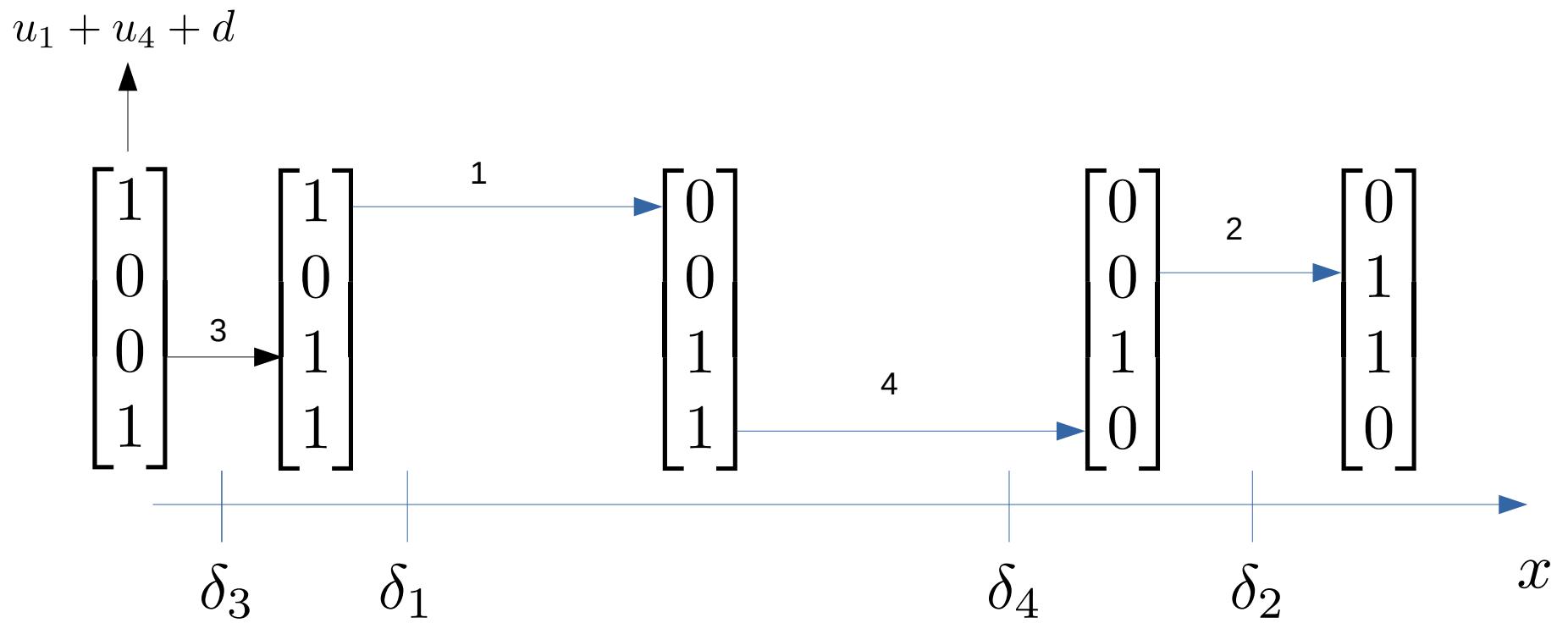


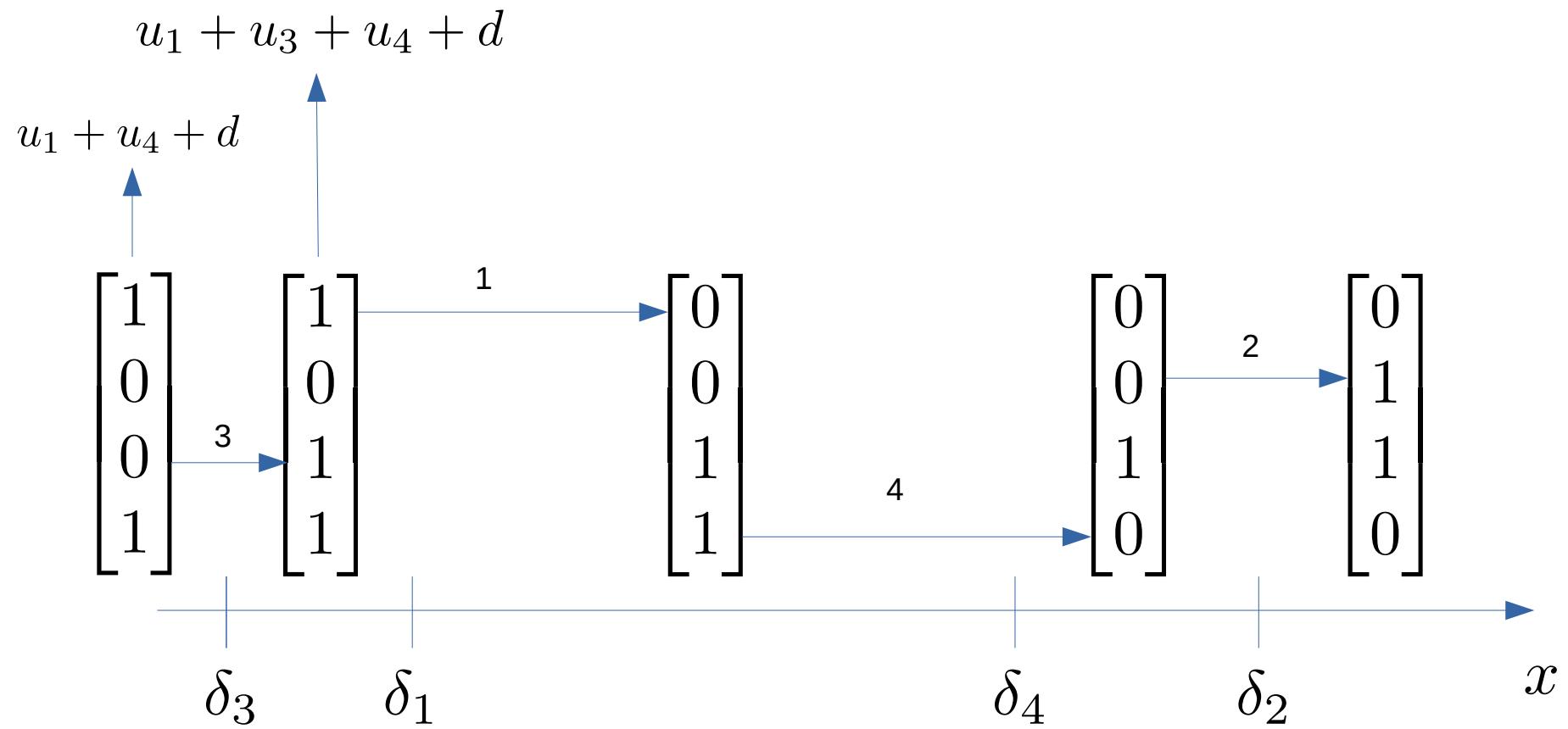


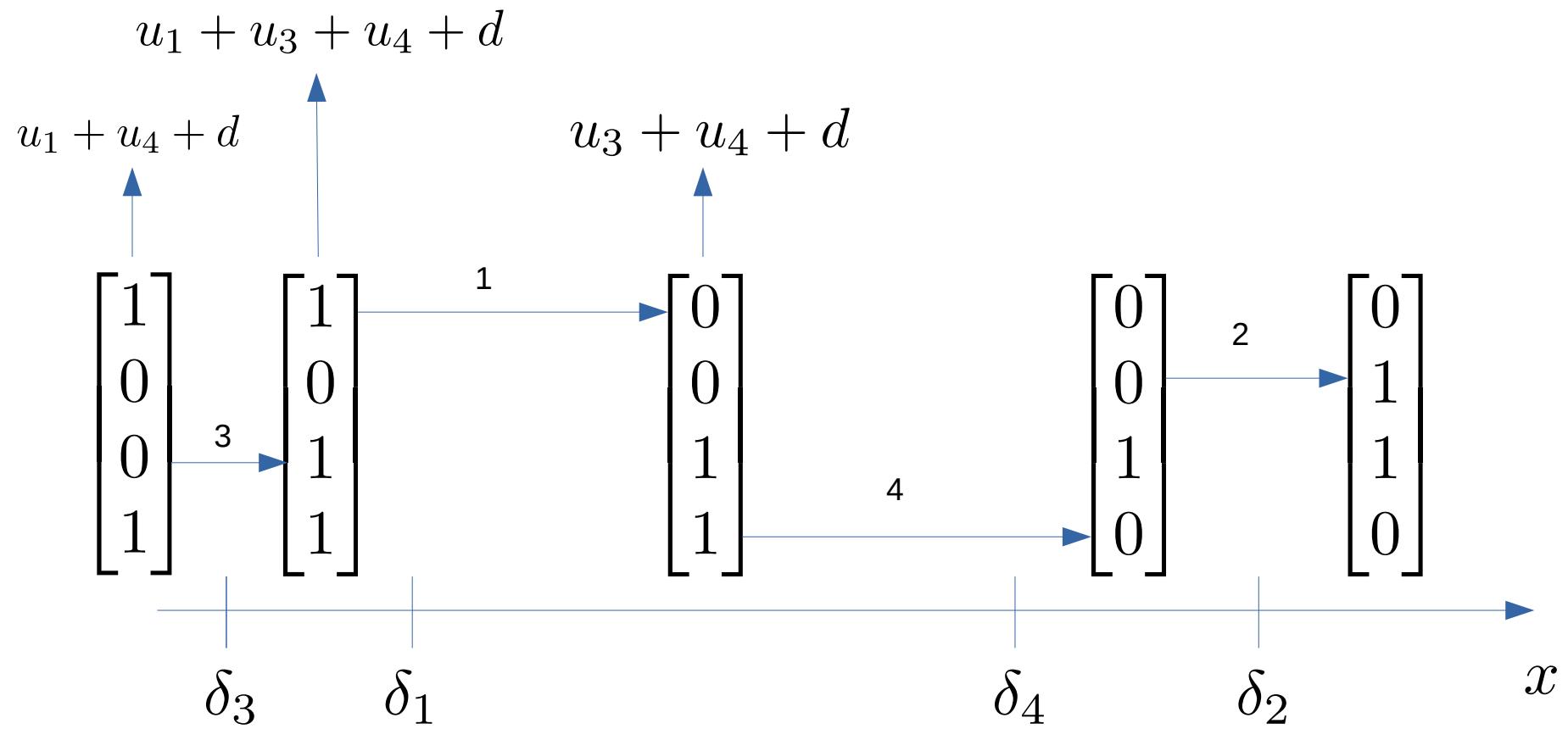


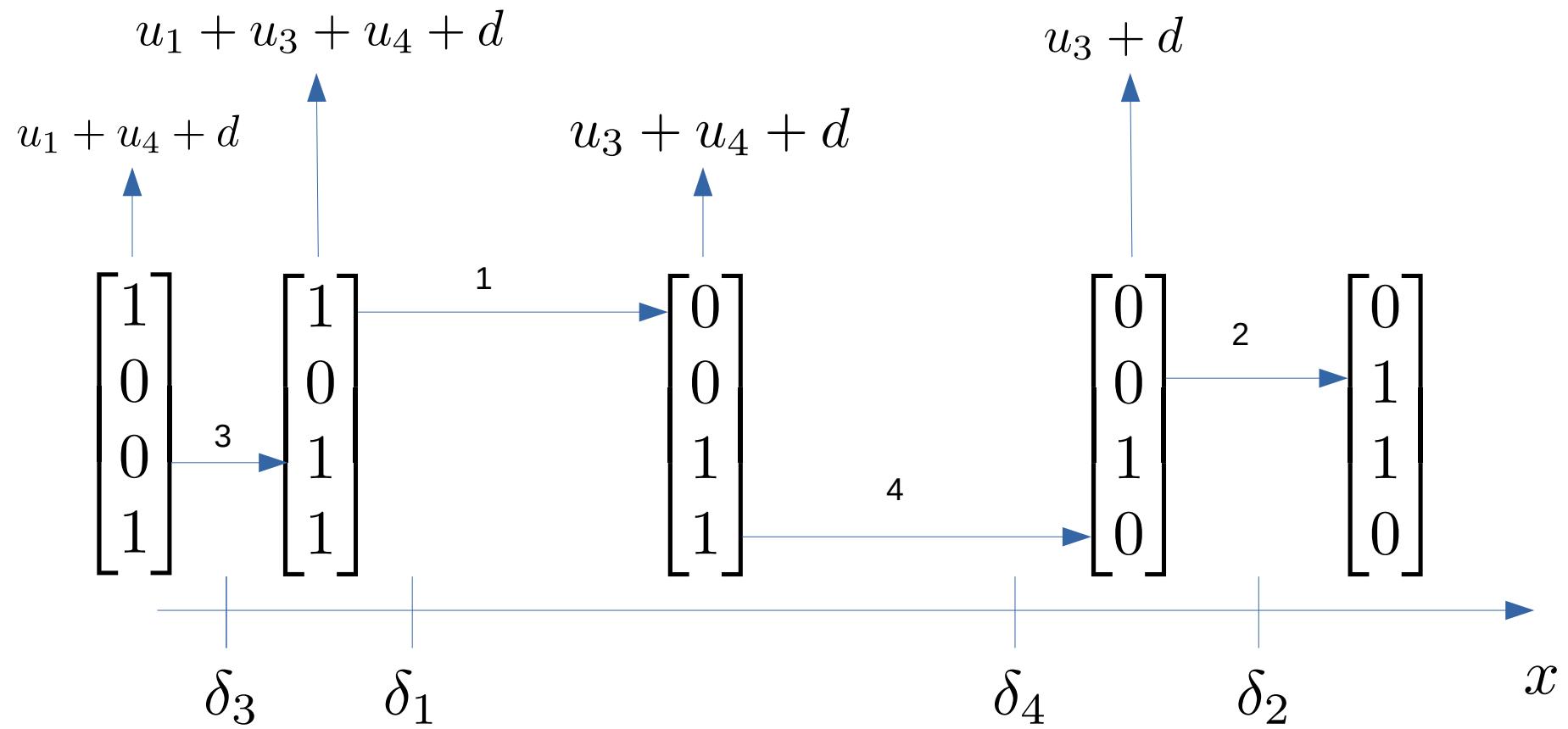


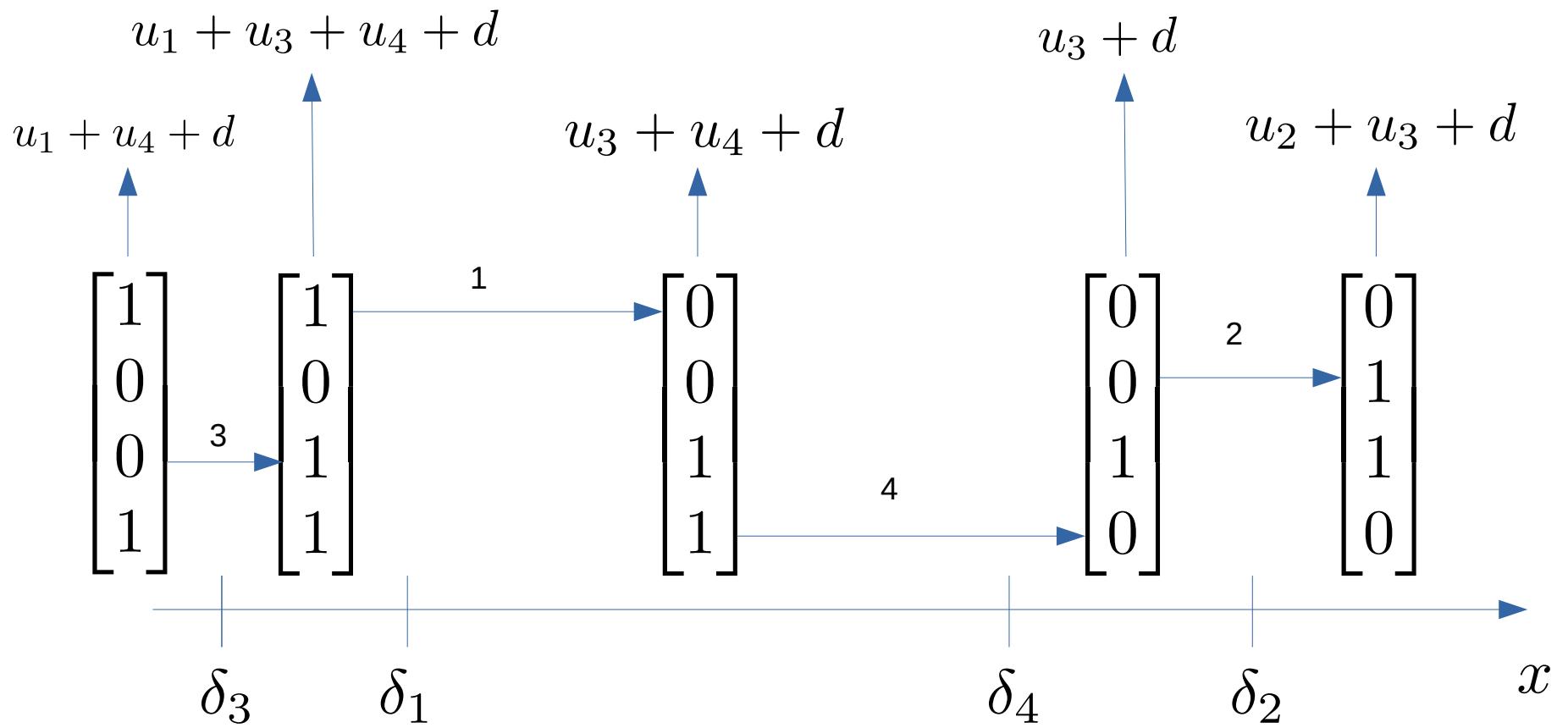
$$\begin{matrix}
 u_1 & \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \\
 u_2 & \\
 u_3 & \\
 u_4 &
 \end{matrix}
 \xrightarrow{\quad 1 * u_1 + 0 * u_2 + 0 * u_3 + 1 * u_4 + d = \boxed{u_1 + u_4 + d} \quad}$$

