# CS 412

FEB 25TH – NEURAL NETWORKS

HTF – CHAPTER 11

# Neural Networks

Networks of processing units (neurons) with connections (synapses) between them

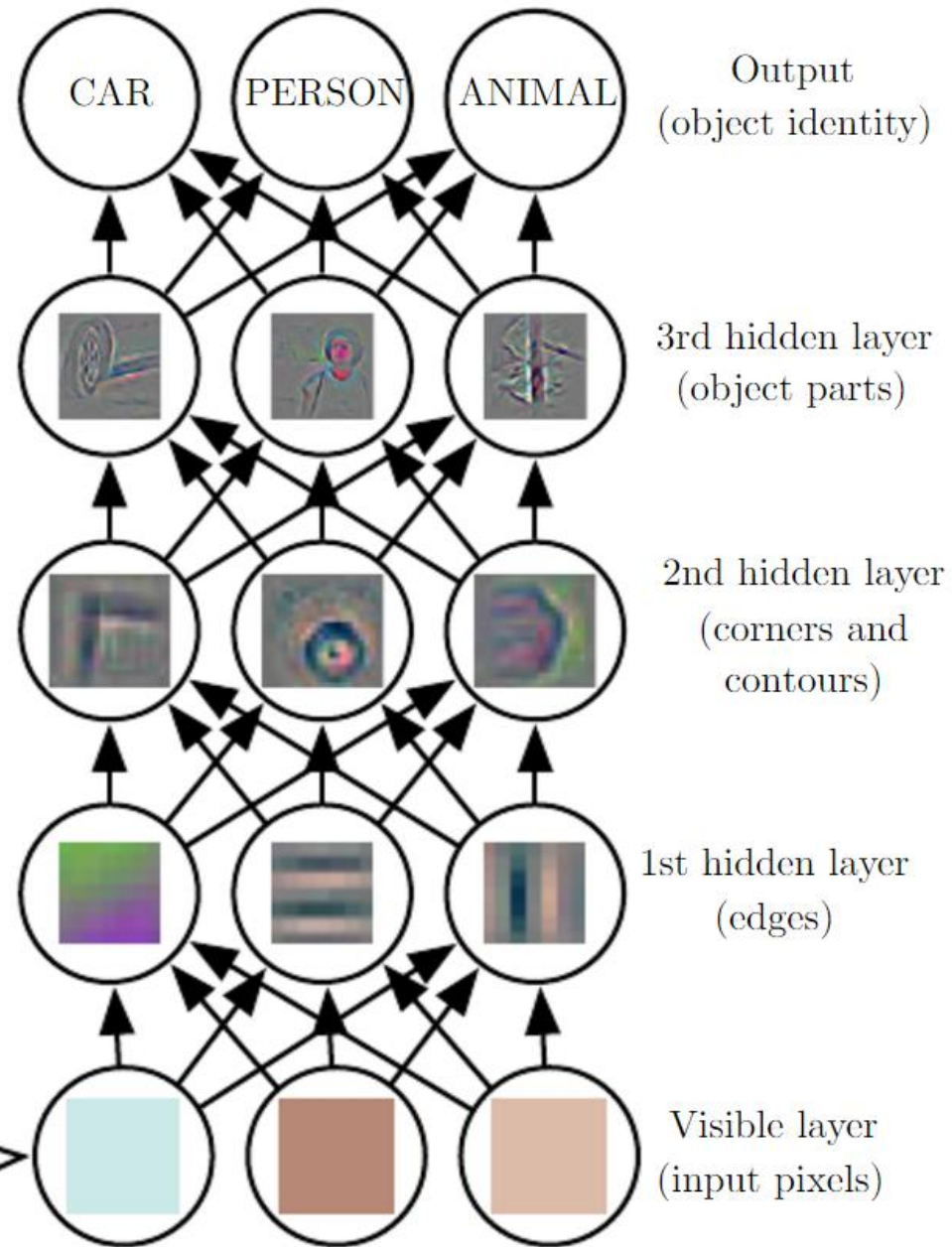Large number of neurons: $10^{10}$

Large connectitivity: $10^5$

Parallel processing

Distributed computation/memory

Robust to noise, failures

Output
(object identity)

3rd hidden layer
(object parts)

2nd hidden layer
(corners and
contours)

1st hidden layer
(edges)

Visible layer
(input pixels)

CAR    PERSON    ANIMAL

# Understanding the Brain

Levels of analysis for an information processing system such as sorting (Marr, 1982)

1. Computational theory: goal of computation and abstract definition of the task
2. Representation and algorithm: how to represent input and output, and how to transform from input to output
3. Hardware implementation

Reverse engineering: From hardware to theory

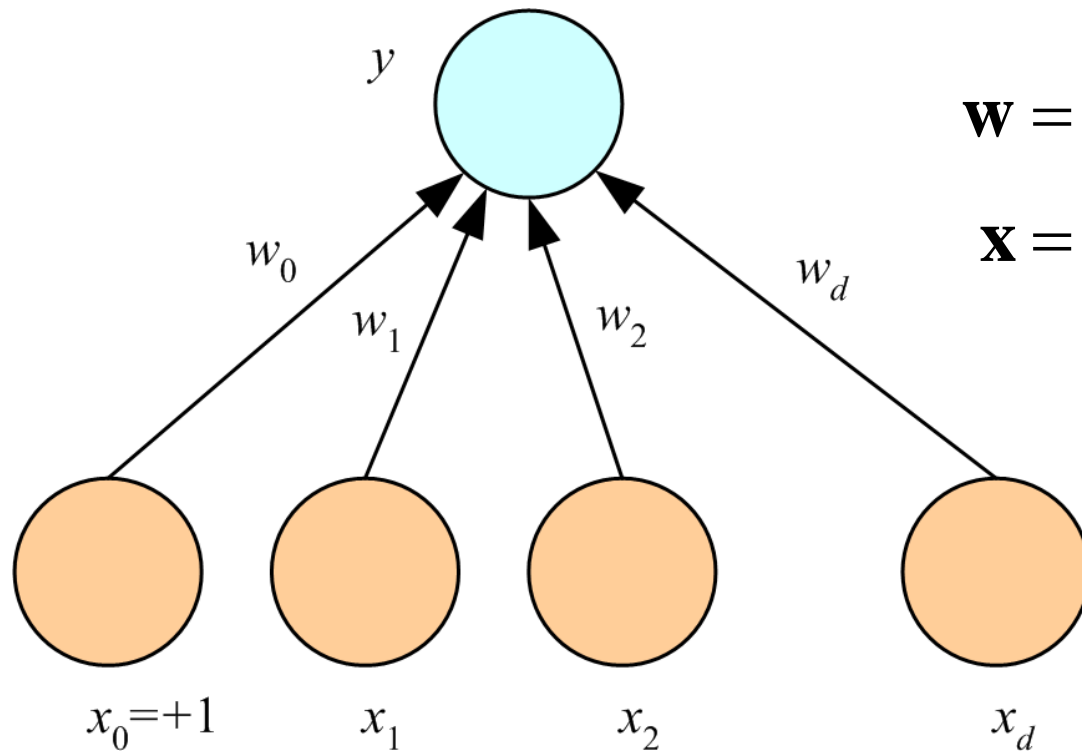Parallel processing: SIMD vs MIMD

Neural net: SIMD with modifiable local memory

Learning: Update by training/experience

GPU ← ↓ ↓ ↳ → CPU

# Perceptron

$$y = \sum_{j=1}^{d} w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = \begin{bmatrix} w_0, w_1, ..., w_d \end{bmatrix}^T$$

$$\mathbf{x} = \begin{bmatrix} 1, x_1, ..., x_d \end{bmatrix}^T$$

(Rosenblatt, 1962)

