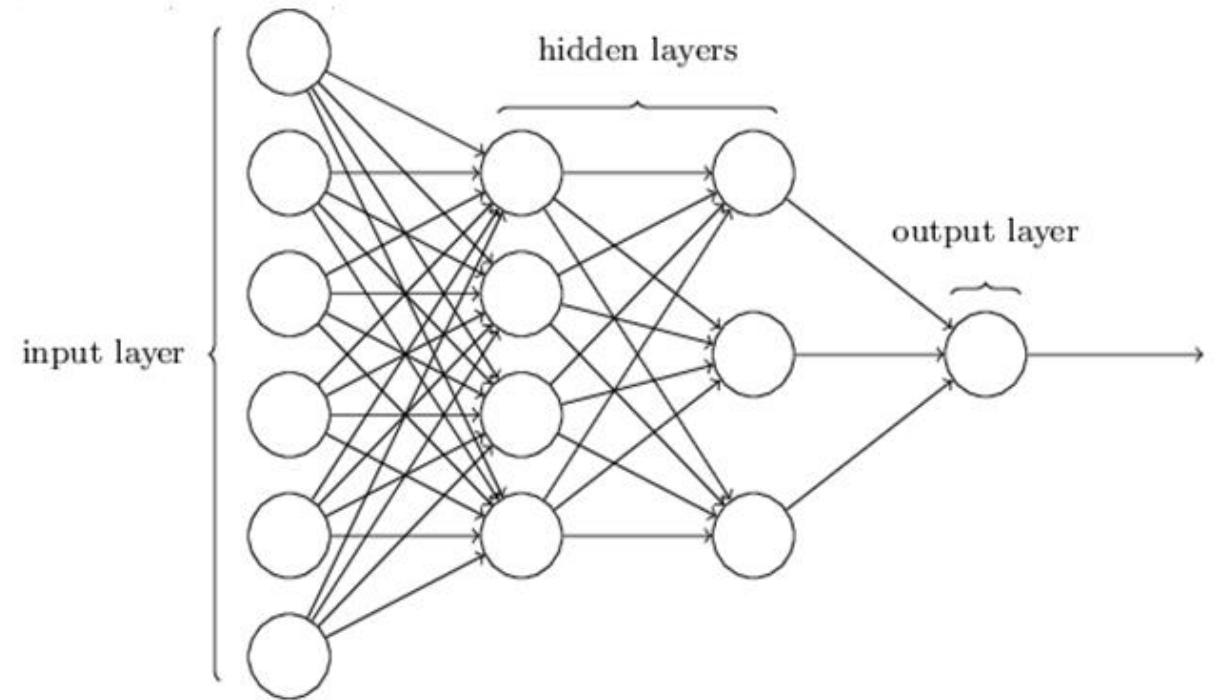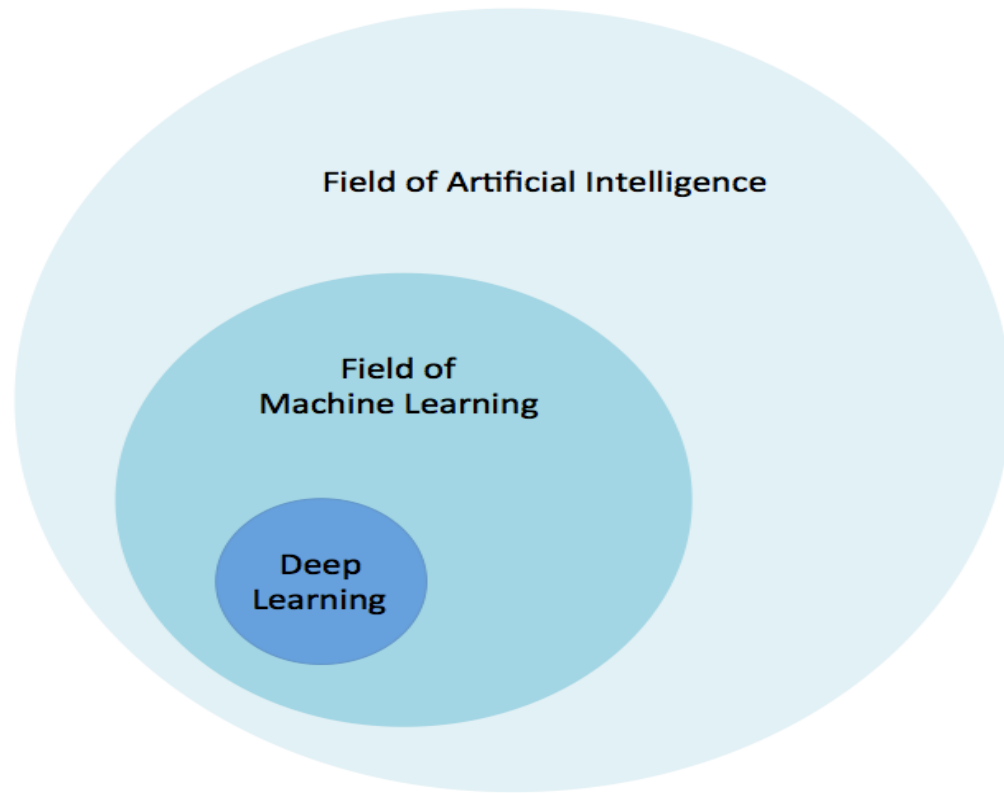# Deep Learning

Mithun Prasad, PhD
miprasad@Microsoft.com

# What Is Deep Learning?

1. Based on Algorithms that attempt to model high level abstractions in data
2. Deep learning is synonymous with artificial neural network (ANN)
3. The "deep" in deep learning refers to the depth of the network. An ANN can be very shallow

# Why Is Deep Learning Popular?

❑ DL models has been here for a long time
- Fukushima (1980) – Neo-Cognitron
- LeCun (1989) – Convolutional Neural Network

❑ DL popularity grew recently
- With growth of Big Data
- With the advent of powerful GPUs

Microsoft

# Motivation: Why Go Deep

- Deep Architectures can be representationally efficient
    Fewer computational units for same function

- Deep Representations might allow for a hierarchy or representation
    Allows non-local generalization
    Comprehensibility

- Multiple levels of latent variables allow combinatorial sharing of statistical strength

- Deep architectures work well (vision, audio, NLP, etc.) !
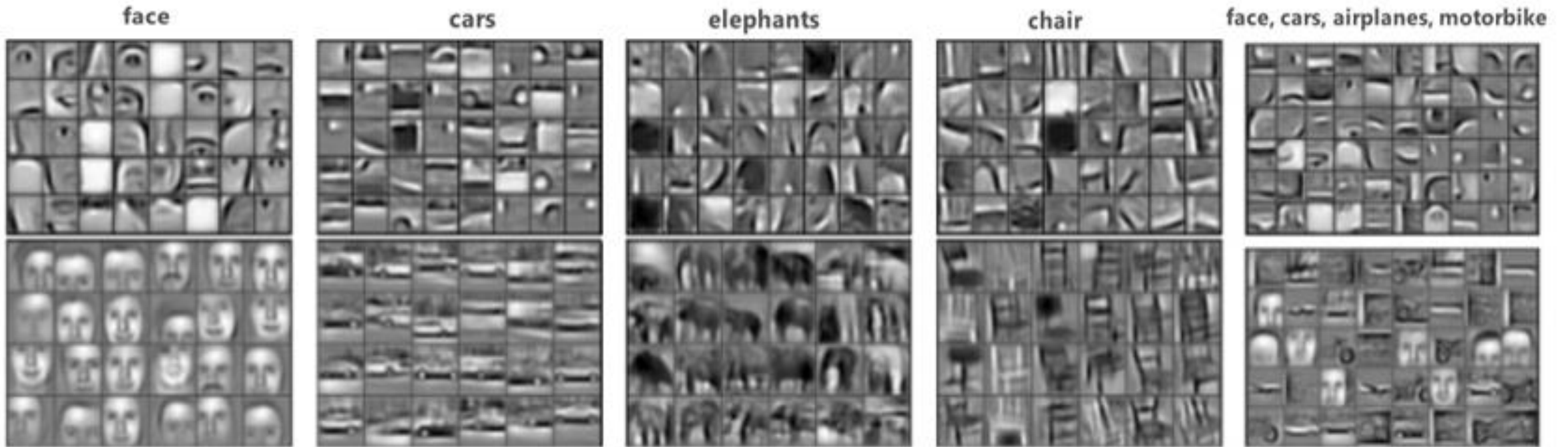
**Microsoft**

# Different Levels Of Abstraction

Hierarchical Learning

- Natural progression from low level to high level structure as seen in natural complexity

- Easier to monitor what is being learnt and to guide the machine to better subspaces

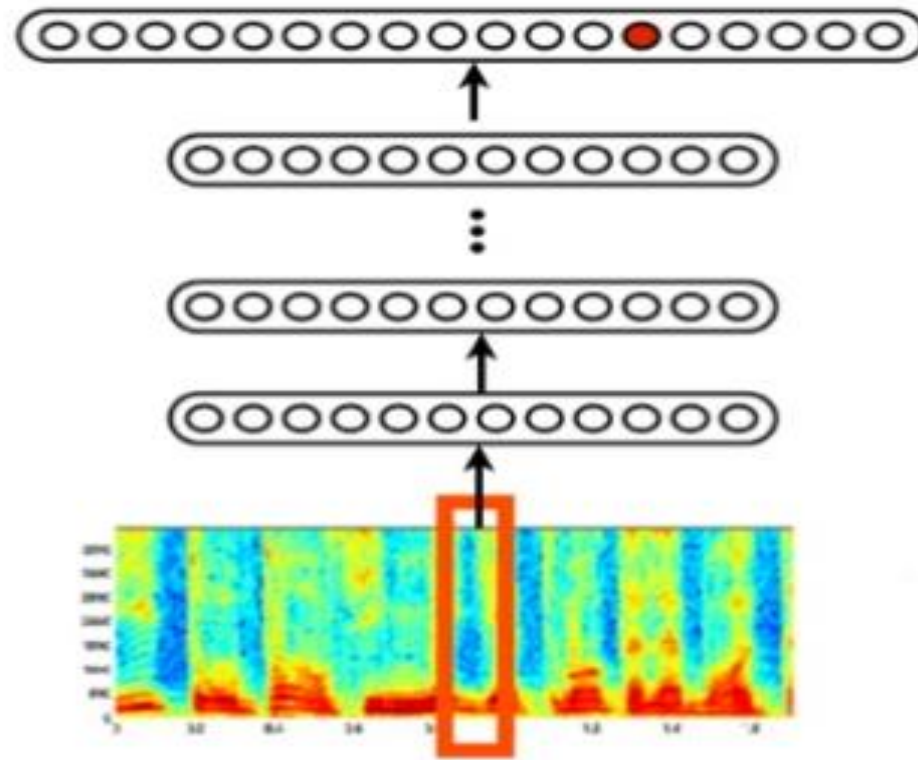- A good lower level representation can be used for many distinct tasks

Microsoft

# Compositional Data

NATURAL DATA
    IS COMPOSITIONAL.