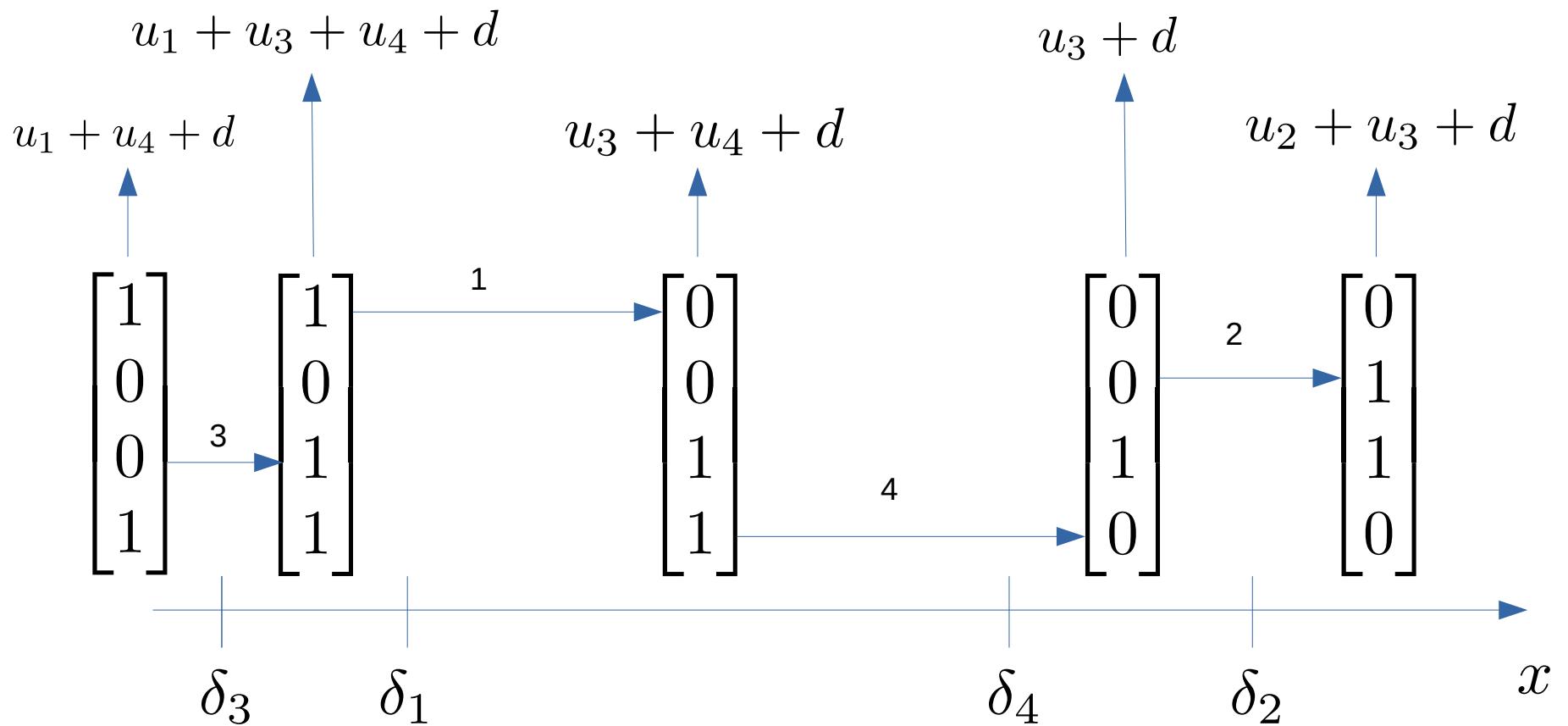






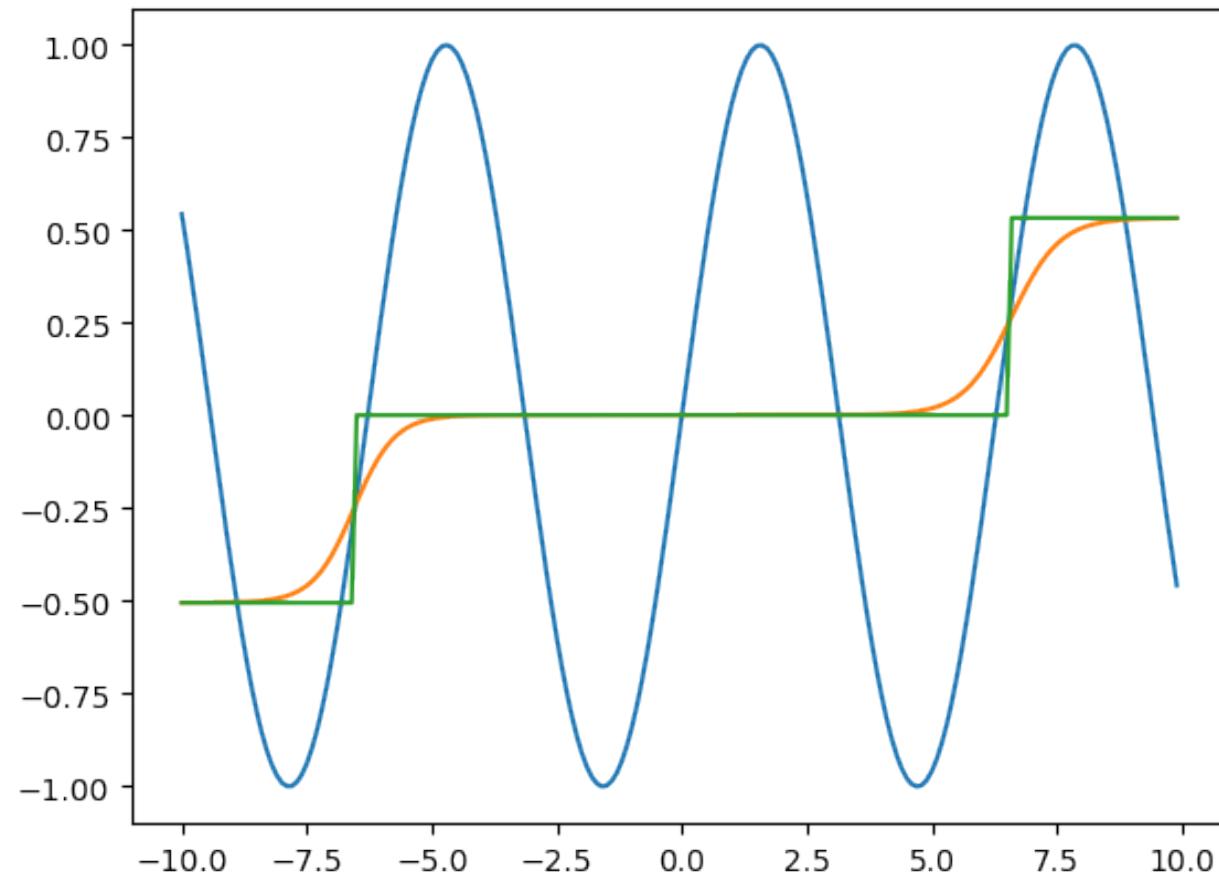




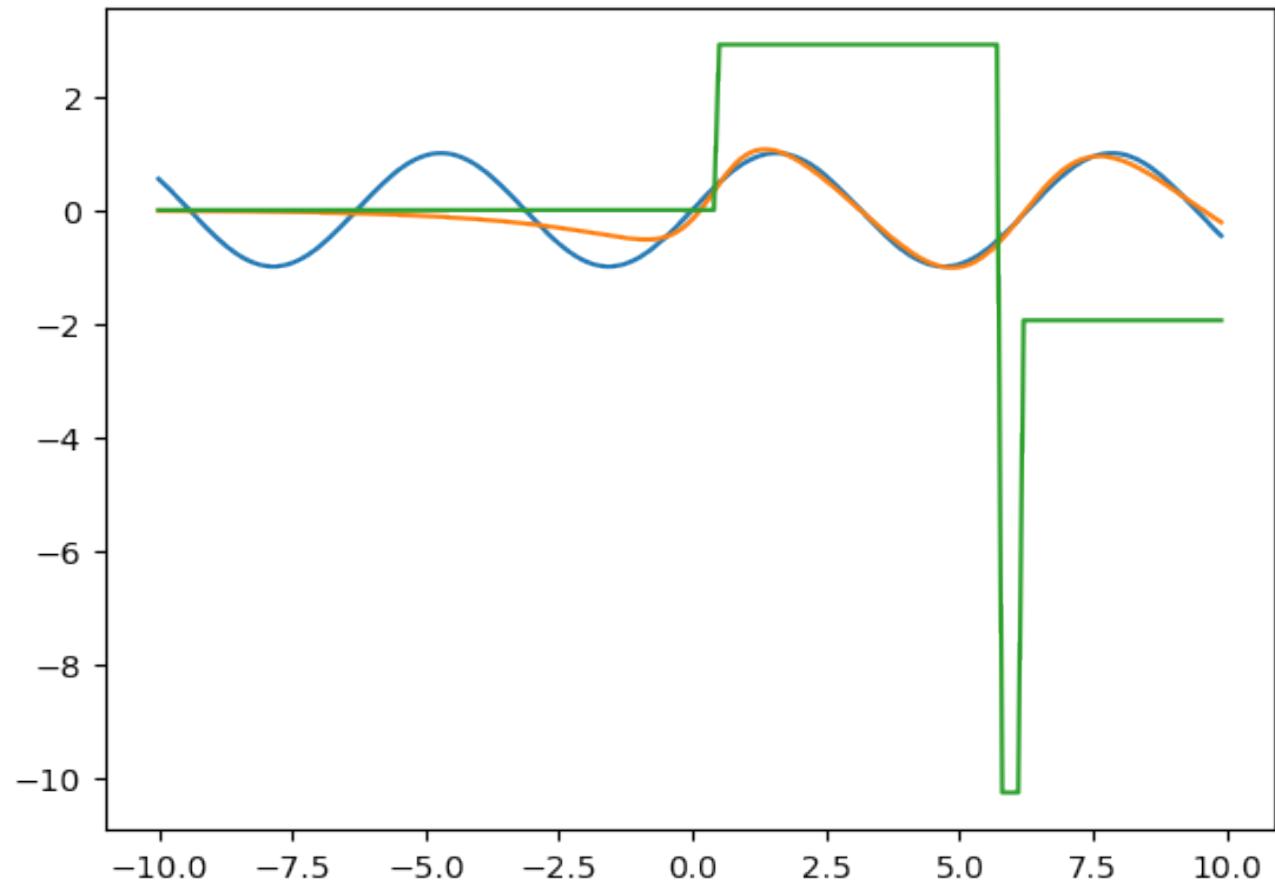


**Five distinct regions where NN predicts constant values  
with sigmoid turn on and turn off**

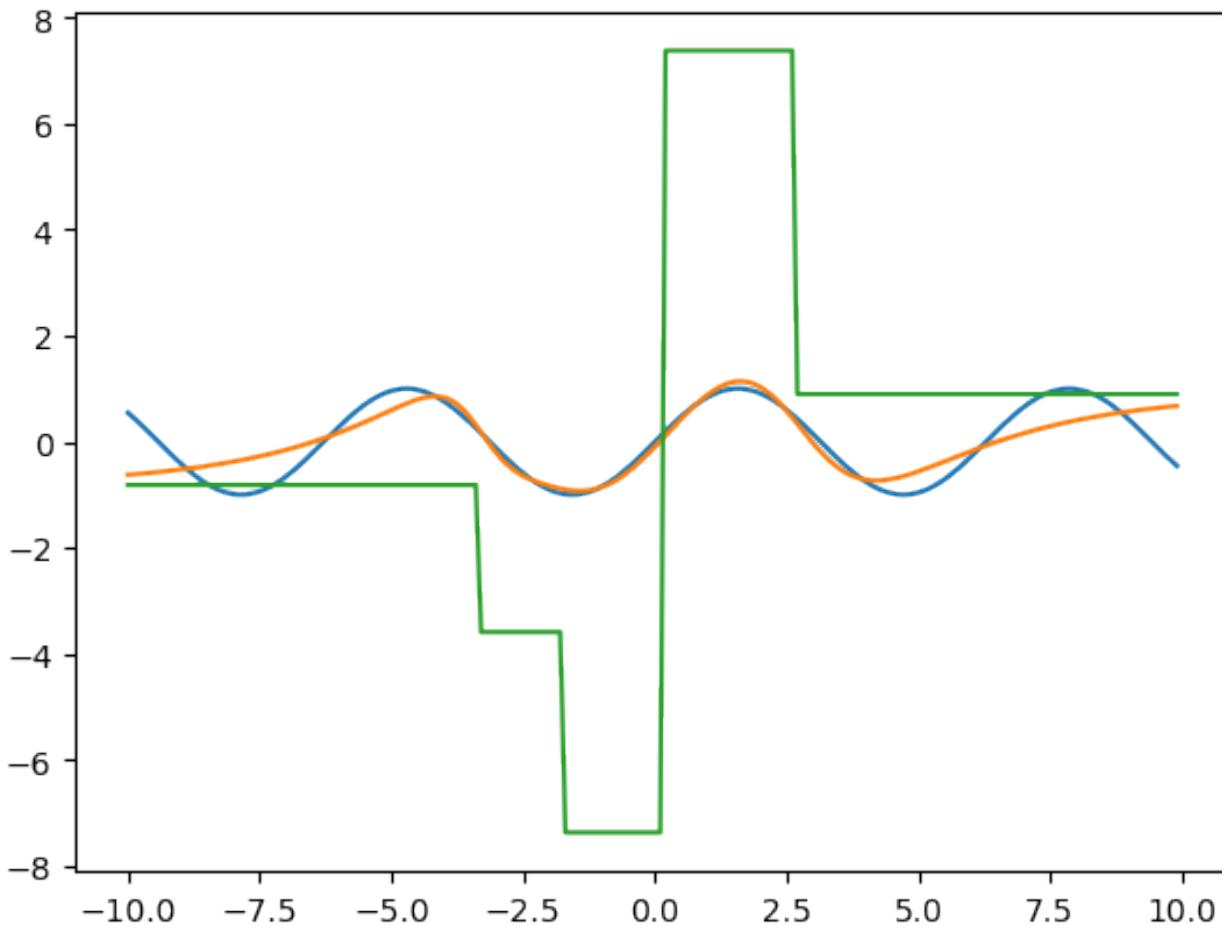
# Hidden Nodes = 2



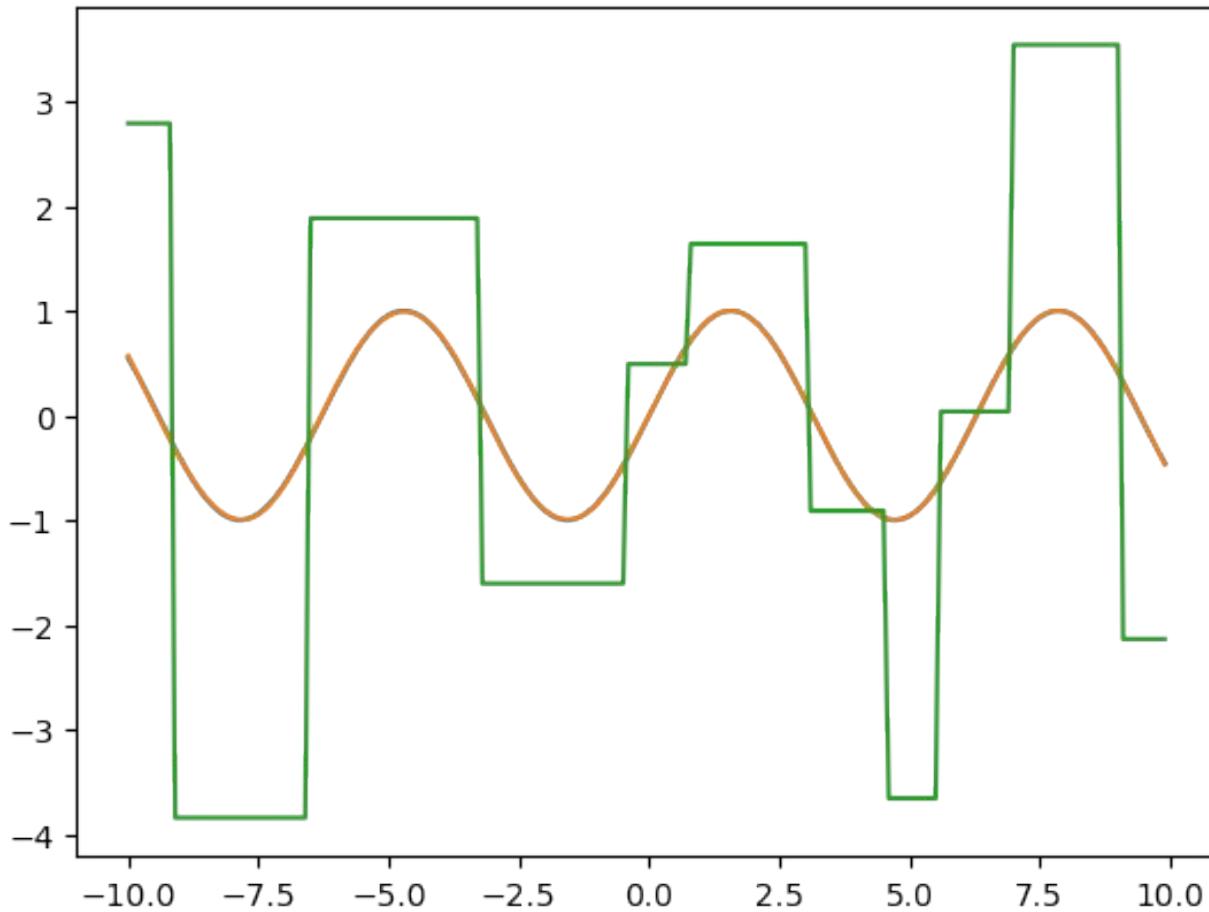