$y$

$w_0$

$w_1$

$w_2$

$w_d$

$x_0 = +1$

$x_1$

$x_2$

$x_d$

# Perceptron

*"adam"*

What is the single-layer perceptron?
- Just the linear discriminator
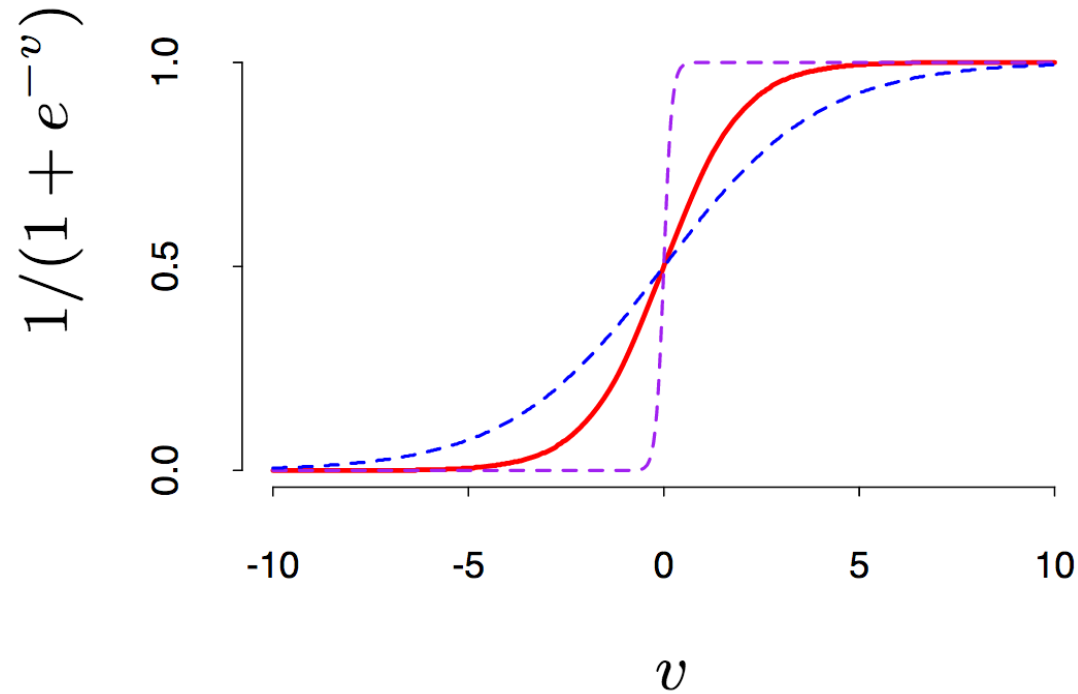- No support vector constraints

*"sgd"*

How do we train it?
- Stochastic gradient descent
  - Small changes based on the data – minimizing loss (what loss should we minimize?)
  - Update = learning factor*(DesiredOutput – Actual Output) * Input

$$\Delta w_{ij}^{t} = \eta \left( r_i^{t} - y_i^{t} \right) x_j^{t}$$

  - Descent is moderated by our learning factor (eta)

# Perceptron



How do we make a regression model into a classification model?

◦ Activation function (here: sigmoid)

◦ Like logistic regression, there is no unique solution, so we also have to consider the rate at which the sigmoid transitions, this is the **activation rate**, s (here: ½,1,10)
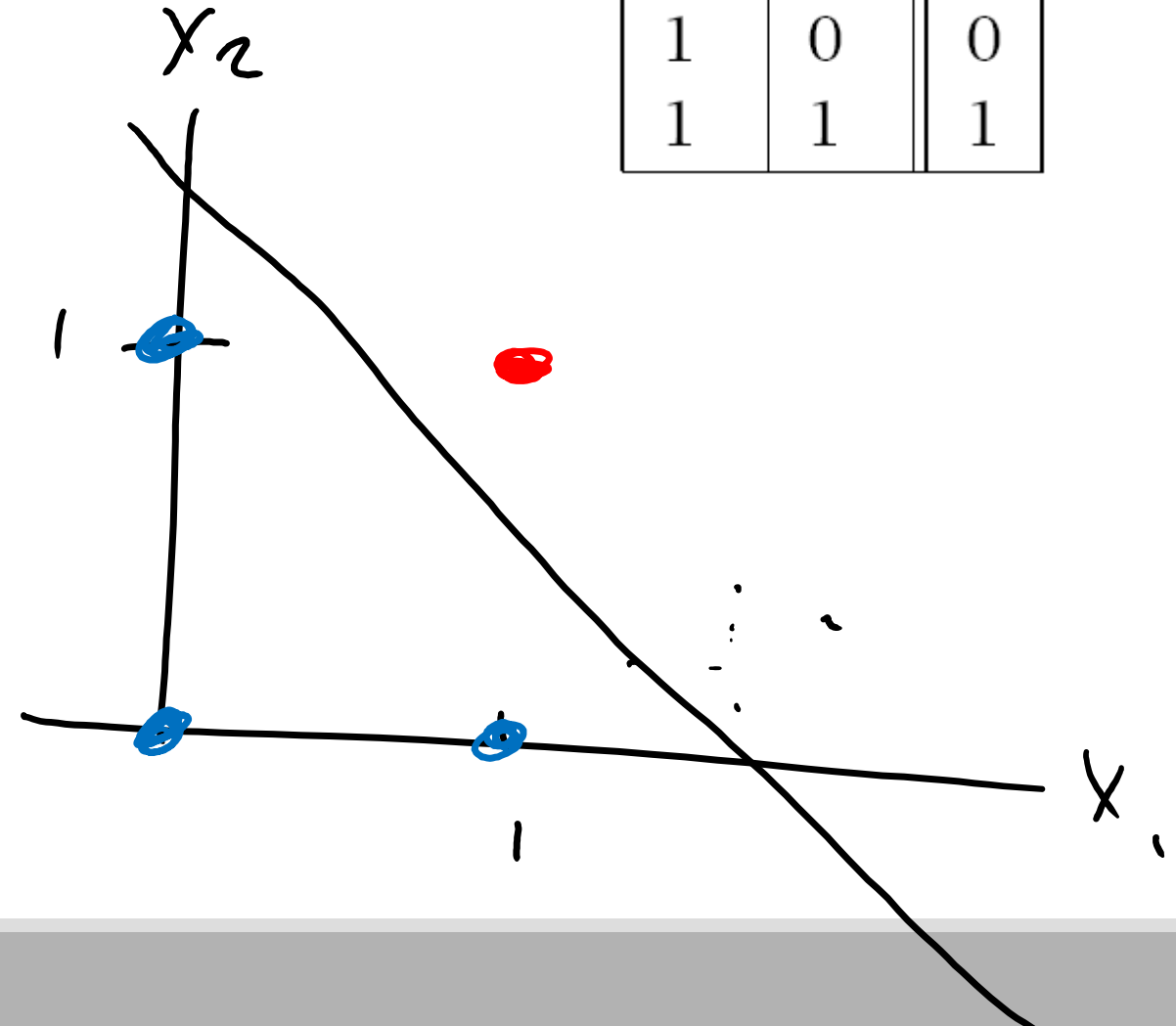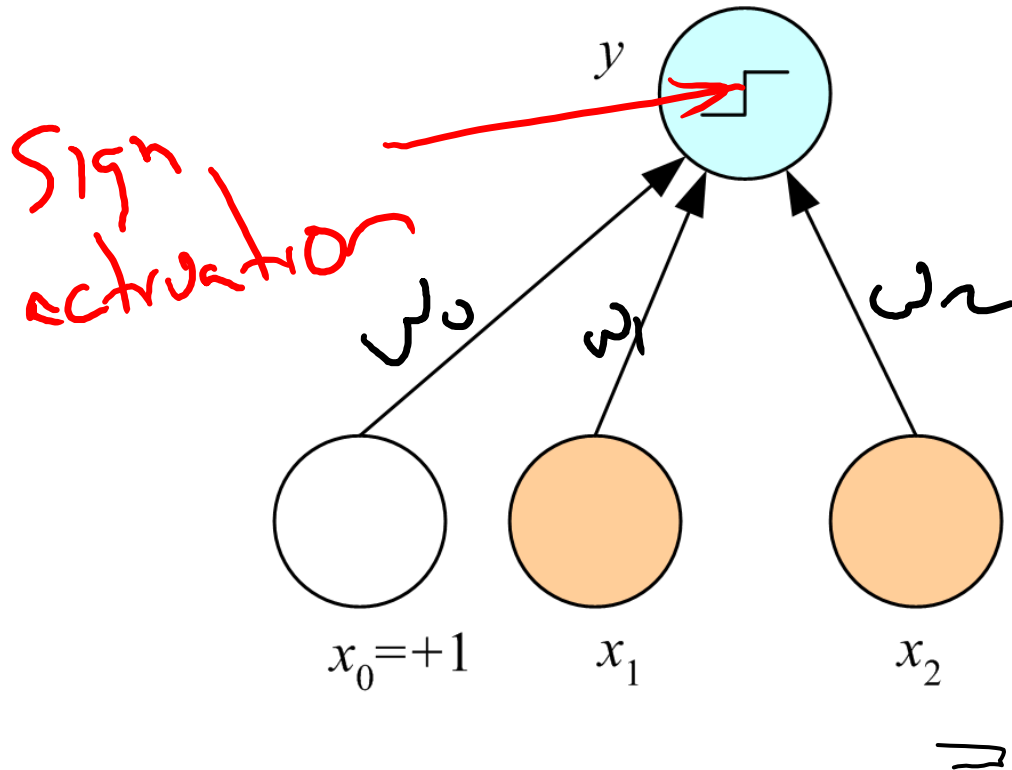
Why do we prefer this to the sign function?