# Compositional Data



face  cars  elephants  chair  face, cars, airplanes, motorbike
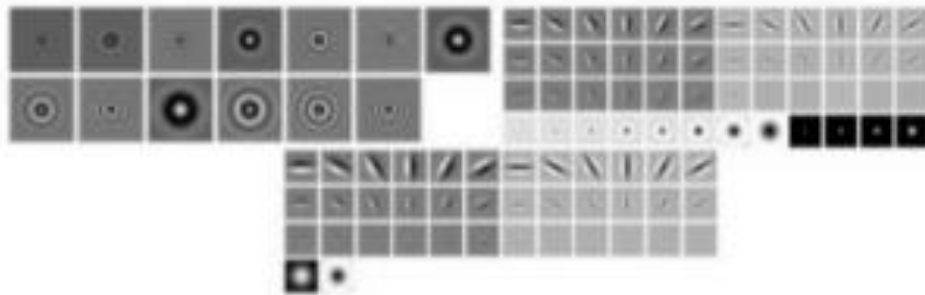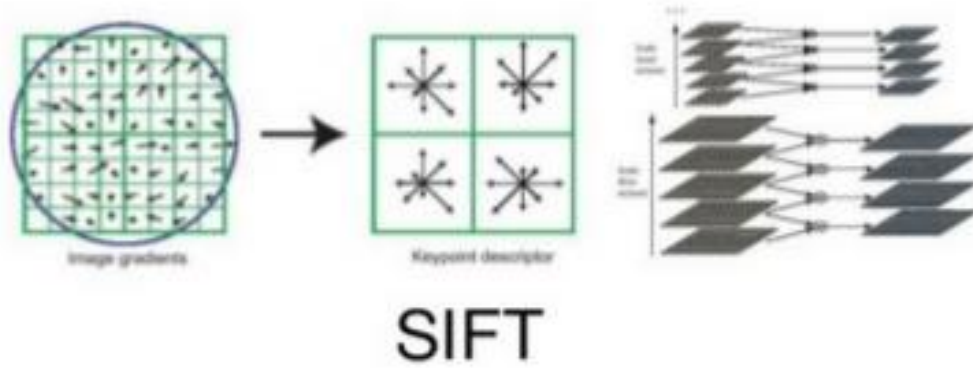
Microsoft

# Compositional Data

## Sound

# Traditional vs Deep Learning



SIFT

Textons

# Traditional vs Deep Learning

Feature extractors, required:
• Expert knowledge
• Time-consuming hand-tuning
• In industrial applications, this 90% of the time
• Sometimes are problem specific

But, what if we could learn feature extractors ?

Microsoft

# Traditional vs Deep Learning

Traditional ML requires manual feature extraction/engineering

Deep learning can automatically learn features in data

Feature extraction for unstructured data is very difficult

Deep learning is largely a "black box" technique, updating learned weights at each layer

Microsoft

# Deep Learning Begins With A Little Function

It all starts with a humble linear function called a perceptron.

$$+ \frac{\begin{array}{l} \text{weight1} \times \text{input1} \\ \text{weight2} \times \text{input2} \\ \text{weight3} \times \text{input3} \end{array}}{\text{sum}}$$

Perceptron:
If sum > threshold: output 1
Else: output 0

Example: The inputs can be your data. Question: Should I buy this car?

$$+ \frac{\begin{array}{l} 0.2 \times \text{gas mileage} \\ 0.3 \times \text{horsepower} \\ 0.5 \times \text{num cup holders} \end{array}}{\text{sum}}$$

Perceptron:
If sum < threshold: buy
Else: walk

Microsoft
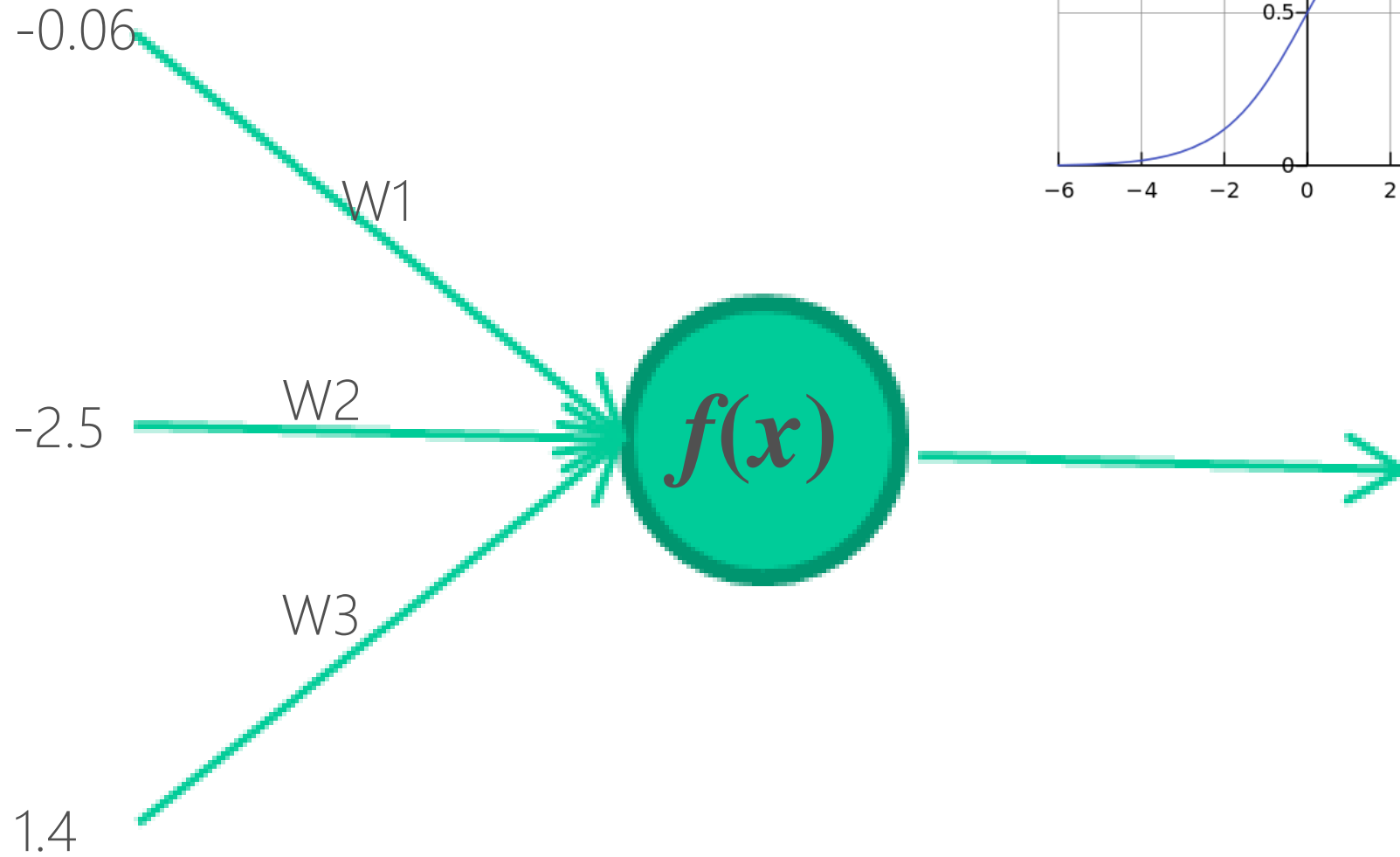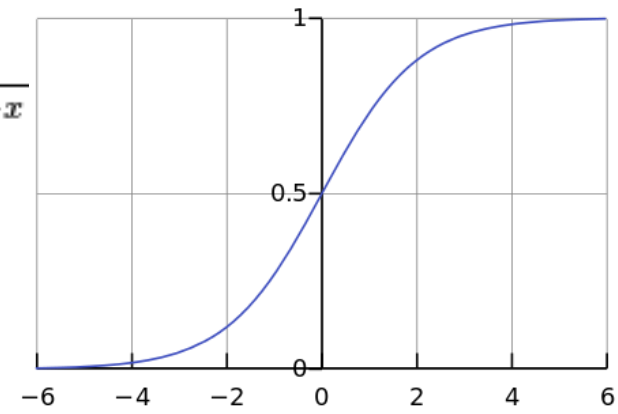
# These Little Functions Are Chained Together