# Hidden Nodes = 3



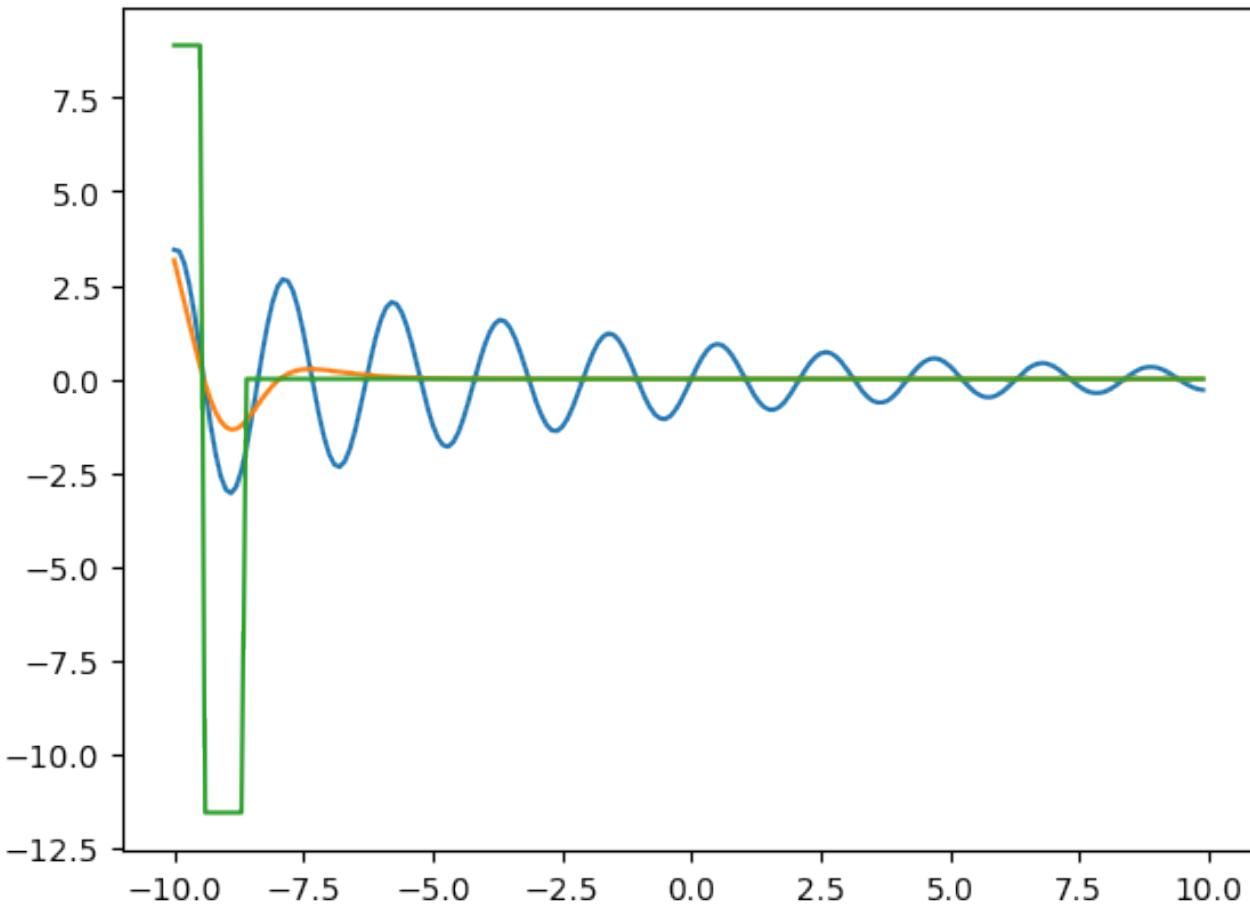
# Hidden Nodes = 4



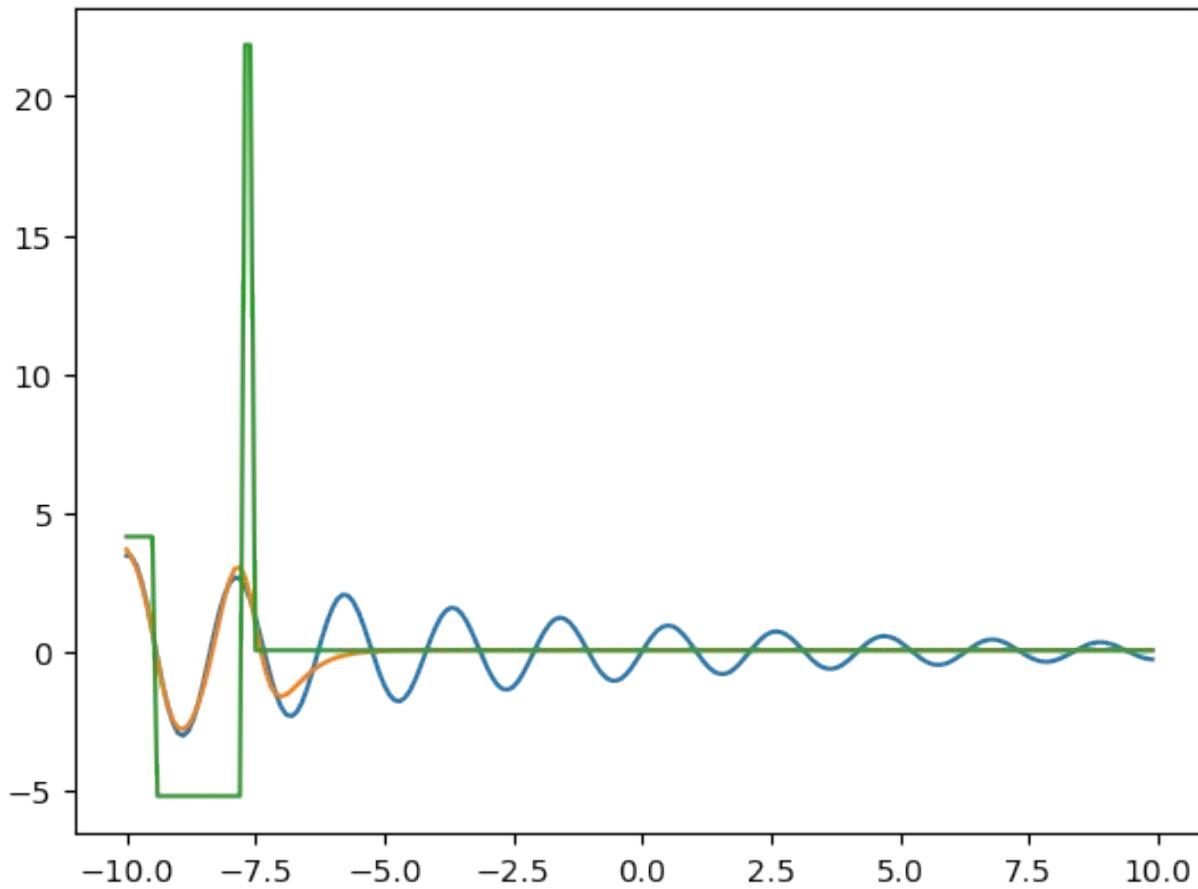
# Hidden Nodes = 10



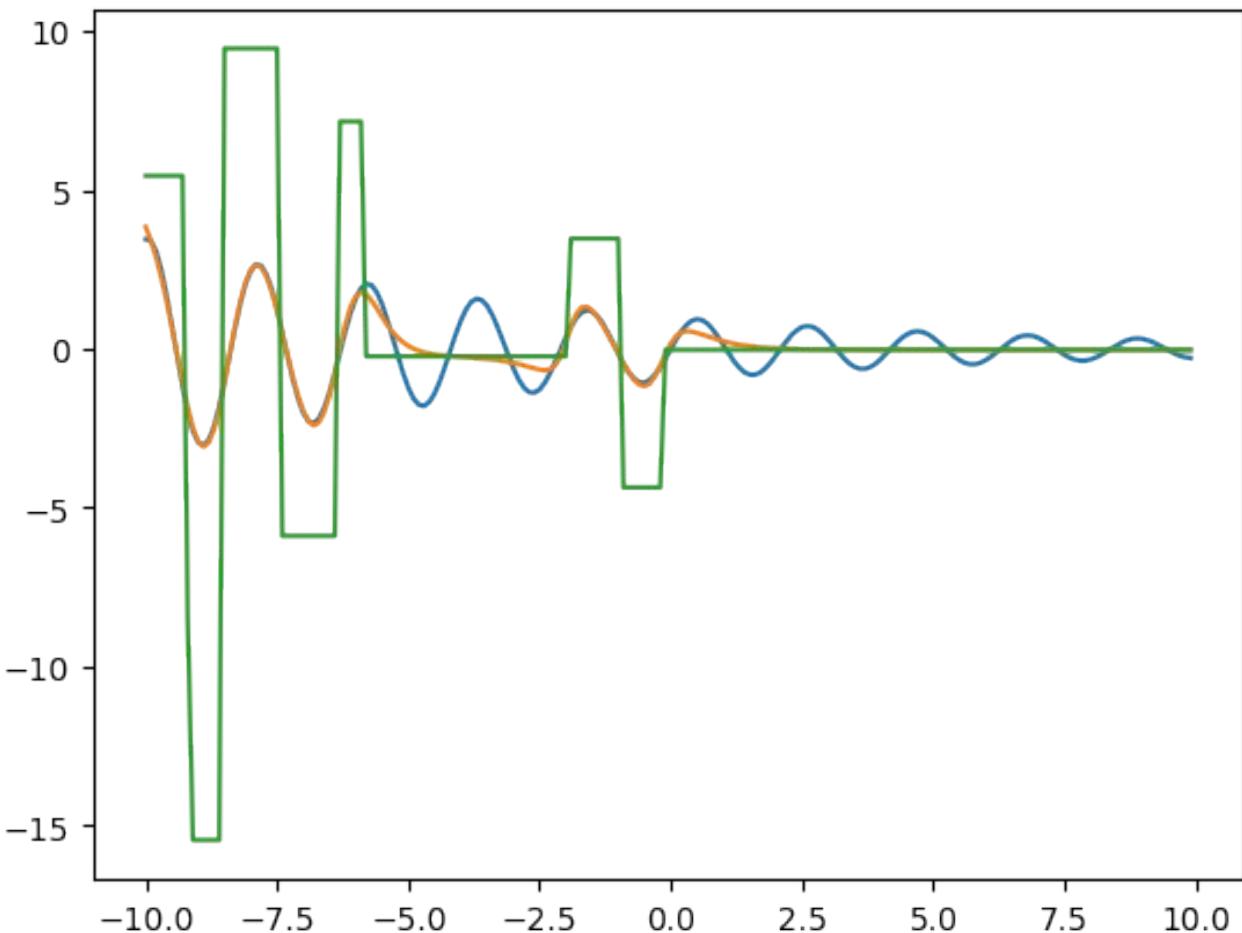
# Hidden Nodes = 2



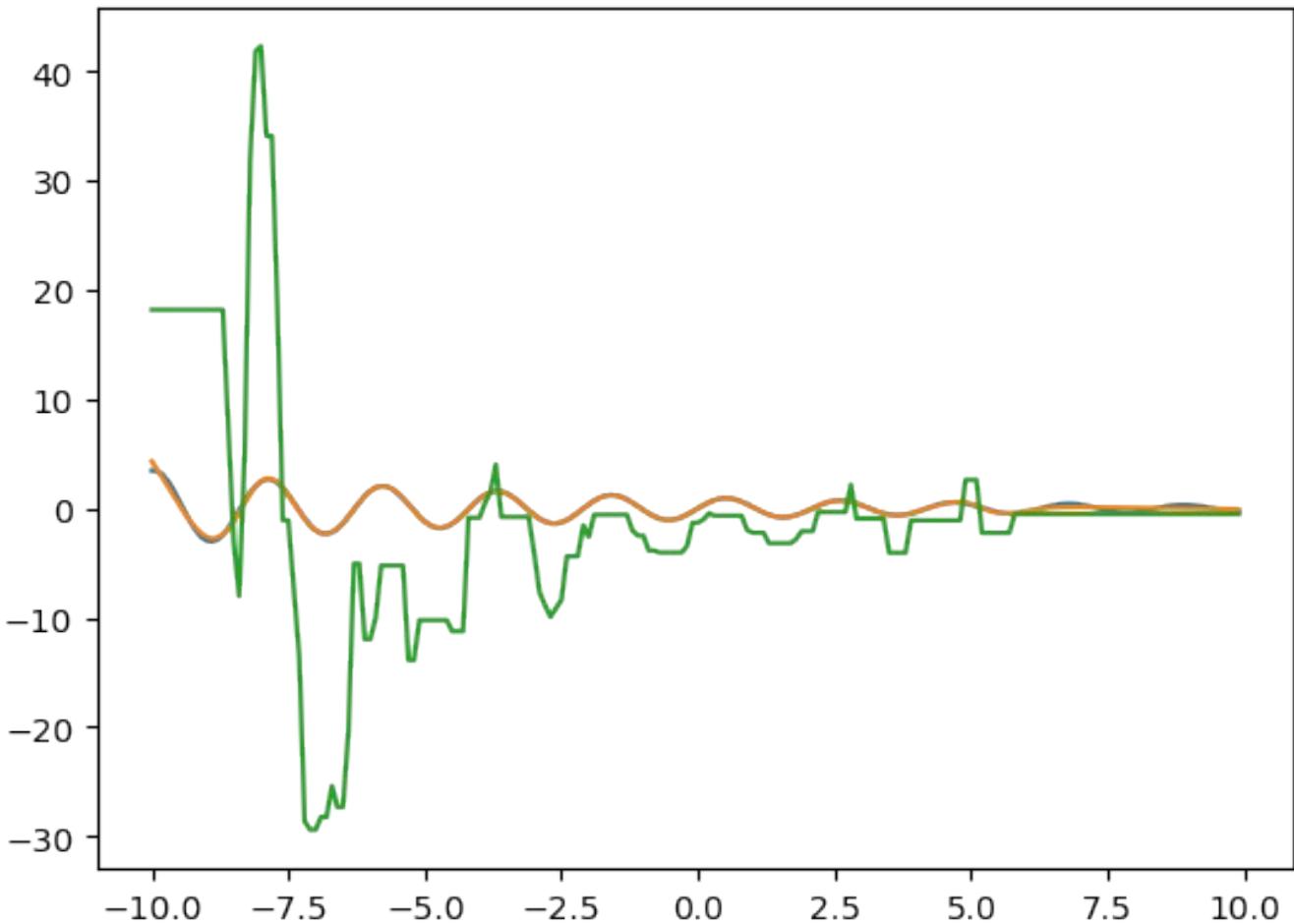
# Hidden Nodes = 3



# Hidden Nodes = 10



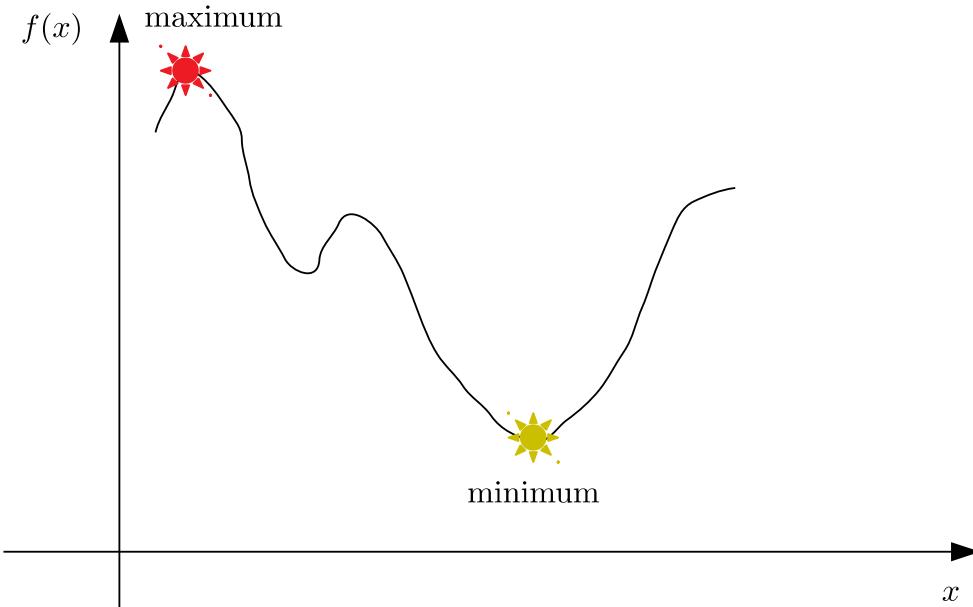
# Hidden Nodes = 100



# Mathematical Optimization

# What is Optimization?

Mathematical and numerical techniques to find the maximum or minimum of a function



# Terminology

Global maximum: the maximum value the function takes across its domain

Local maximum: the maximum value of the function in a small neighborhood around an  $x$  (input value)