◦ They are differentiable and non-linear

# Learning Boolean AND

| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

What are the weights for this perceptron?
(purple line)

# Learning Boolean AND

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

What are the weights for this perceptron?
(purple line)



$w_0 + 1.5 = 0$

$w_2 = ?$

$+1$ $+1$

$x_0 = +1$ $x_1$ $x_2$

# Learning Boolean AND

| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 0   |
| 1     | 0     | 0   |
| 1     | 1     | 1   |

What are the weights for this perceptron?
       (purple line)

# XOR

| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# XOR

| $x_1$ | $x_2$ | $r$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

No $w_0$, $w_1$, $w_2$ satisfy:

$$w_0 \leq 0$$
$$w_2 + w_0 > 0$$
$$w_1 + w_0 > 0$$
$$w_1 + w_2 + w_0 \leq 0$$

# Perceptron

What is the problem with the simple perceptron?

# Perceptron

What is the problem with the simple perceptron?
  ◦ It can't model non-linear data

How do we fix this?

# Perceptron

What is the problem with the simple perceptron?
◦ It can't model non-linear data


How do we fix this?
◦ SVM fixed this by using the kernel methods
◦ Can the perceptron?

# Perceptron

What is the problem with the simple perceptron?
- ◦ It can't model non-linear data


How do we fix this?
- ◦ SVM fixed this by using the kernel methods
- ◦ Can the perceptron? **Yes**, but that's not what the neural network does

# Perceptron
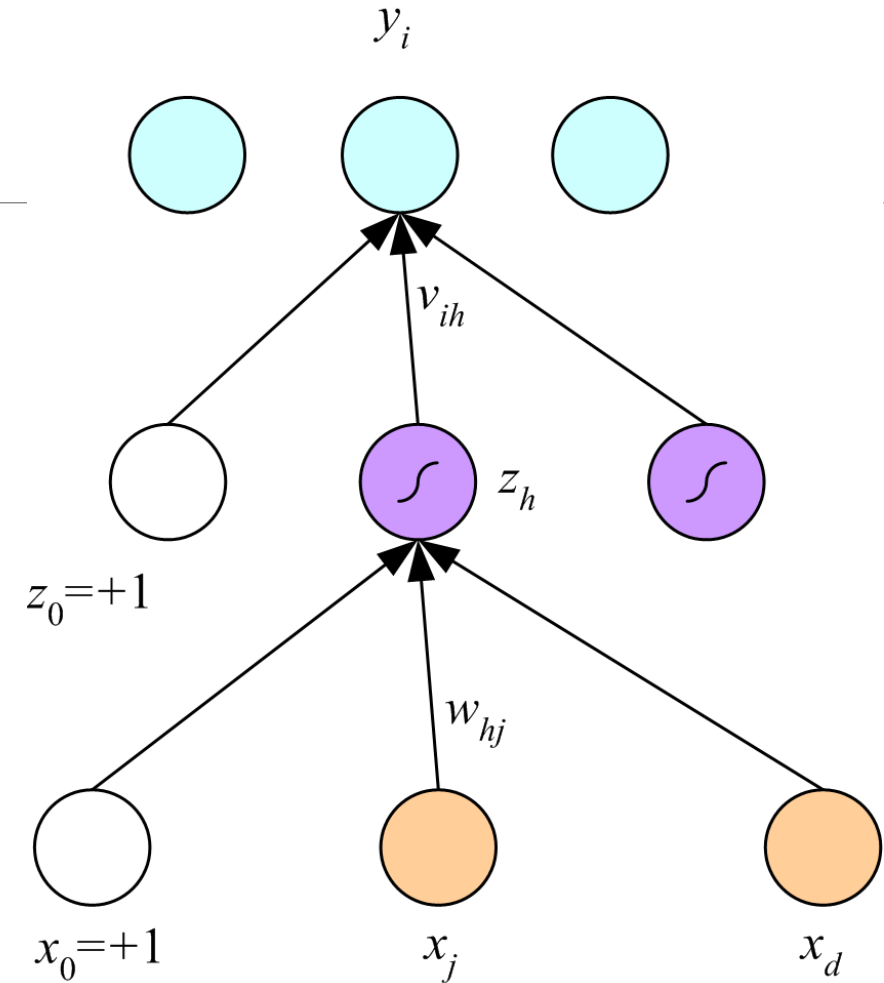
What is the problem with the simple perceptron?

- It can't model non-linear data

How do we fix this?