- Deep learning comes from chaining a bunch of these little functions together
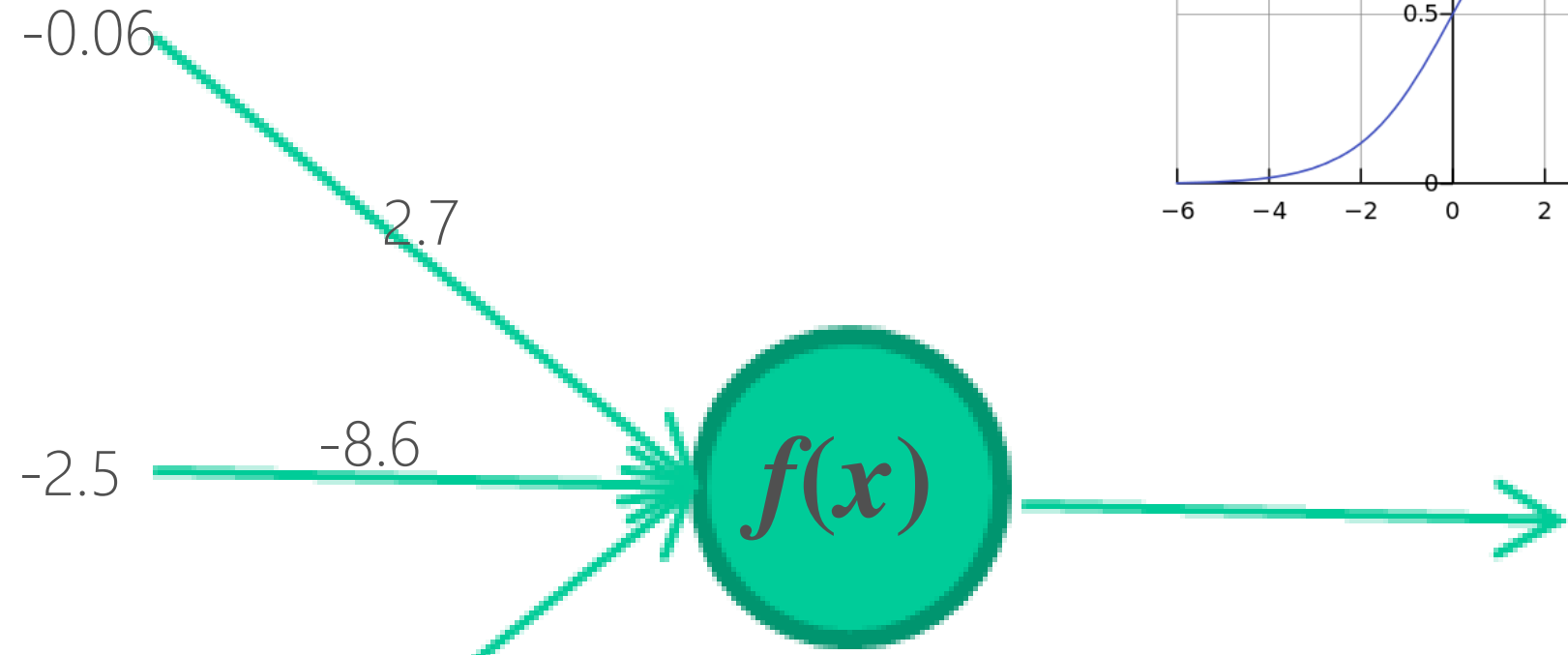- Chained together, they are called neurons

weight1 ✖ input1

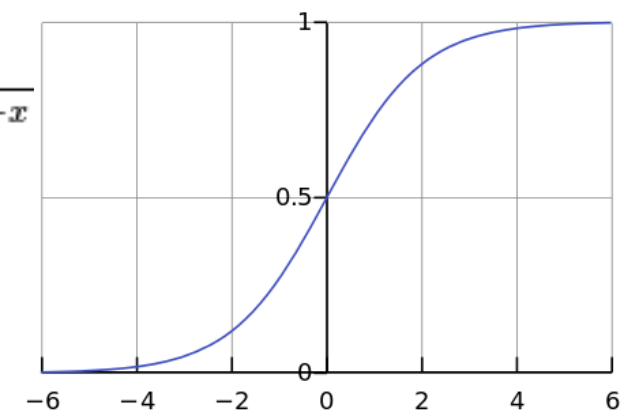weight2 ✖ input2

weight3 ✖ input3

Microsoft

# Deep Neural Network (DNN)



"Non-deep" feedforward neural network

input layer

hidden layer

output layer

Deep neural network

input layer

hidden layer 1    hidden layer 2    hidden layer 3

output layer

Microsoft

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

-0.06

W1

-2.5

W2

$f(x)$

W3

1.4

Microsoft

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

-0.06

2.7

-2.5

-8.6

0.002

1.4

$f(x)$

$x = $ -0.06×2.7 + 2.5×8.6 + 1.4×0.002  = 21.34

Microsoft

A dataset

Fields                 class
1.4   2.7   1.9          0
3.8   3.4   3.2          0
6.4   2.8   1.7          1
4.1   0.1   0.2          0
etc ...

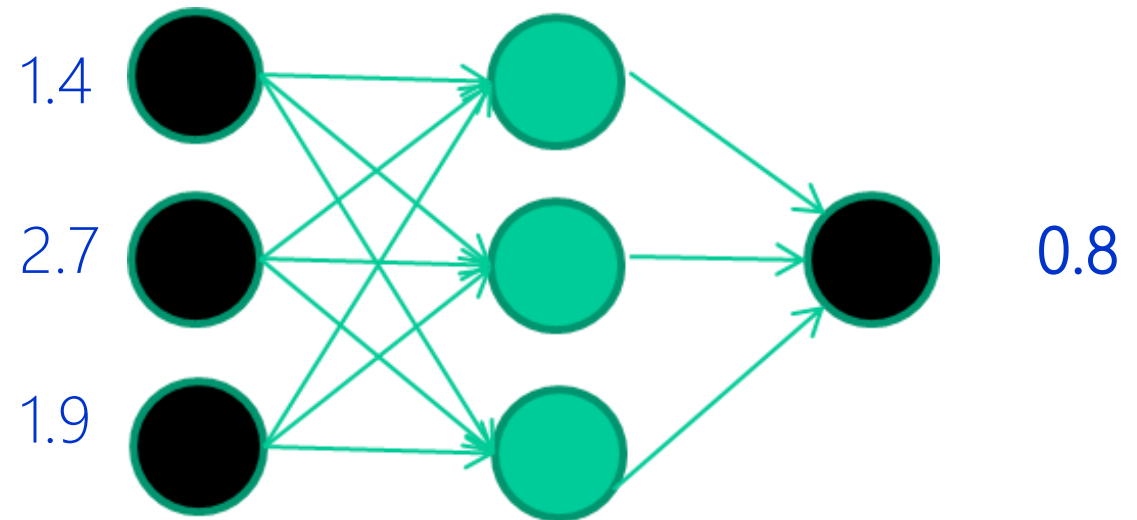*Training the neural network*

Fields             class

1.4   2.7   1.9        0

3.8   3.4   3.2        0

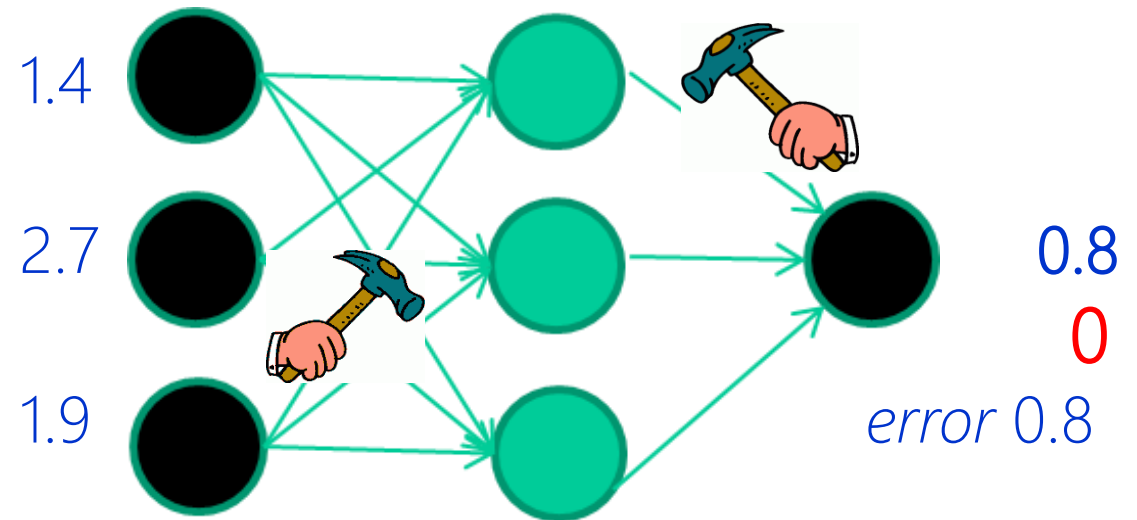6.4   2.8   1.7        1

4.1   0.1   0.2        0

etc ...

*Training the neural network*

*Fields*          *class*
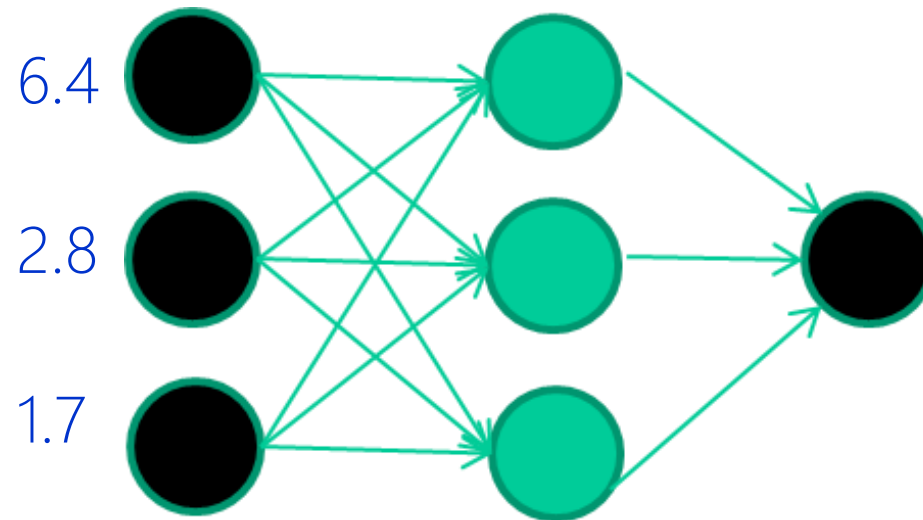
1.4  2.7  1.9       0

3.8  3.4  3.2       0

6.4  2.8  1.7       1

4.1  0.1  0.2       0

etc ...