Similar definitions for minima

# Terminology

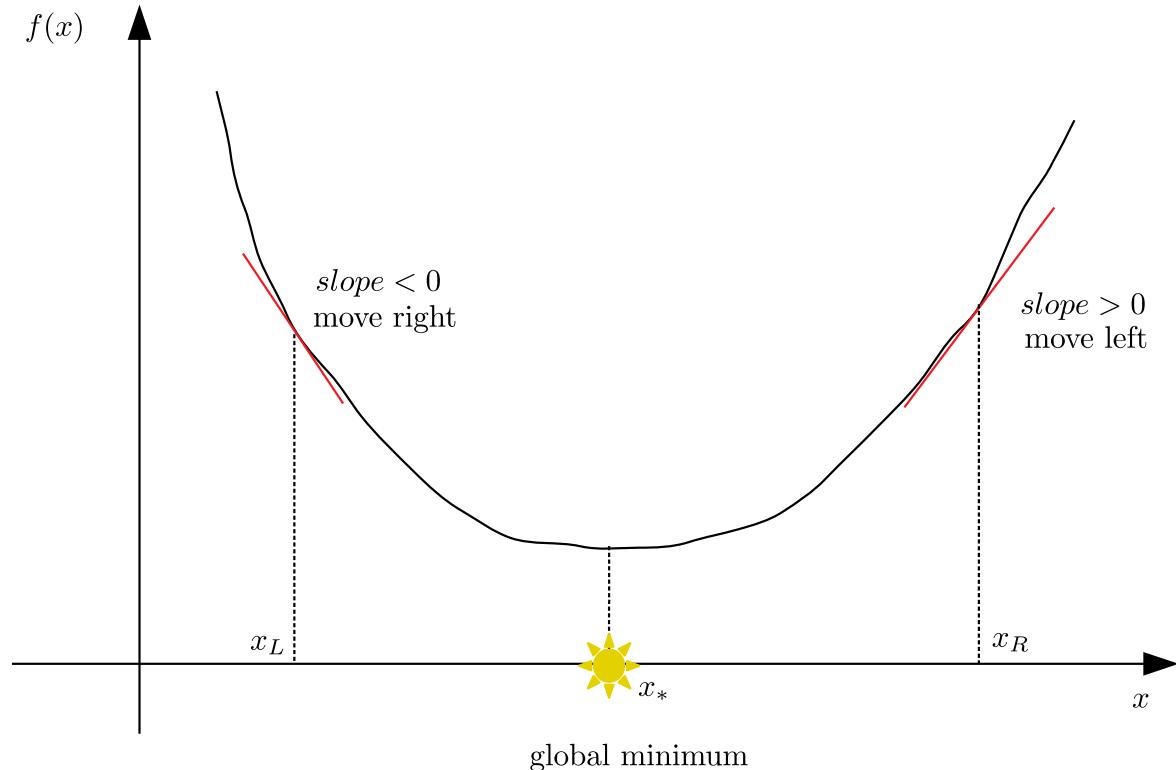
**Training:** Finding weights that minimize the loss function on the test set (“out of sample”) for a neural network is an optimization problem

**Hyperparameter tuning:** Finding the number of layers, the activation function, the number of nodes in each layer, and other parameters is also an optimization problem

# Obstacles

- $f$  might be computationally expensive to evaluate.
- $f$  might be discrete so no way to evaluate derivatives which can guide search for minima
- Even if  $f$  is continuous and well-behaved, evaluating higher derivatives (second, third etc.) is very expensive.
- $f$  might have very complex structure with multiple (possibly infinite) local minima
- $f$  might be very high-dimensional i.e. it has a large number of inputs and hence we are searching for the minima in a high-dimensional space with many more directions to explore.