- SVM fixed this by using the kernel methods
- Can the perceptron? **Yes**, but that's not what the neural network does

Let's add multiple layers to the perceptron

- At each level we have a **regression** model defined by the activation function and **always** a constant $w_0$
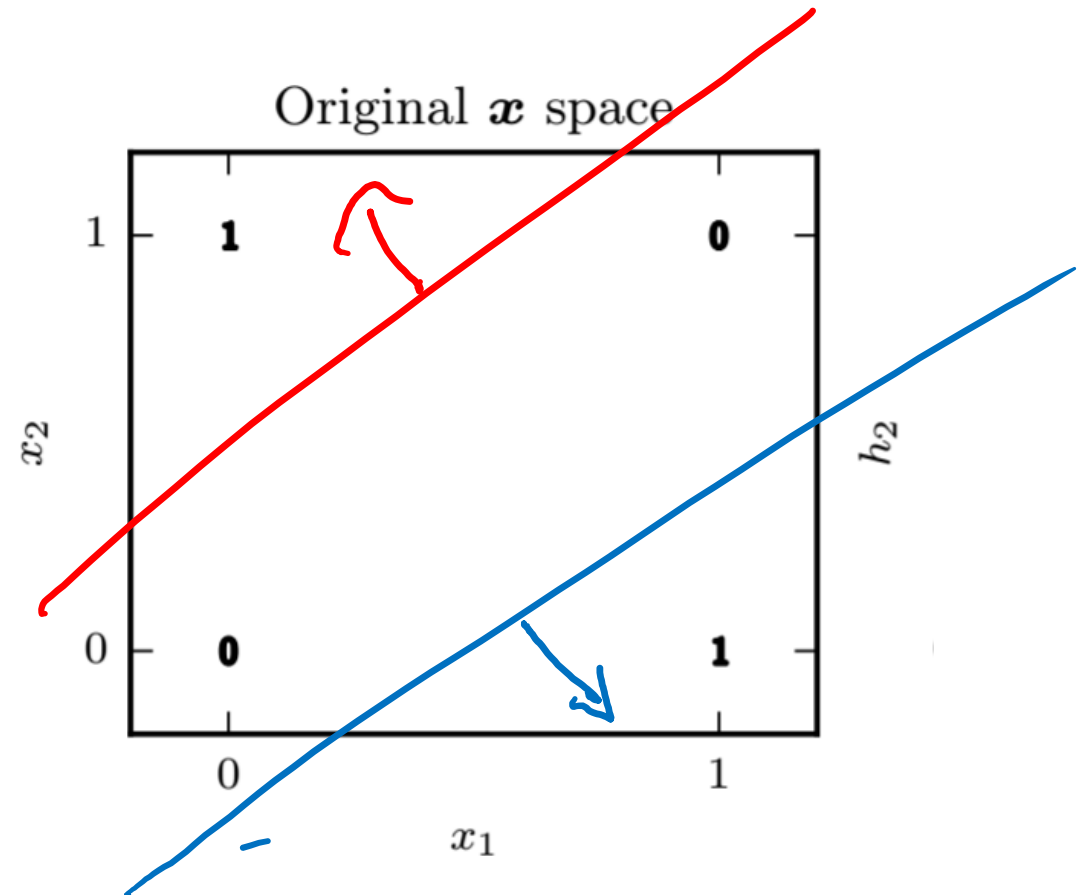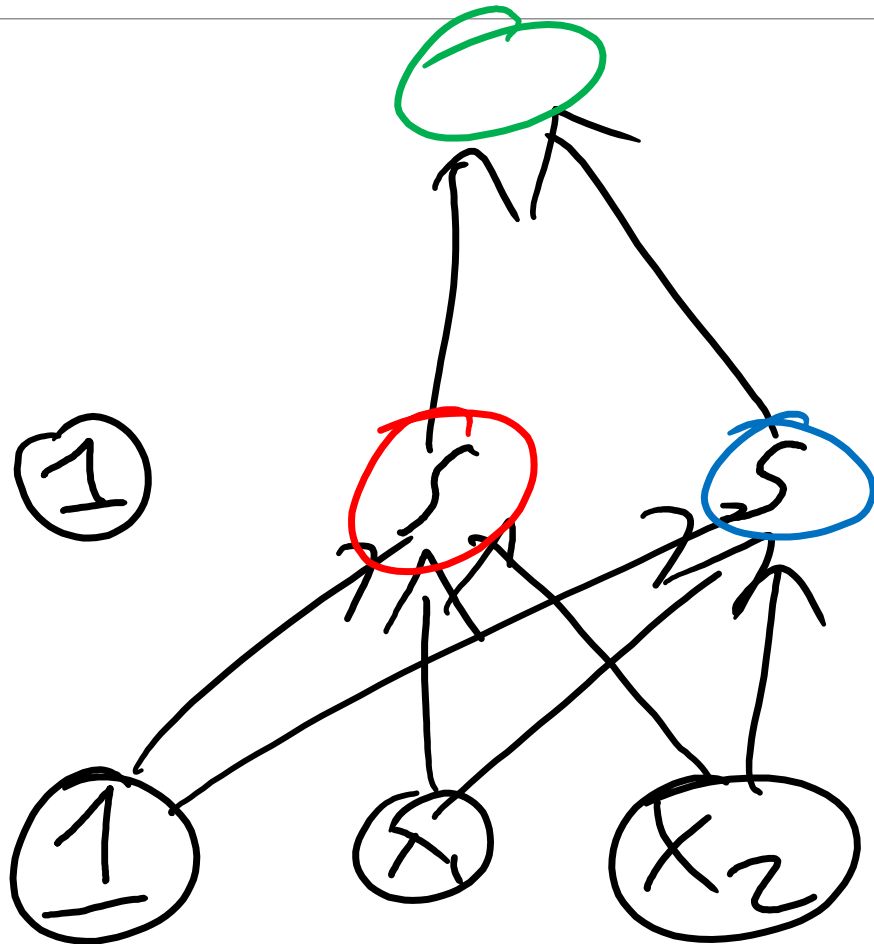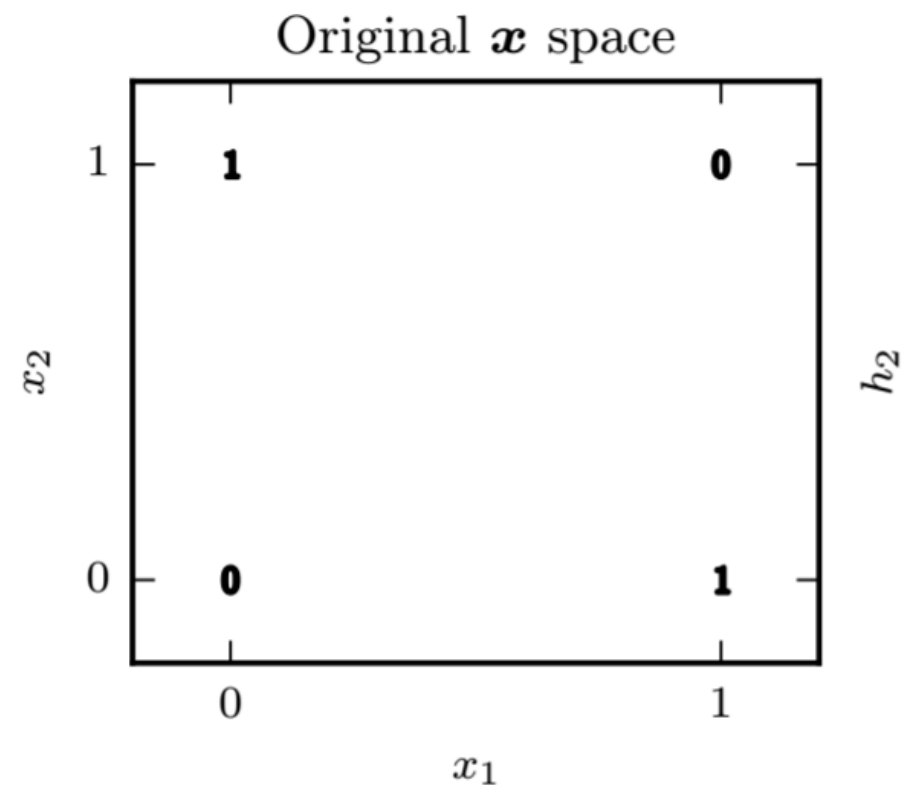
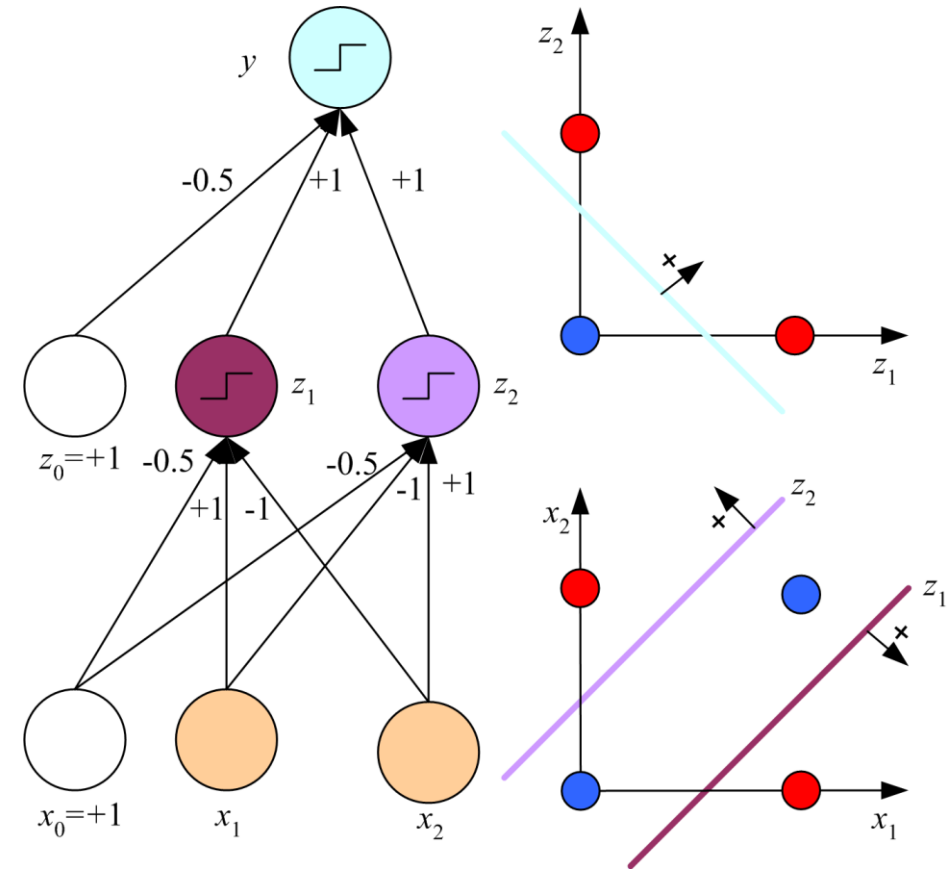# Perceptron