Initialise with random weights

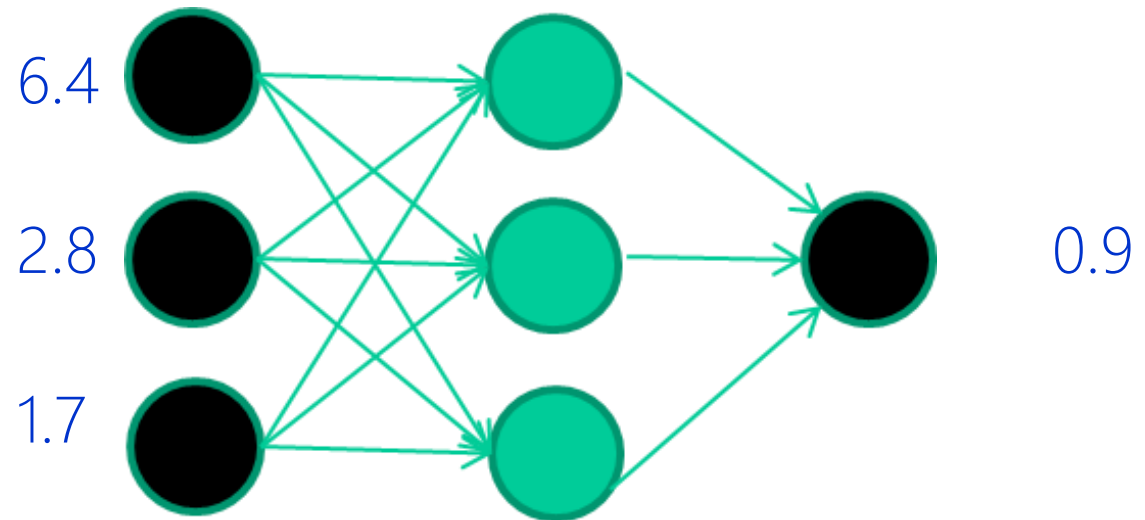Fields          class
1.4  2.7  1.9     0
3.8  3.4  3.2     0
6.4  2.8  1.7     1
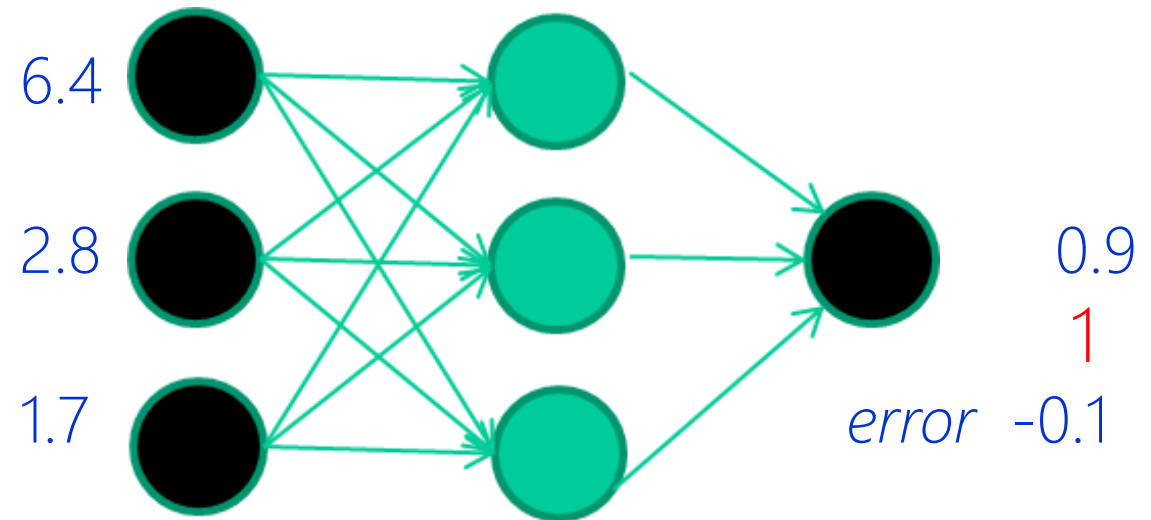4.1  0.1  0.2     0
etc ...

Present a training pattern

1.4
2.7
1.9

Fields          class
1.4  2.7  1.9      0
3.8  3.4  3.2      0
6.4  2.8  1.7      1
4.1  0.1  0.2      0
etc ...

Microsoft

Fields      class

| | | | |
|---|---|---|---|
| 1.4 | 2.7 | 1.9 | 0 |
| 3.8 | 3.4 | 3.2 | 0 |
| 6.4 | 2.8 | 1.7 | 1 |
| 4.1 | 0.1 | 0.2 | 0 |

etc ...

**Feed it through to get output**



6.4

2.8

1.7

0.9
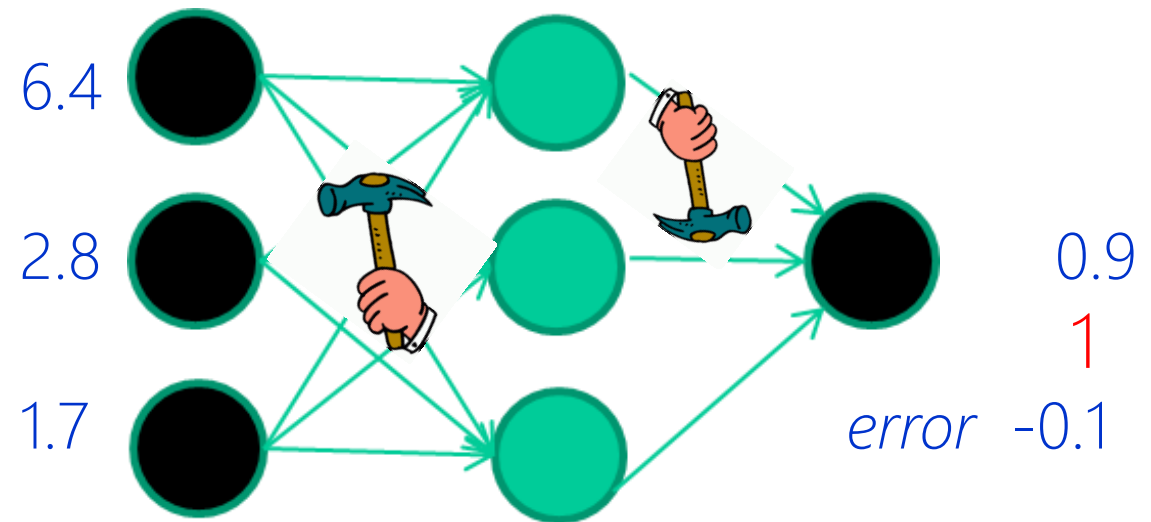
Fields                 class
1.4   2.7   1.9          0
3.8   3.4   3.2          0
6.4   2.8   1.7          1
4.1   0.1   0.2          0
etc ...

And so on ....

6.4

2.8                                         0.9
                                             1
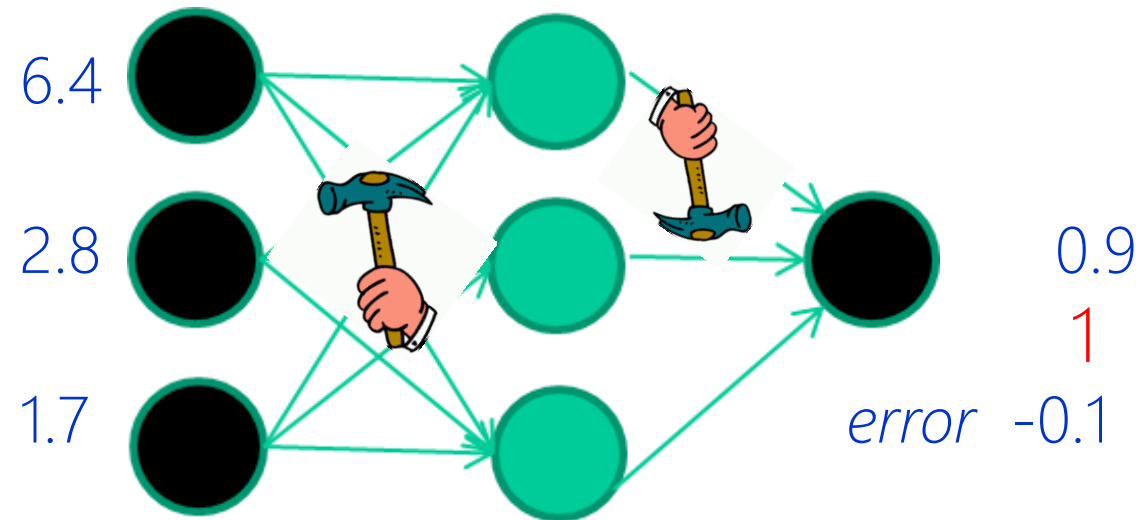1.7                                  error  -0.1