# Gradient Descent



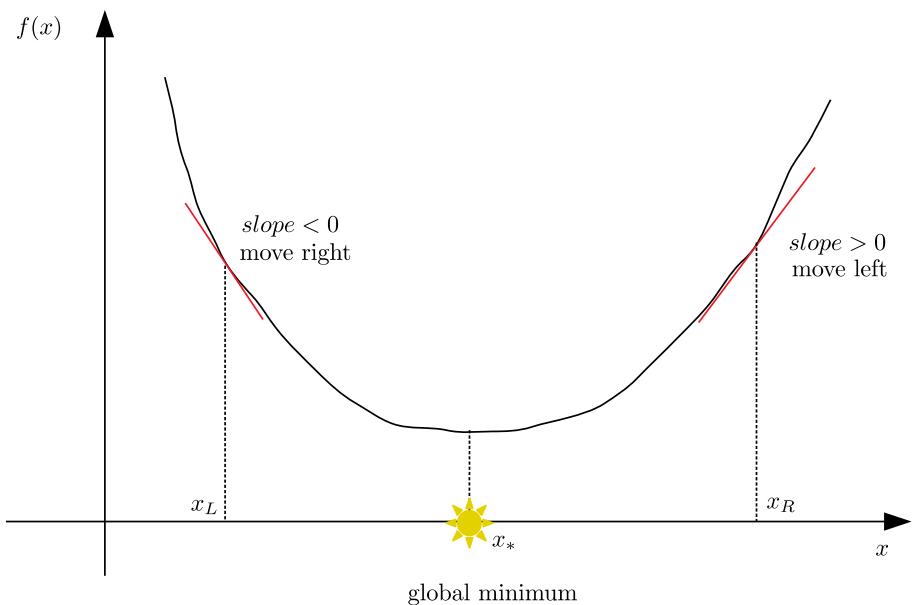
# Gradient Descent

Iterative method

Start at some reasonable point  
and keep updating

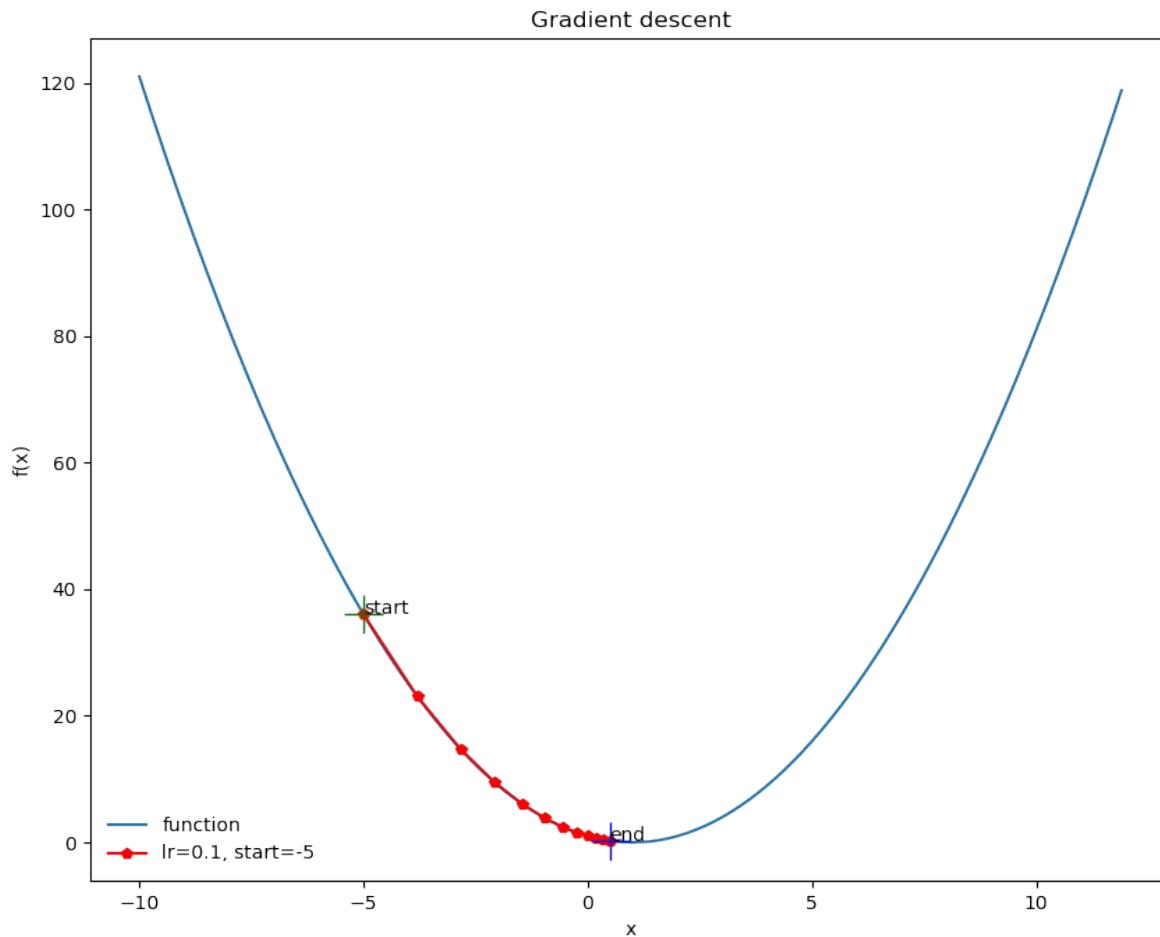
$$x^{(t+1)} = x^{(t)} + \text{update}$$

$$x^{(t+1)} = x^{(t)} - \eta \frac{df}{dx}(x^{(t)})$$

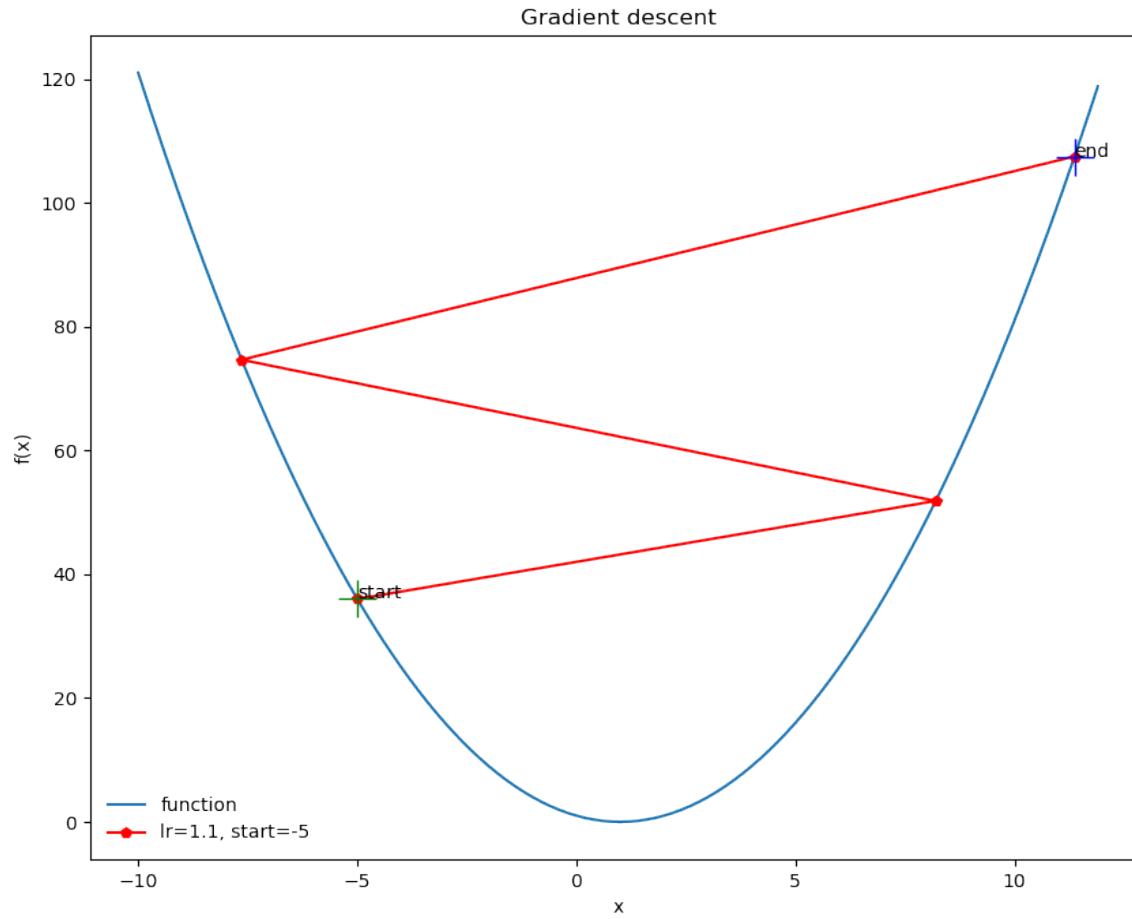


$\eta$  = learning rate