How do we use this to solve the XOR problem?
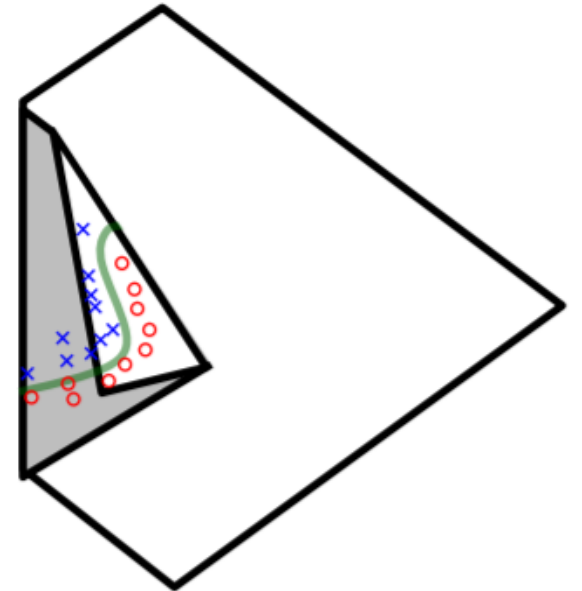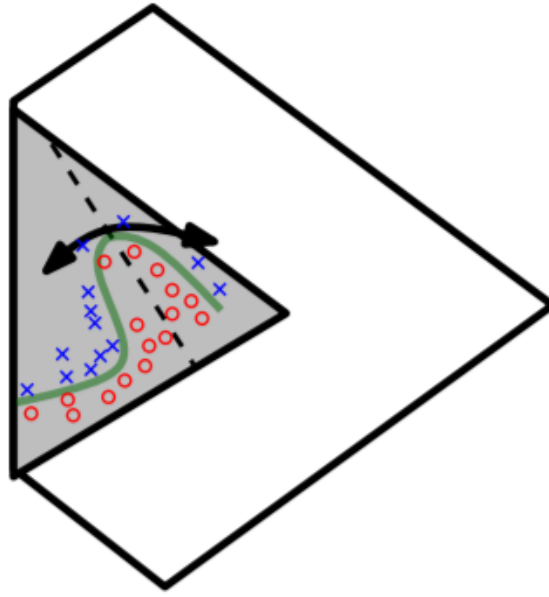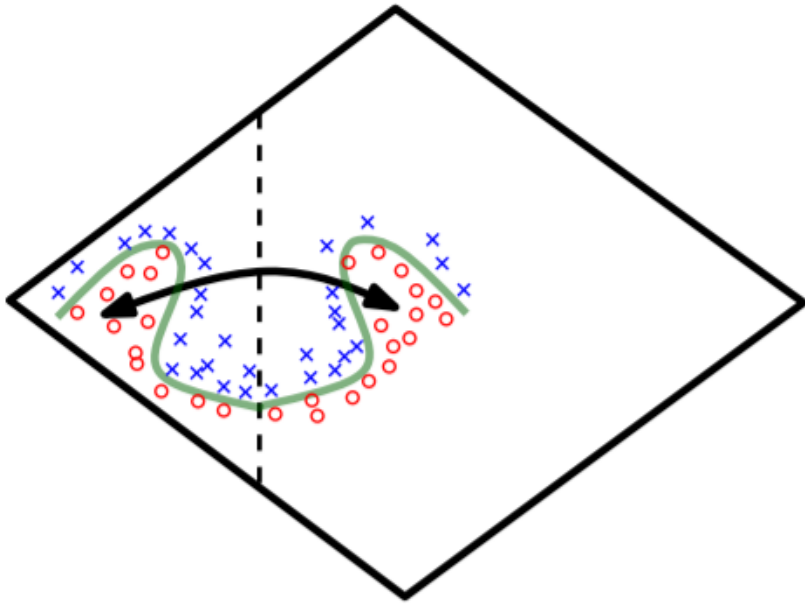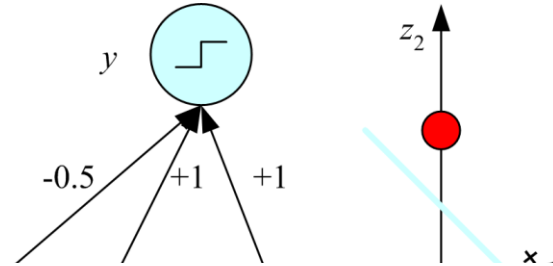
# XOR



Original $x$ space

# Perceptron

How do we use this to solve the XOR problem?