Repeat this thousands, maybe millions of times – each time
taking a random training instance, and making slight
weight adjustments
  Algorithms for weight adjustment are designed to make
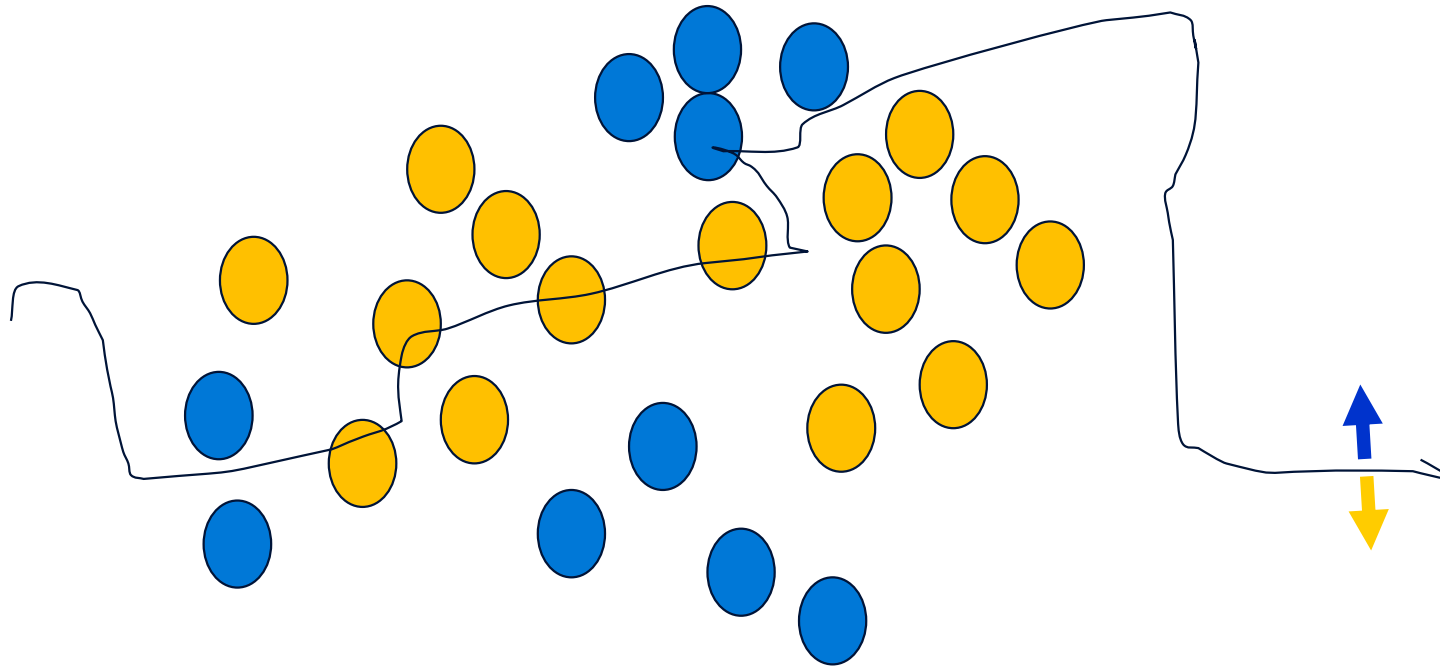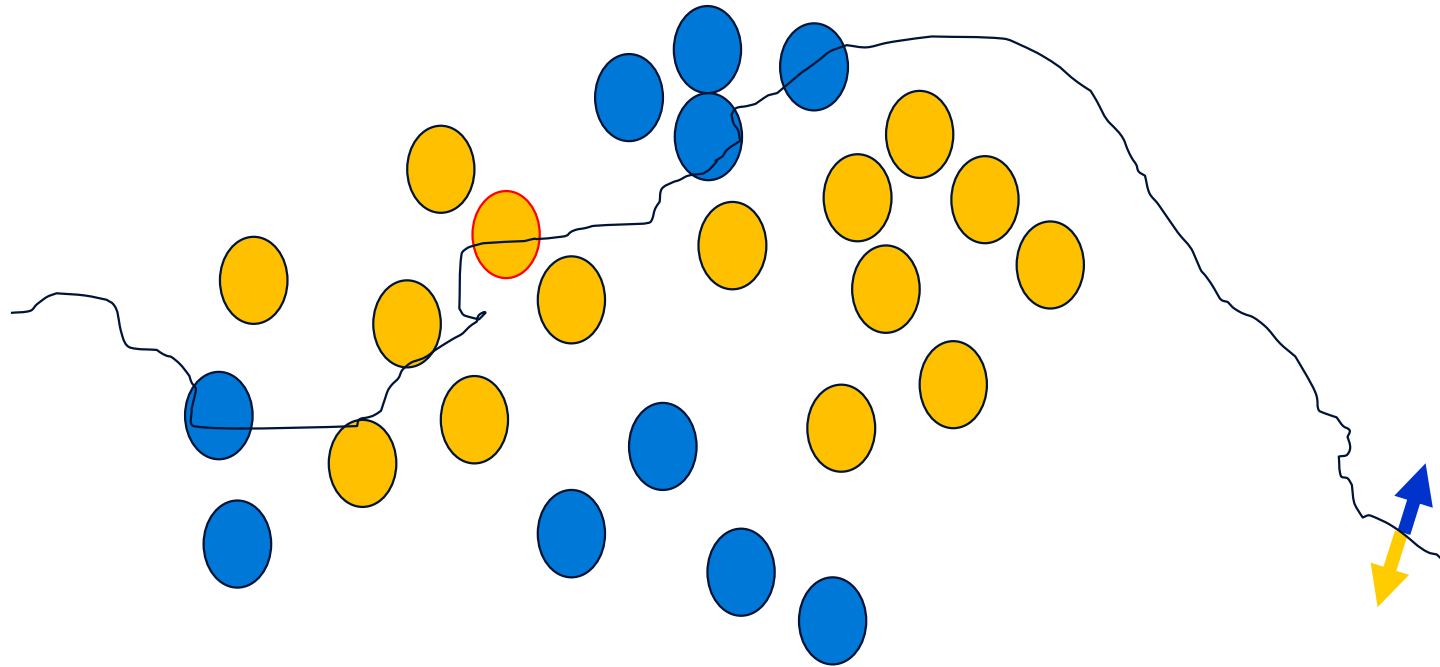changes that will reduce the error

Microsoft

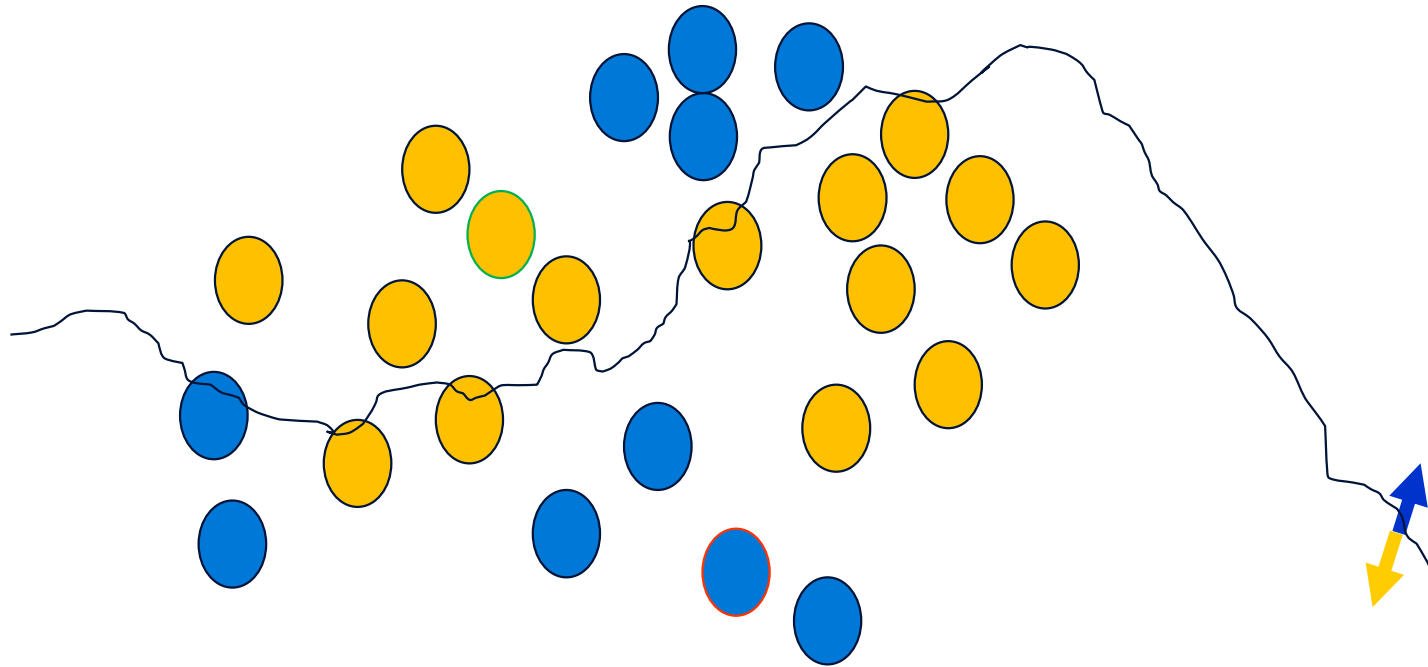# The Decision Boundary Perspective...
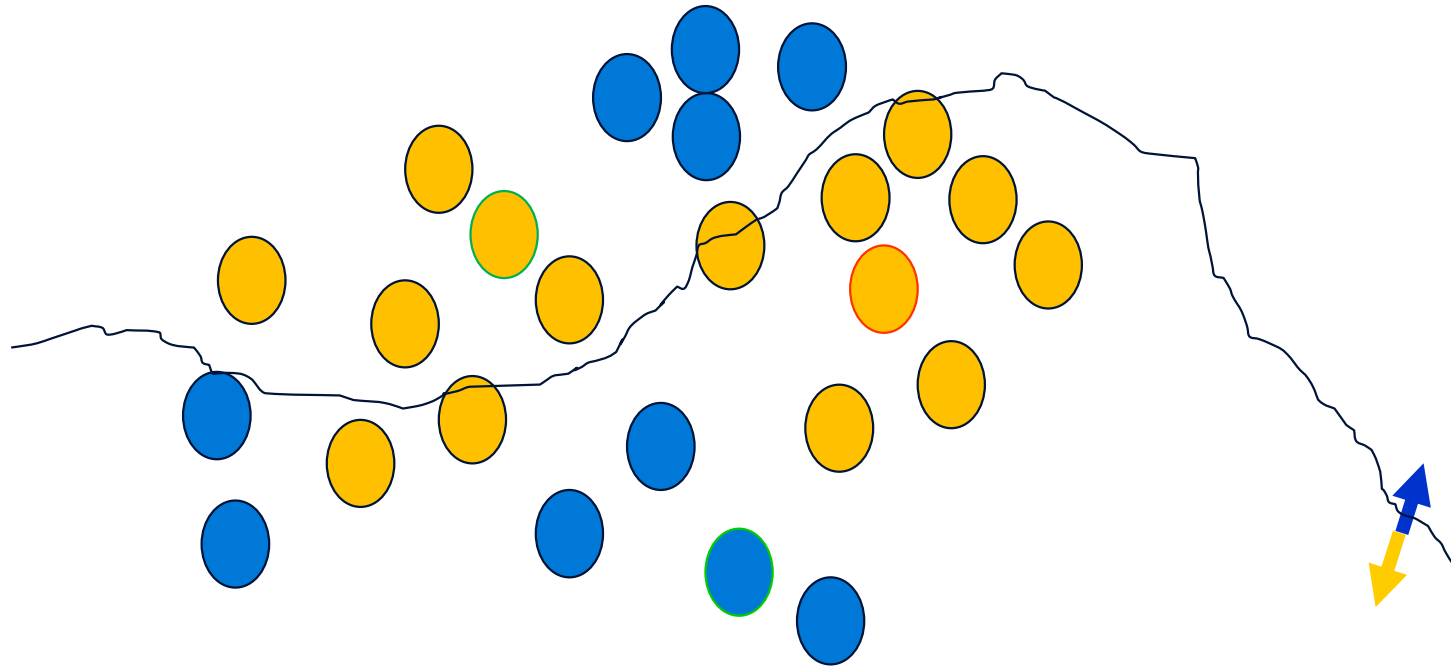
Initial random weights