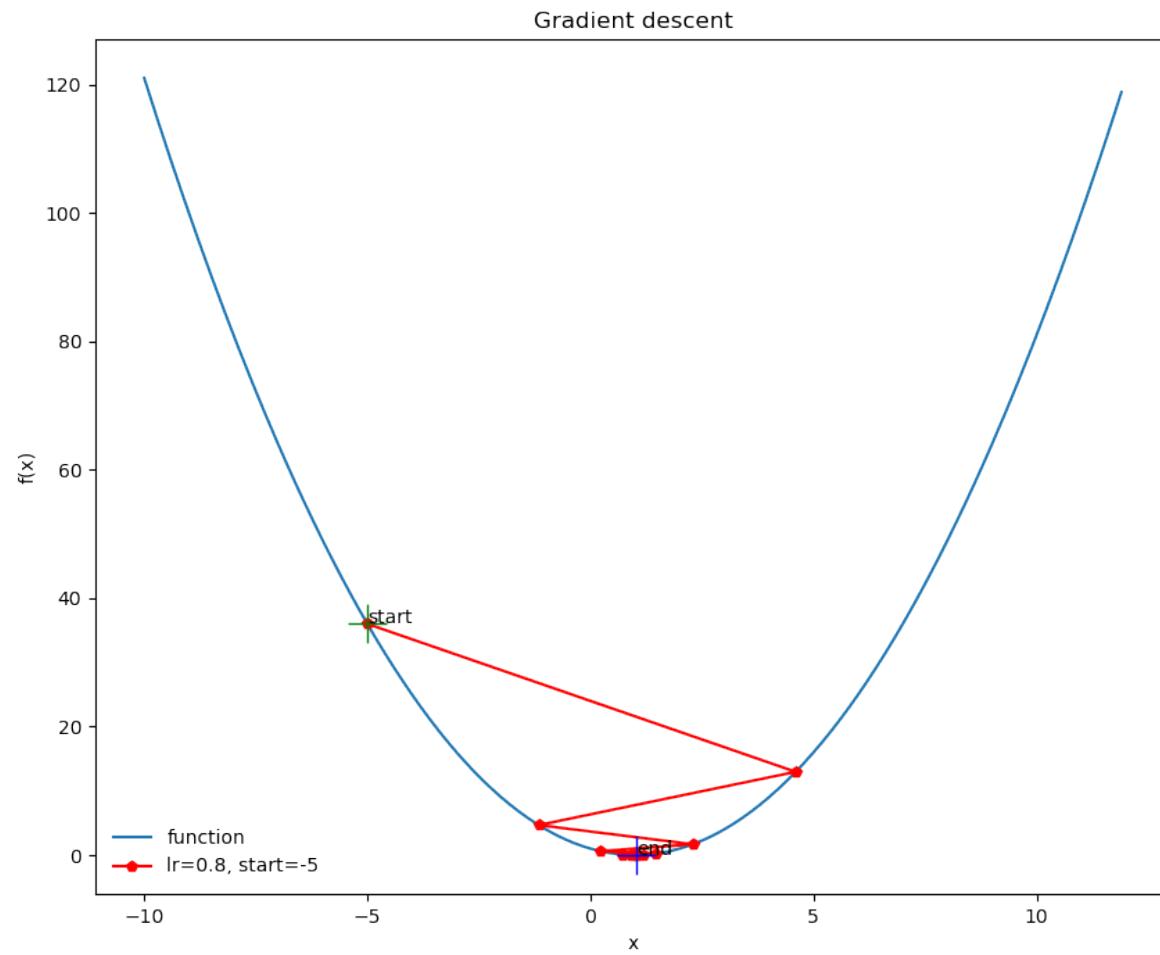
# Intuition



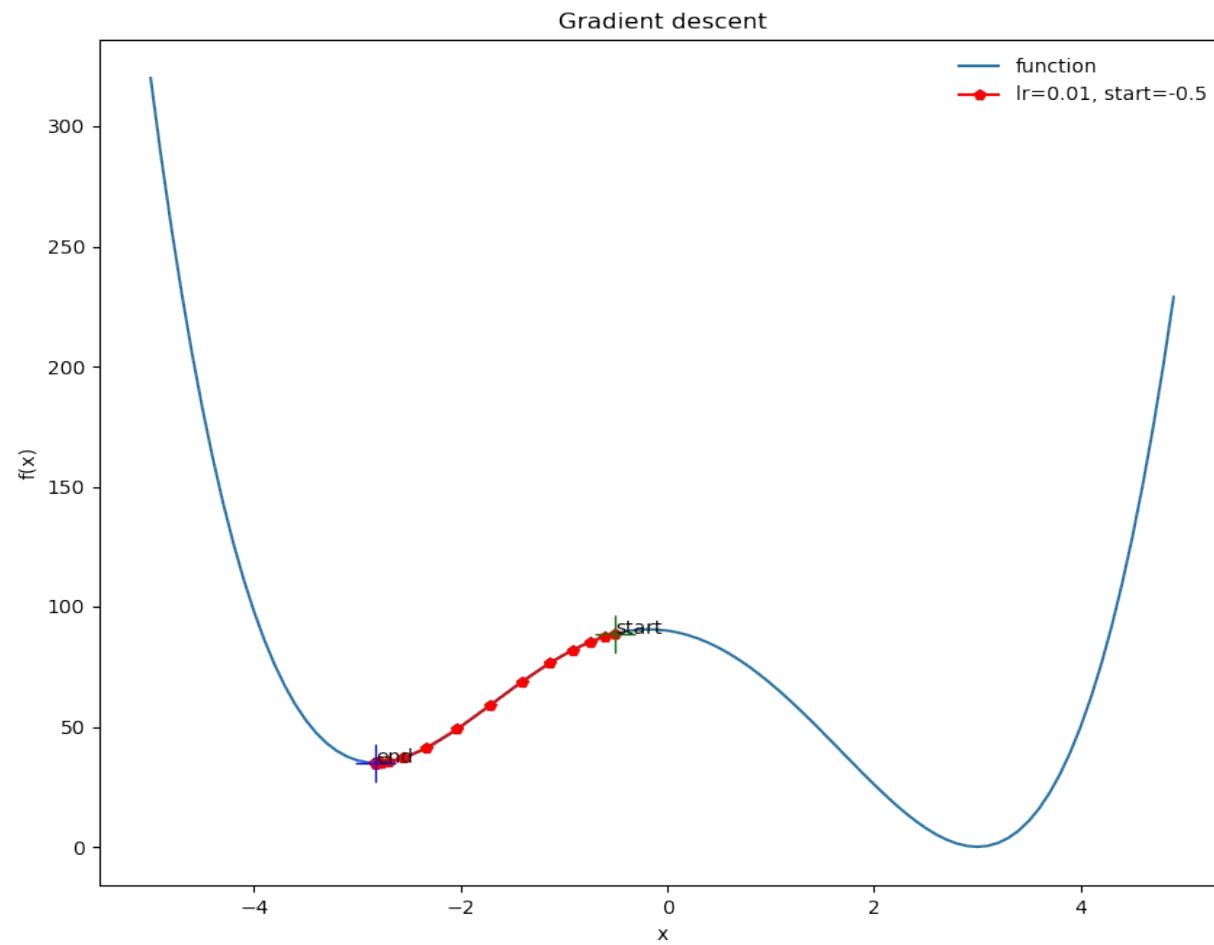
# Intuition



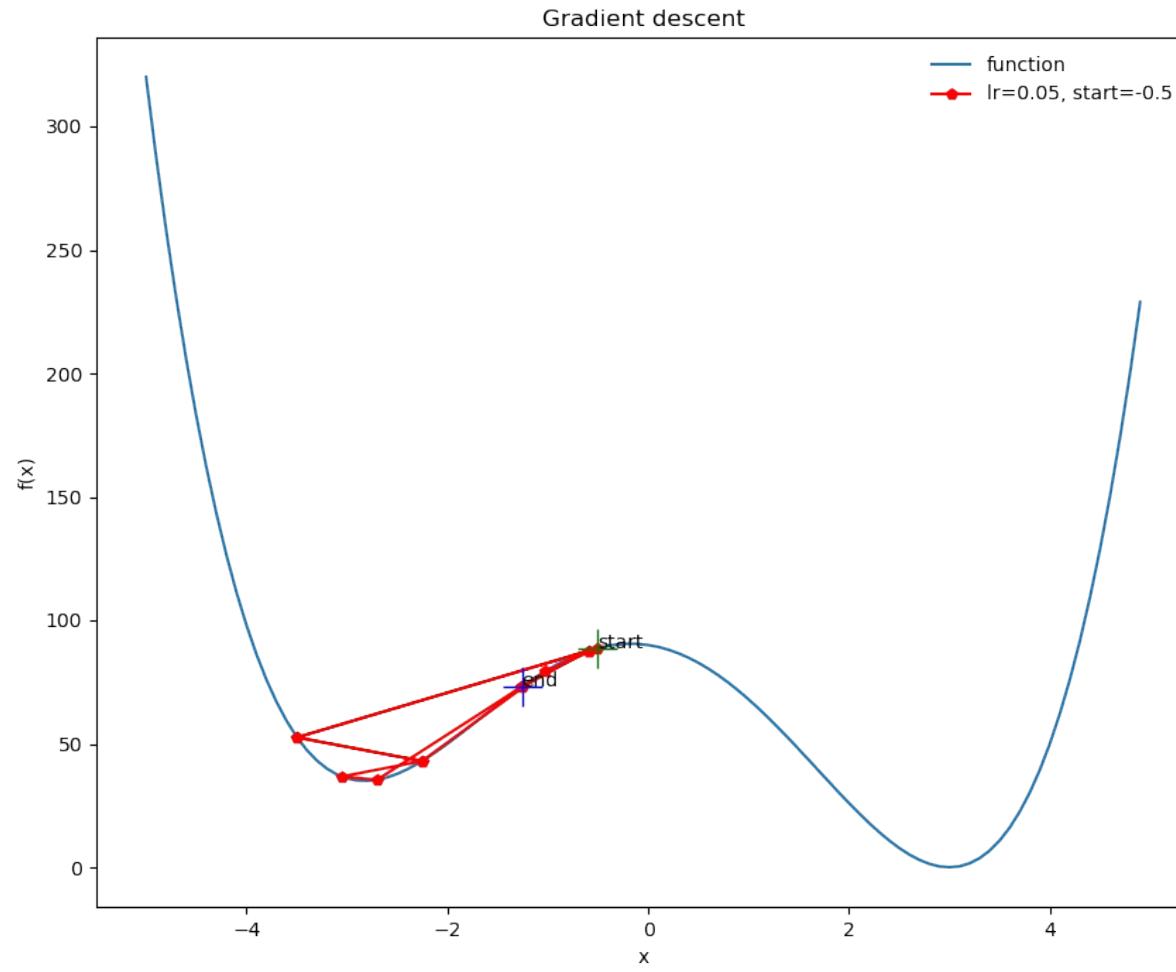
# Intuition



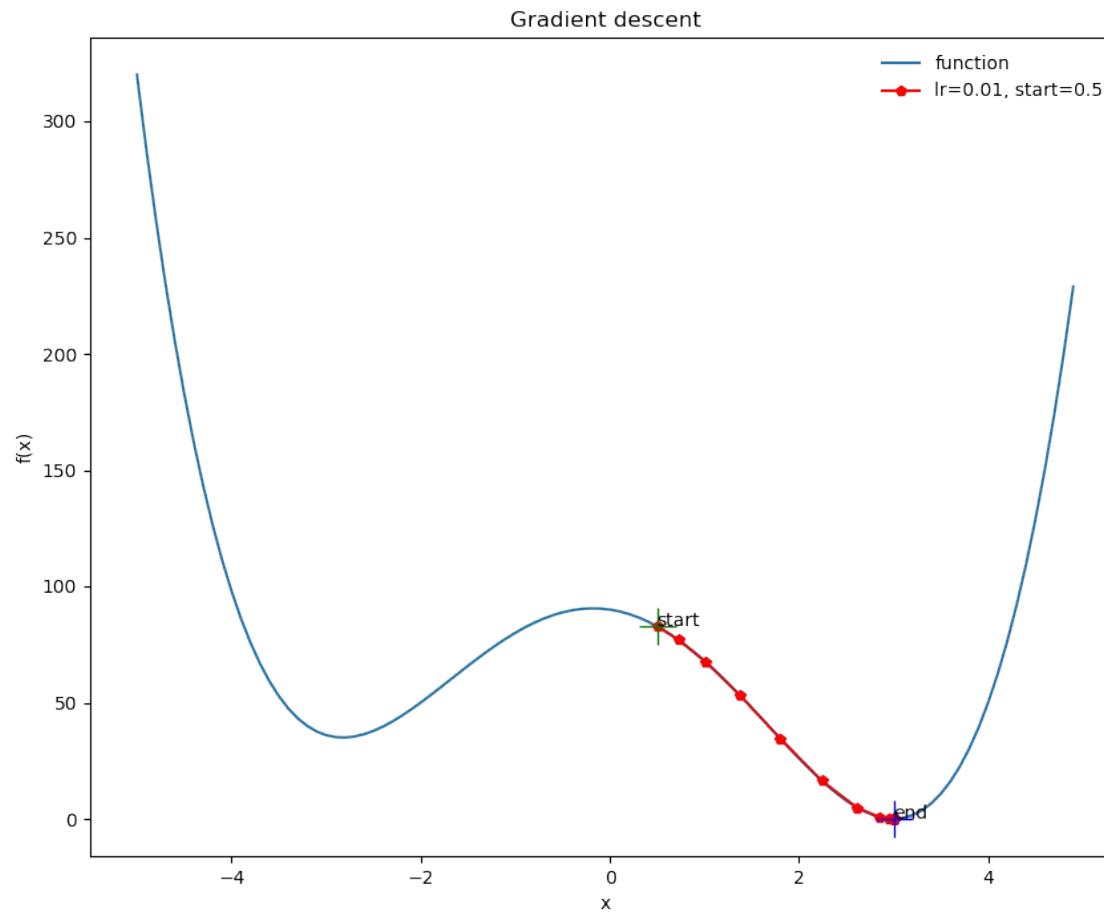
# Intuition



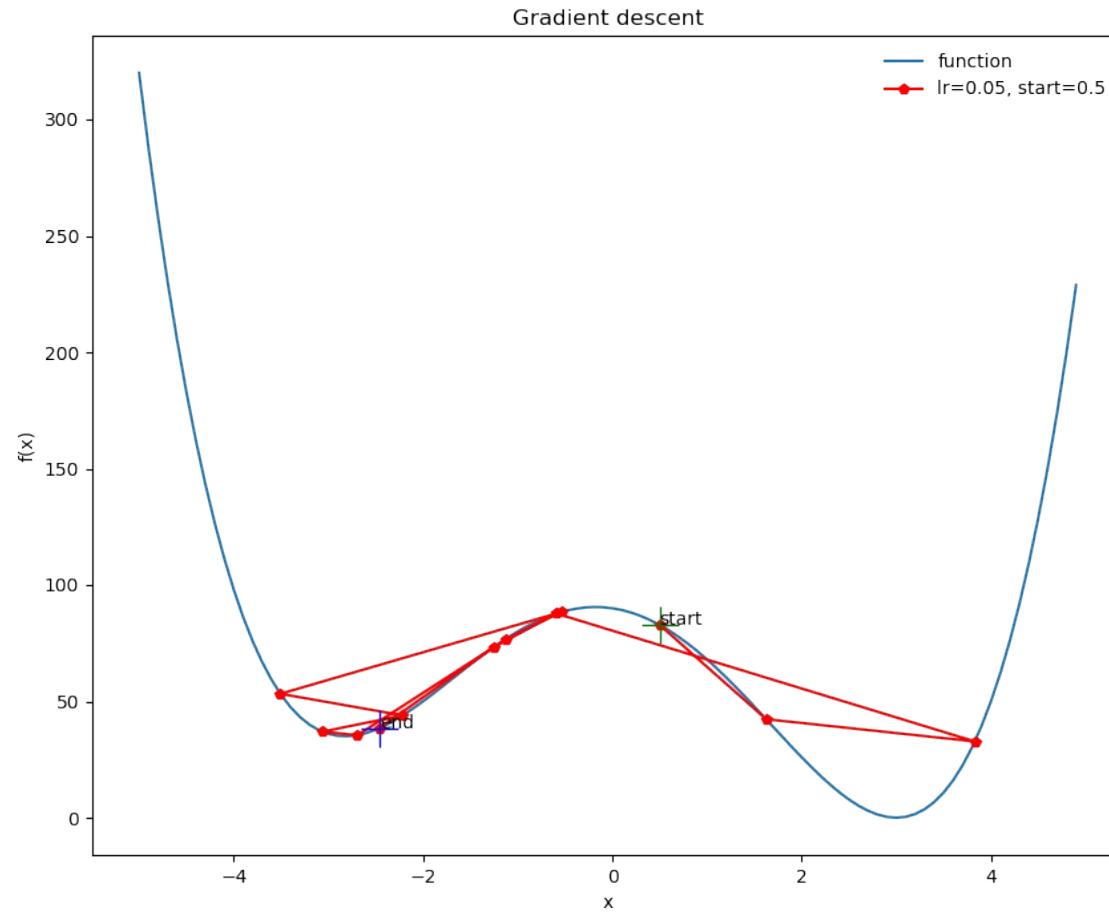
# Intuition



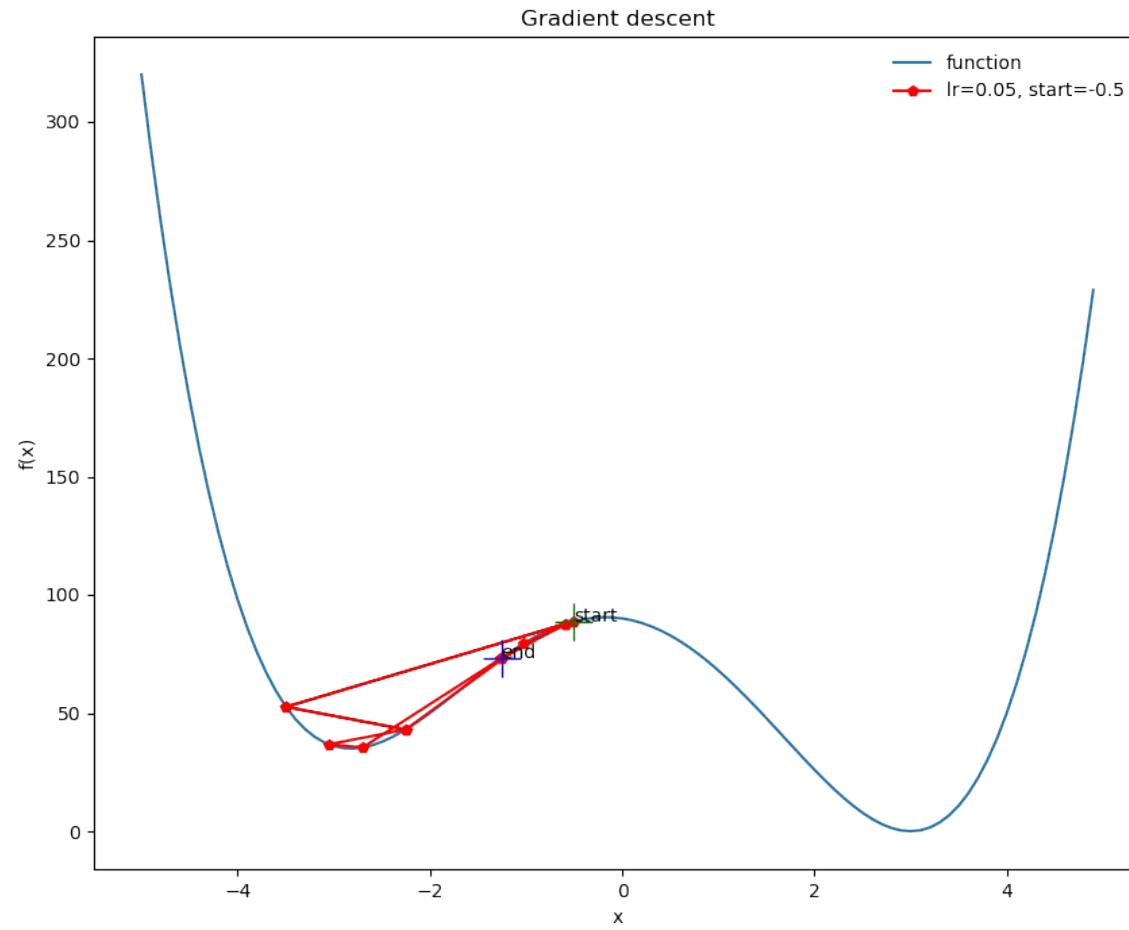
# Intuition



# Intuition



# Intuition



# Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.



[linkedin.com/company/red-hat](https://www.linkedin.com/company/red-hat)



[youtube.com/user/RedHatVideos](https://www.youtube.com/user/RedHatVideos)



[facebook.com/redhatinc](https://www.facebook.com/redhatinc)



[twitter.com/RedHat](https://twitter.com/RedHat)