# Perceptron

How do we use this to solve the XOR problem?
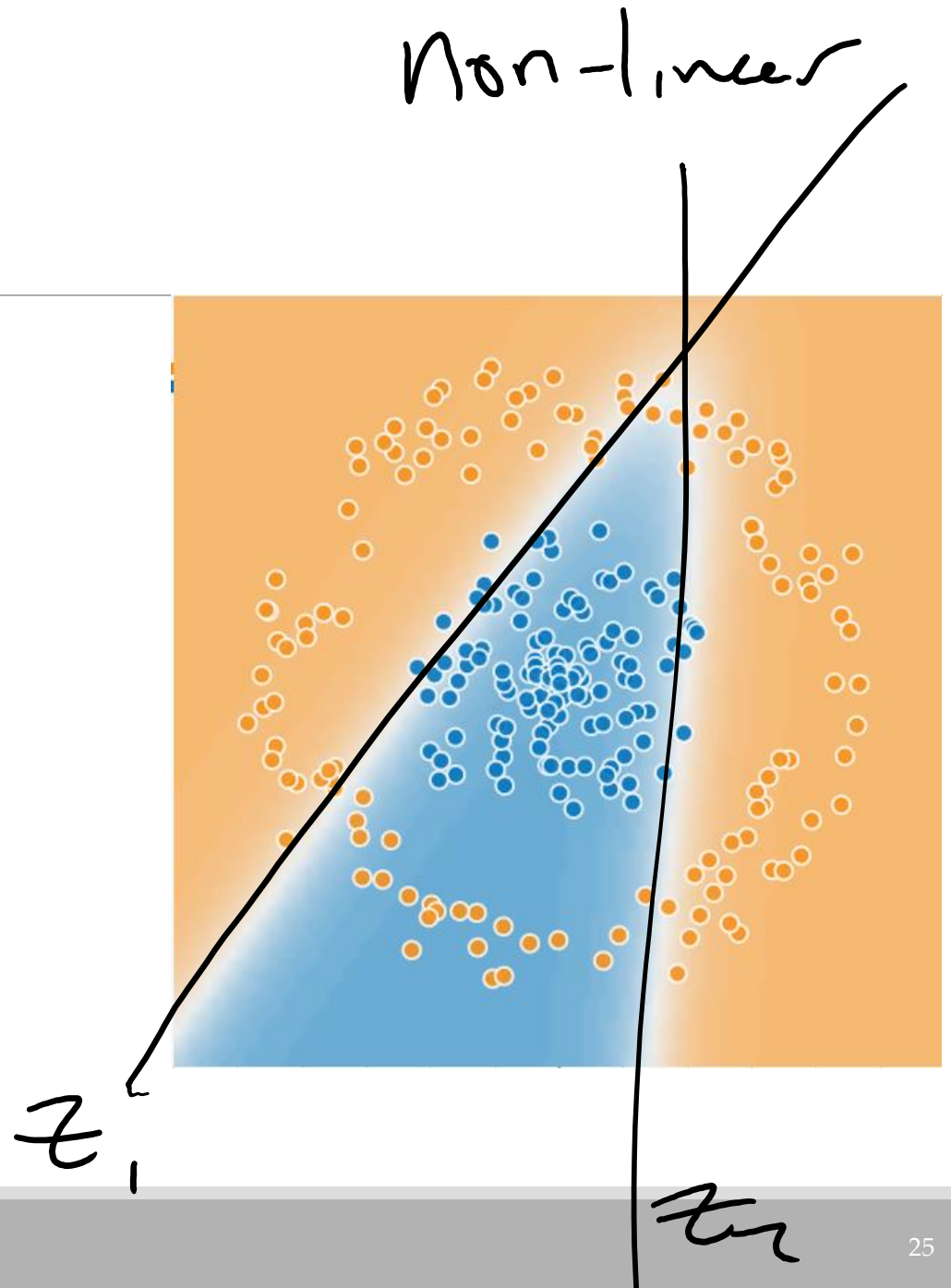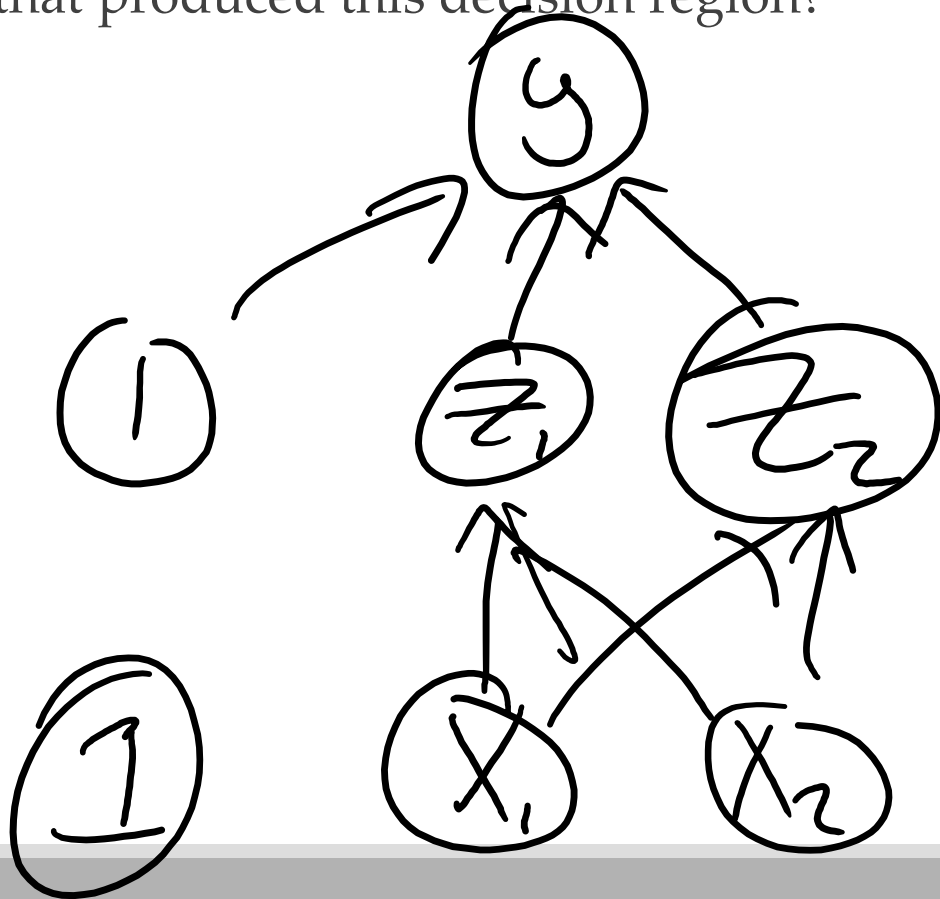◦ What is happening here?