# The Decision Boundary Perspective...

Present a training instance / adjust the weights

# The Decision Boundary Perspective...



Present a training instance / adjust the weights
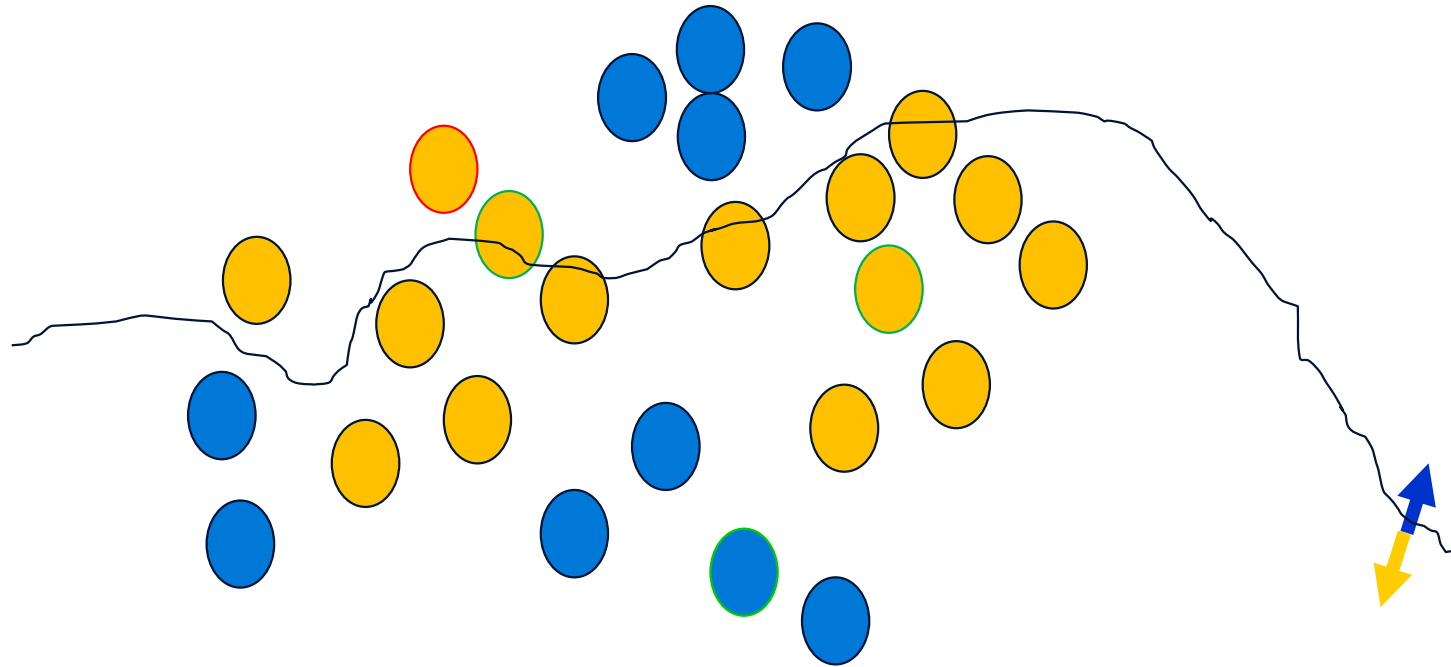
# The Decision Boundary Perspective...

Present a training instance / adjust the weights

# The Decision Boundary Perspective...

Present a training instance / adjust the weights
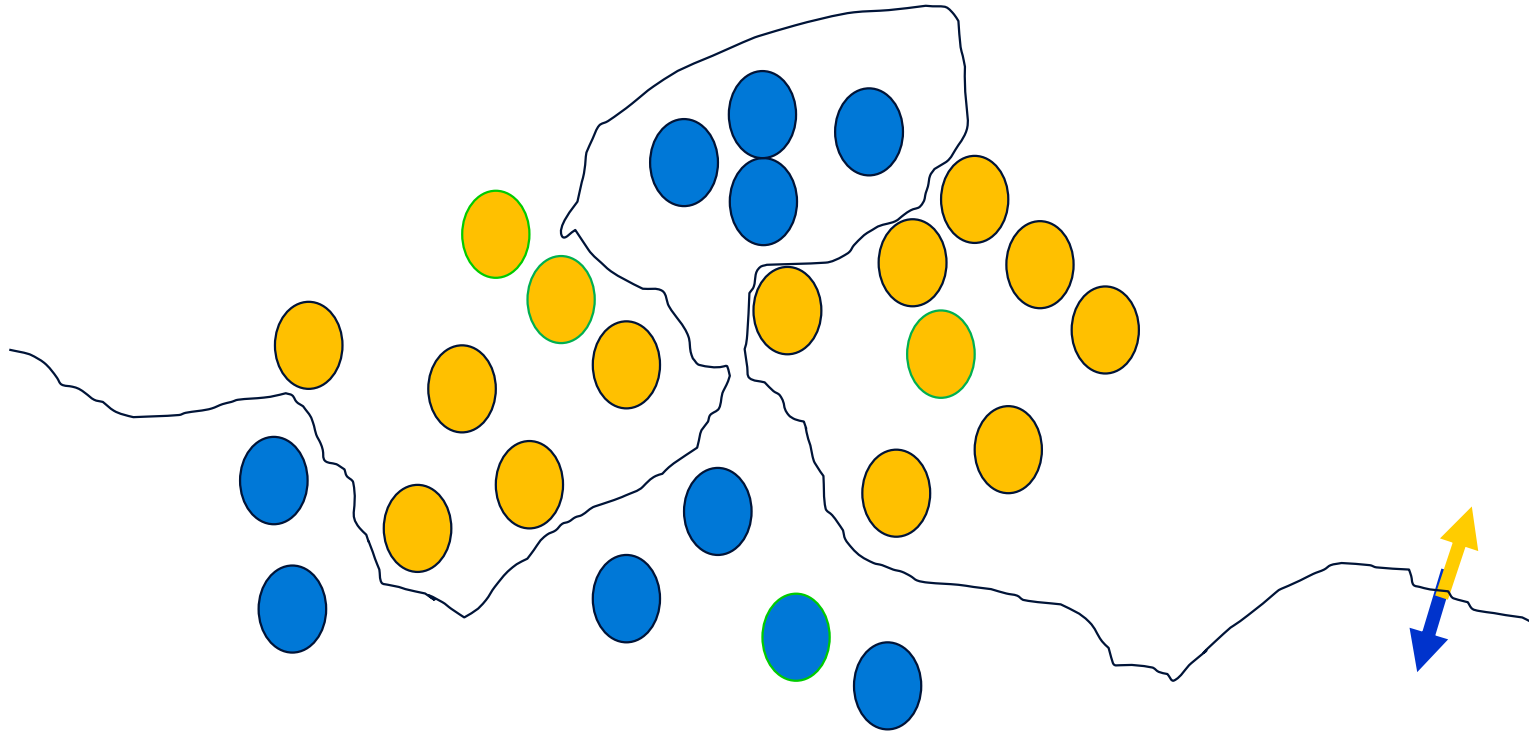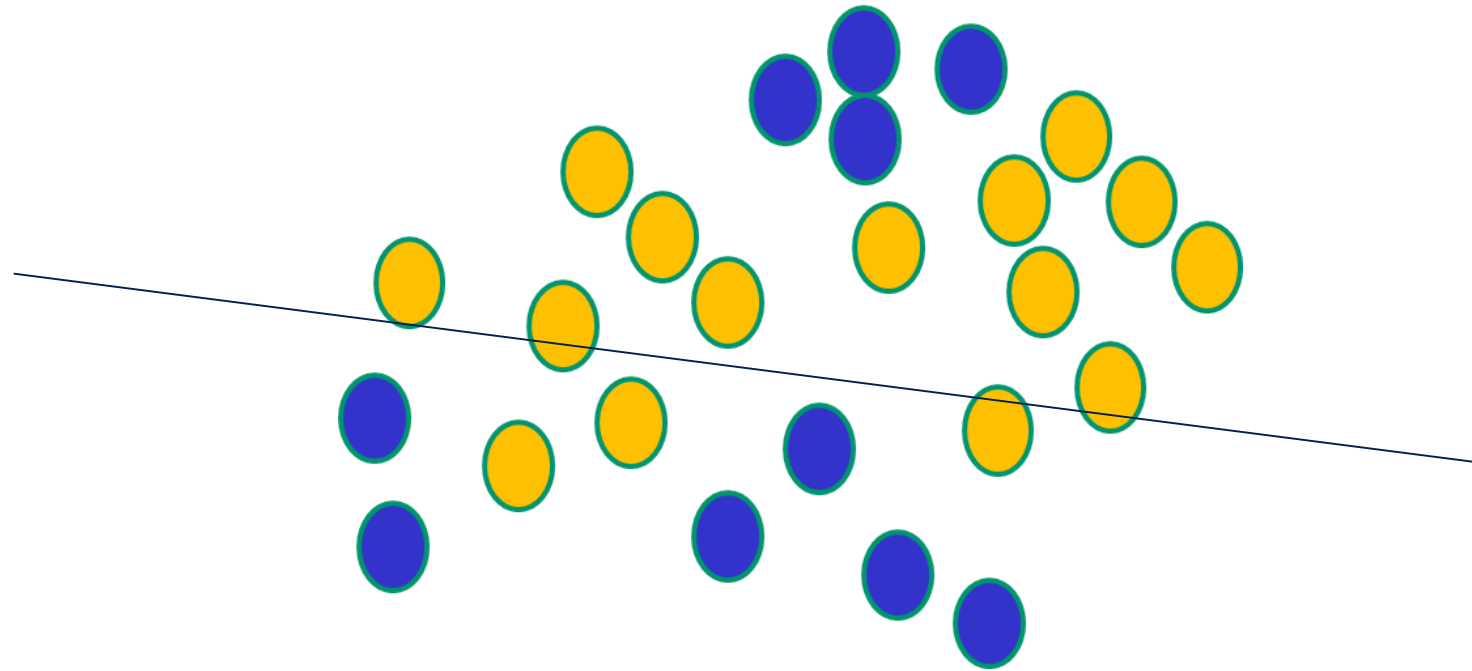
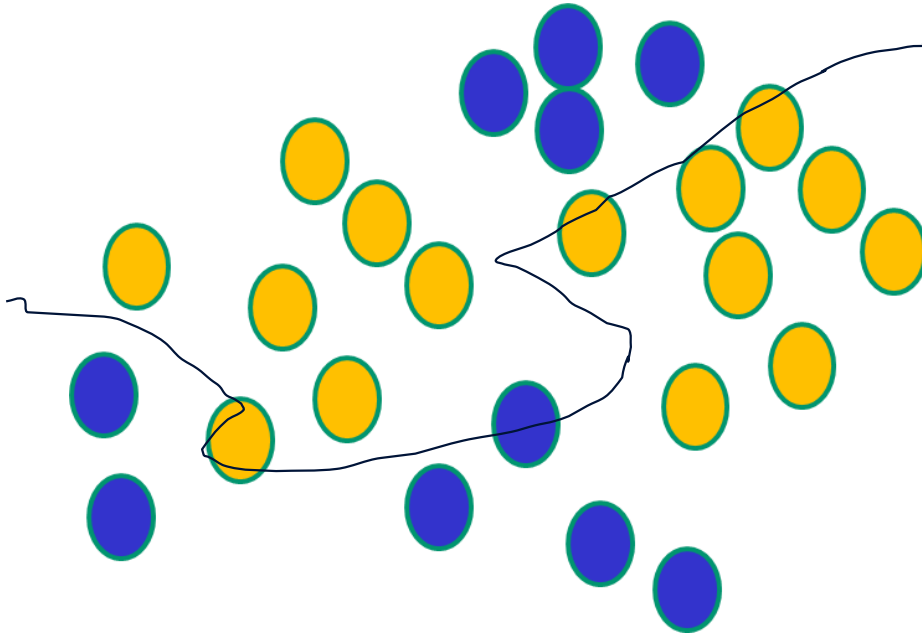# The Decision Boundary Perspective...

Eventually ....

# Some Other 'By The Way' Points

If $f(x)$ is linear, the NN can **only** draw straight decision boundaries (even if there are many layers of units)
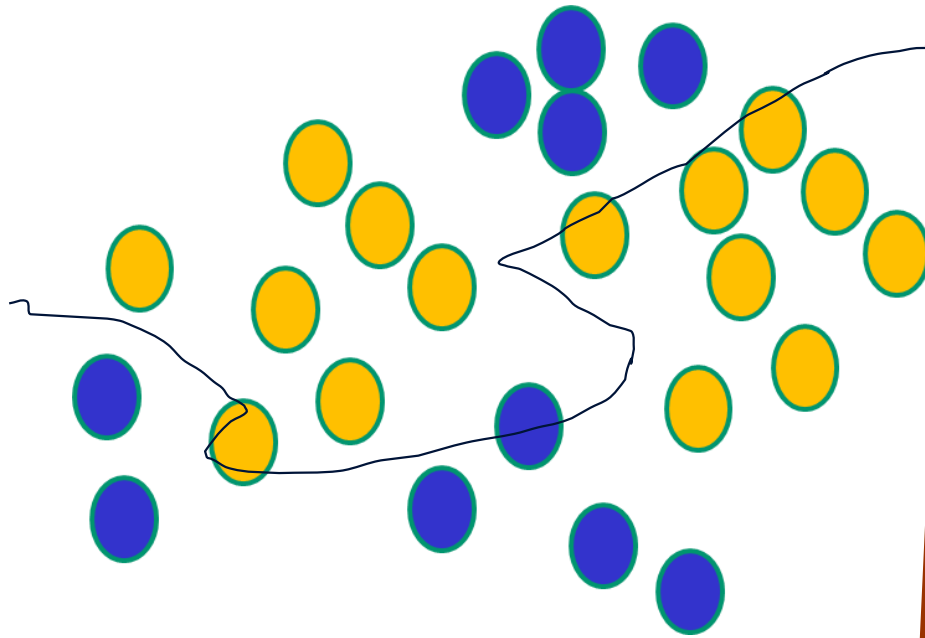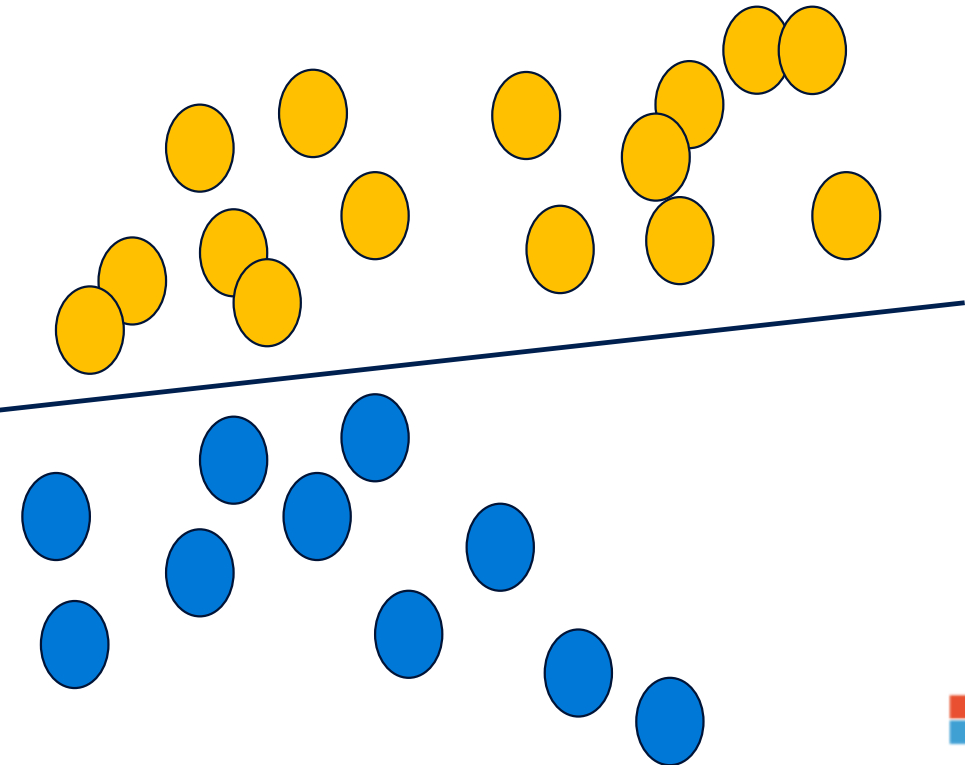


Microsoft

# Some Other 'By The Way' Points

NNs use nonlinear $f(x)$ so they can draw complex boundaries, but keep the data unchanged

# Some Other 'By The Way' Points
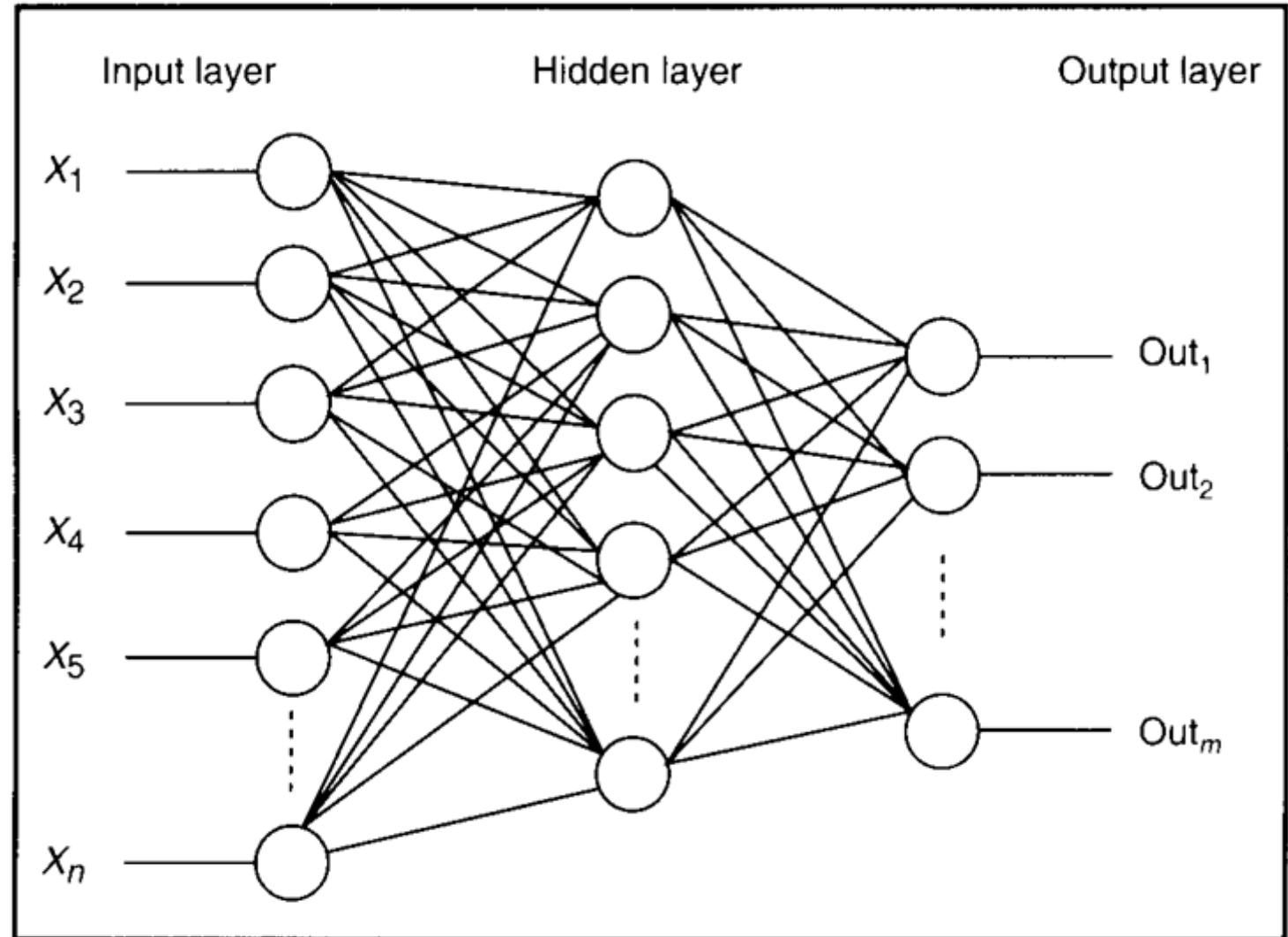
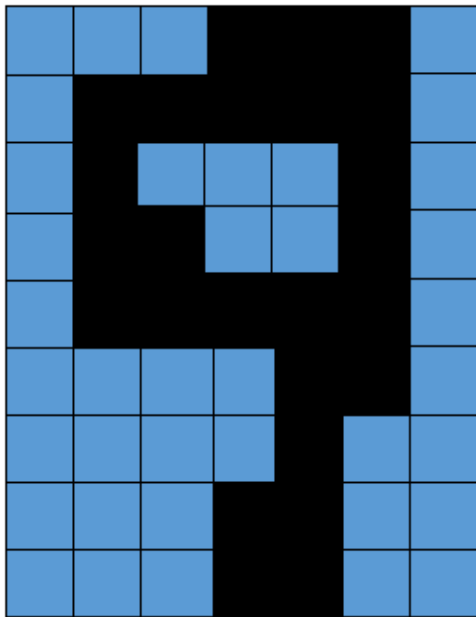NNs use nonlinear $f(x)$ so they can draw complex boundaries, but keep the data unchanged

SVMs only draw straight lines, but they transform the data first in a way that makes that OK



Microsoft

# Feature Detectors

# What Is This Unit Doing?

# Hidden Layer Units Become Self-Organised Feature Detectors

# What Does This Unit Detect?



strong +ve weight

low/zero weight

Strong signal for a horizontal line in the top row, ignoring everywhere else

What features might you expect a good NN to learn, when trained with data like this?