# How to train?

We (hopefully) have some idea of how a particular set of weights causes the neural network to make a decision

- We have some vector of inputs $X_d$ that are all fed as parameters to some number of nodes
- Each of these nodes outputs a sigmoid function to the next hidden layer
- This process eventually leads to the final layer, which makes the final prediction

# How to train?

What is the structure of the neural network
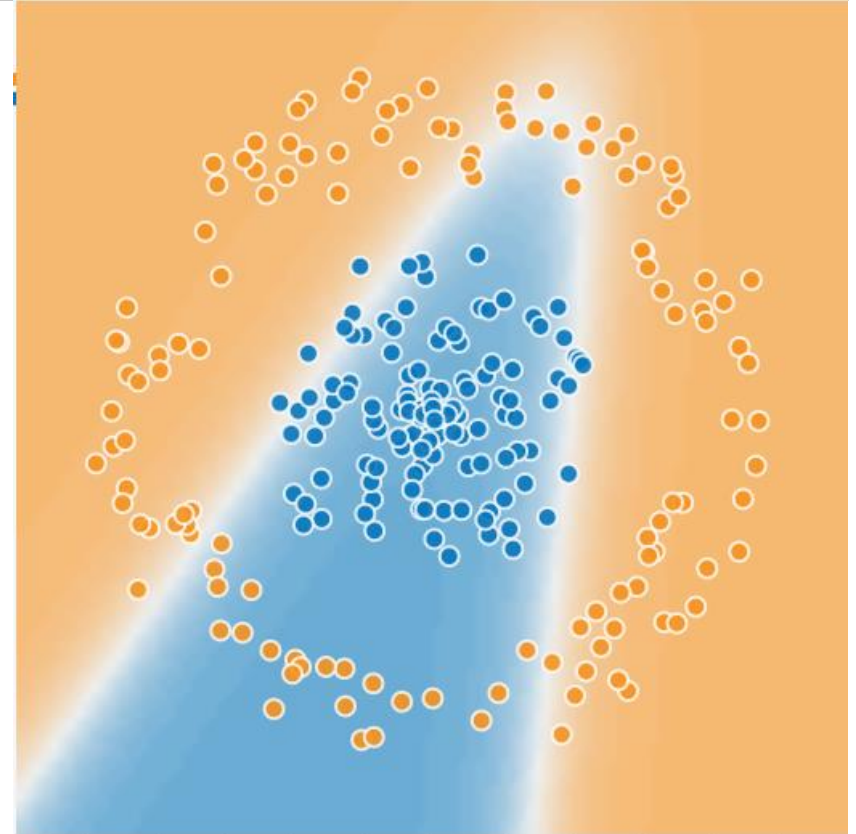that produced this decision region?

# How to train?

0.5

What is the structure of the neural network
that produced this decision region?
- Blue is positive, orange is negative
- What do the white regions represent?
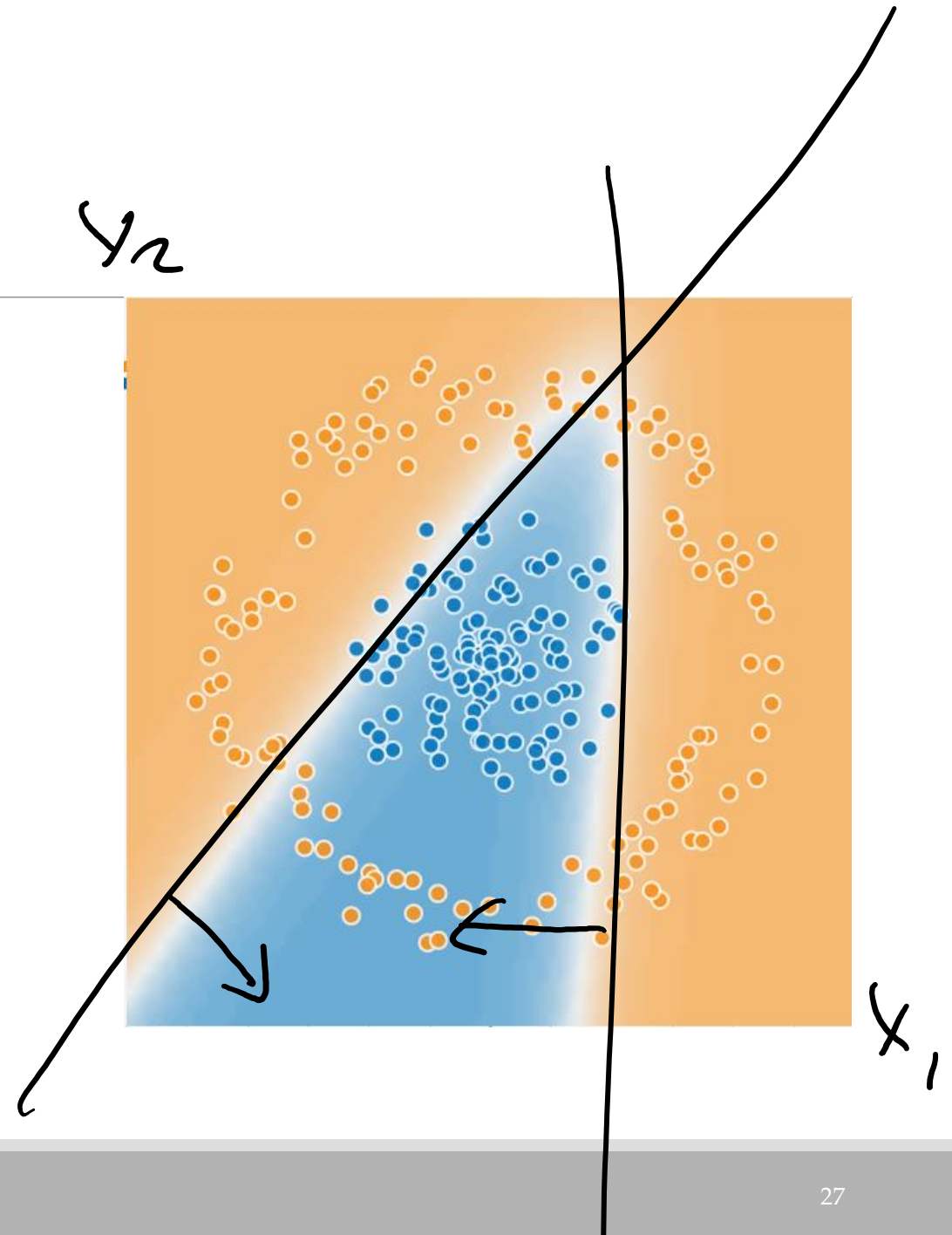
regions of low
certainty

# How to train?

What is the structure of the neural network that produced this decision region?
- Blue is positive, orange is negative
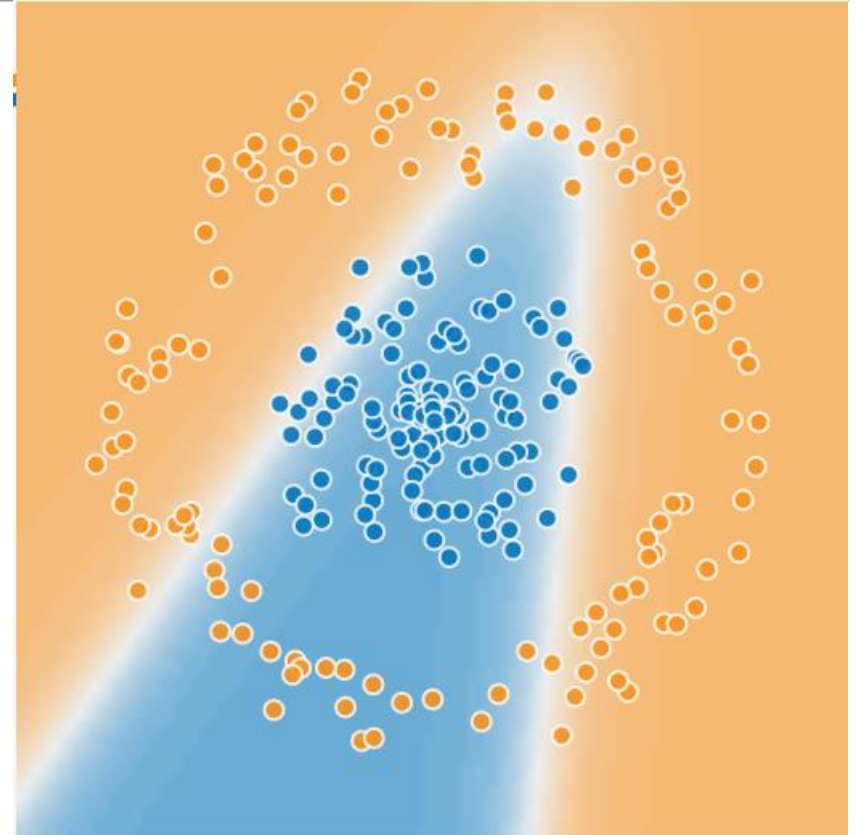- What do the white regions represent?

Two lines from two middle "hidden nodes" with sigmoid behavior

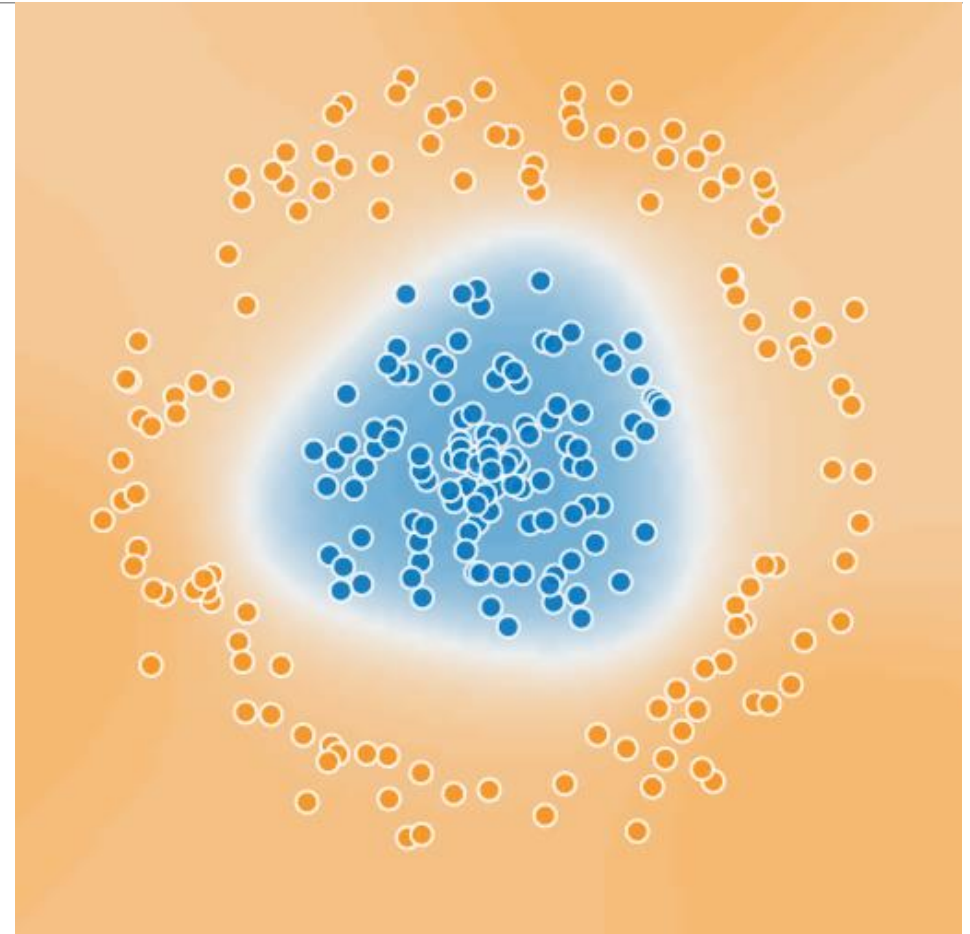What might the weights look like for each of these nodes?

# How to train?

Which would help more? Increasing the number of nodes in our hidden layer or increasing the number of hidden layers?

# How to train?

With three internal nodes, we can now generate three linear models to separate the data.
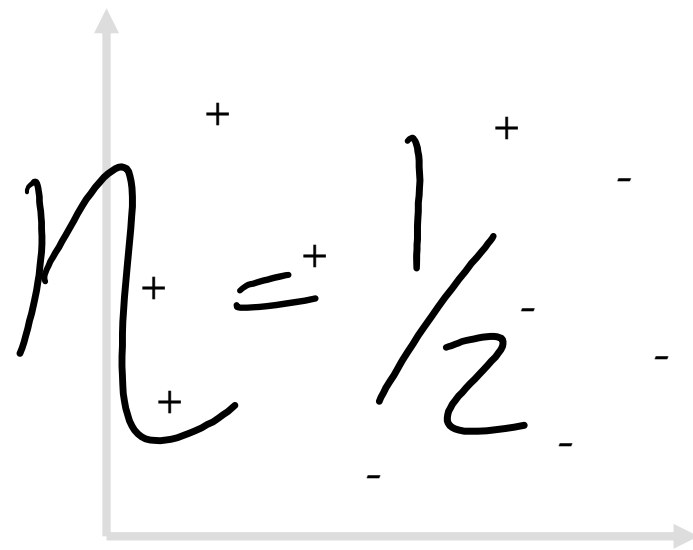
## Algorithm 8.4: Perceptron algorithm

1  Input: linearly separable data set $\mathbf{x}_i \in \mathbb{R}^D$, $y_i \in \{-1, +1\}$ for $i = 1 : N$;

2  Initialize $\boldsymbol{\theta}_0$;

3  $k \leftarrow 0$;

4  **repeat**

5     $k \leftarrow k + 1$;

6     $i \leftarrow k \bmod N$;

7     **if** $\hat{y}_i \neq y_i$ **then**

8         $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + y_i \mathbf{x}_i$

9     **else**

10         no-op

11  **until** *converged*;

$$\eta = \frac{1}{2}$$

$$\eta(\hat{y} - y)x$$

# What's wrong with this approach?

It's an easy way to separate linear data. But what's wrong?

# What's wrong with this approach?

It's an easy way to separate linear data. But what's wrong?
- Data isn't always linearly separable
- Nodes of the neural network work in conjunction with each other
  - This approach would mean that all internal nodes convert to the same point!

# What's wrong with this approach?

It's an easy way to separate linear data. But what's wrong?
- Data isn't always linearly separable
- Nodes of the neural network work in conjunction with each other
  - This approach would mean that all internal nodes convert to the same point!

How do we get around this?

# What's wrong with this approach?

It's an easy way to separate linear data. But what's wrong?
- ◦ Data isn't always linearly separable
- ◦ Nodes of the neural network work in conjunction with each other
  - ◦ This approach would mean that all internal nodes convert to the same point!

How do we get around this?
- ◦ Random weights
- ◦ Feedback between nodes