Microsoft

vertical lines

Microsoft

Horizontal lines

Small circles

Small circles

But what about position invariance ???
our example unit detectors were tied to
specific parts of the image

Microsoft

# Successive Layers Can Learn Higher-Level Features



Detect lines in specific positions

Higher level detectors ( horizontal line, "RHS vertical lune" "upper loop", etc…

etc …

etc …

# Multiple Layers Make Sense

# Multiple Layers Make Sense

- Deep Learning = Brain "inspired"
- Audio / Visual Cortex has multiple stages = Hierarchical



Object Recognition

Driver Assistance (ADAS)

Artificial Intelligence (AI)

Video Analytics

Augmented Reality / Virtual Reality

Microsoft

# Multiple Layers Make Sense

Many-layer neural network architectures should be capable of learning the true underlying features and 'feature logic', and  therefore generalise very well …

# The New Way To Train Multi-Layer NNs...

# The New Way To Train Multi-Layer NNs...



Train this layer first

# The New Way To Train Multi-Layer NNs...



Train this layer first

then this layer

Microsoft

# The New Way To Train Multi-Layer NNs...



Train this layer first

then this layer

then this layer

# The New Way To Train Multi-Layer NNs...



Train this layer first

then this layer

then this layer

then this layer

Microsoft

# The New Way To Train Multi-Layer NNs...



Train this layer first

then this layer

then this layer

then this layer

finally this layer

Microsoft

# The New Way To Train Multi-Layer NNs...



Each layer can be thought of as a set of features

# Idea Behind Deep Learning

- There are many types of deep learning
- Different kinds of autoencoder, variations on architectures and training algorithms, etc.
- It's a growing area

# Common DNNs

❑ Deep Convolutional Neural Network (DCNN)
  ▪ To extract representation from images

❑ Recurrent Neural Network (RNN)
  ▪ To extract representation from sequential data

❑ Deep Belief Neural Network (DBN)
  ▪ To extract hierarchical representation from a dataset

❑ Deep Reinforcement Learning (DQN)
  ▪ To prescribe how agents should act in an environment in order to maximize future cumulative reward (e.g., a game score)

We will cover DCNN today

Microsoft

# Open Source Deep Learning Frameworks

DL4J
- JVM-based
- Distrubted
- Integrates with Hadoop and Spark

Theano
- Very popular in Academia
- Fairly low level
- Interfaced with via Python and Numpy

Torch
- Lua based
- In house versions used by Facebook and Twitter
- Contains pretrained models

Microsoft

# Open Source Deep Learning Frameworks

TensorFlow
- Google written successor to Theano
- Interfaced with via Python and Numpy
- Highly parallel
- Can be somewhat slow for certain problem sets

Caffe
- Not general purpose. Focuses on machine-vision problems
- Implemented in C++ and is very fast
- Not easily extensible
- Has a Python interface

Microsoft

# ConvNet

e.g. Google Photos search

Face Verification, Taigman et al. 2014 (FAIR)

[Goodfellow et al. 2014]

Self-driving cars

# Image Classification

- Task of taking an input image and outputting a class
- Probability of classes that best describes the image
- For humans, effortless task



What We See



What Computers See

Microsoft

# Input Image

- An image is an an array of pixel values
- A JPG color image with size 480 x 480:
  - The representative array will be 480 x 480 x 3. Each number is given a value from 0 to 255 which is the pixel intensity
- Grey scale image contains a single sample (intensity value) for each pixel
- Image Classification:
  - Given an array of numbers, produce probabilities of the image being a certain class

Microsoft

# Convolutional Neural Networks

Three basic ideas:
- Local receptive fields
- Shared weights
- Pooling

Microsoft

# Local Receptive Fields

- Connections are from small, localized regions of the input image to hidden layers
- A little window on the input pixels
- Each neuron in the first hidden layer is connected to a small region of the input neurons. For example, a 5x5 region

# Local Receptive Fields

# Pooling Layers

Often used immediately after convolutional layers

- Simplify the information in the output from the convolutional layer
- Takes each feature map output from the convolutional layer and prepares a condensed feature map
- Max-pooling:

    A pooling unit simply outputs the maximum activation in the $p \times p$ region



hidden neurons (output from feature map)

max-pooling units

Microsoft

# Max Pooling



Single depth slice

max pool with 2x2 filters and stride 2

# Pooling Layers

- Smaller representations and more manageable
- Operates over each activation map independently

# Pooling Layers

Combined convolutional and max-pooling layers:

# Shared Weights And Biases

1. Each hidden neuron has a bias and pxp weights connected to its local receptive field

2. The same weights and bias for each of the hidden neurons

3. In other words, for the $j, k^{th}$ hidden neuron, the output is:

$$\sigma\left(b + \sum_{l=0}^{n}\sum_{m=0}^{n} w_{l,m}a_{j+l,k+m}\right)$$

where $\sigma$ is the neural activation function - perhaps the sigmoid function
$b$ is the shared value for the bias
$w_{l,m}$ is a $n \times n$ array of shared weights
$a_{x,y}$ to denote the input activation at position $x, y$

Microsoft

# Shared Weights And Biases

- Convolutional networks are well adapted to the translation invariance of images
- Greatly reduces the number of parameters involved in a convolutional network ($p \times p + b$)

Terminology

| Feature map/Activation map | Map from the input layer to the hidden layer |
|---|---|
| Shared weights | Weights defining the feature map |
| Shared bias | Bias defining the feature map |
| Kernel/Filter | Shared weights and bias |

Microsoft

# First Layer – High Level Perspective

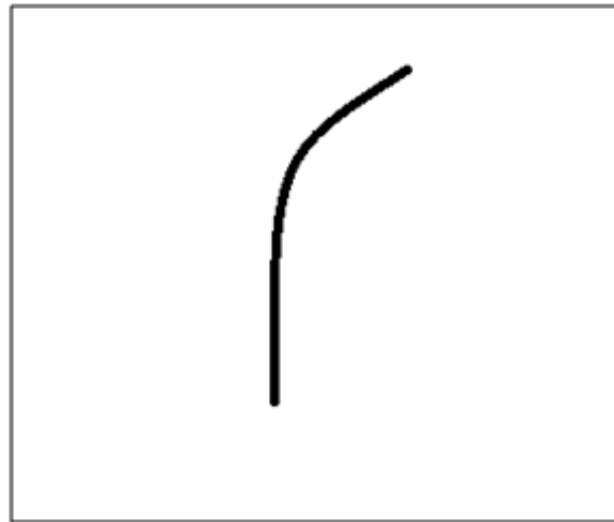- Filters can be thought of as feature identifiers (straight edges, simple colors, and curves)
- In the simple case of a one filter convolution and a curve detector filter, activation map results in regions that are most likely curves in the picture

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Pixel representation of filter**

**Visualization of a curve detector filter**

- More filters mean greater the depth of the activation map
- This results in more information about the input

Microsoft

# Going Deeper Through The Network

- Many layers are interspersed between convolution layers (example: ReLu and Dropout)
- Introduction of nonlinearities
- Improve the robustness of the network and control overfitting

Input -> Conv -> ReLU -> Conv -> ReLU -> Pool -> ReLU -> Conv -> ReLU -> Pool ->Fully Connected
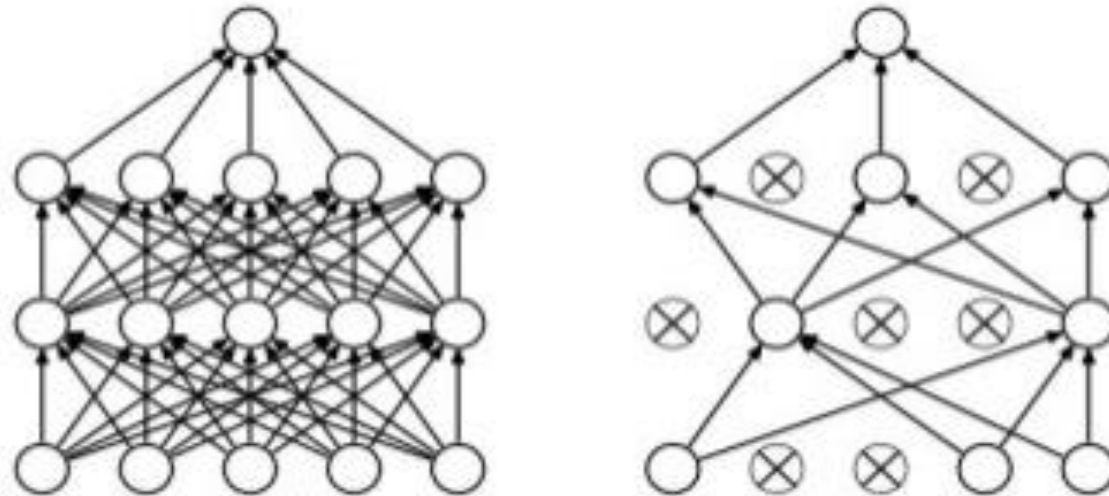
Microsoft

# RelU

- The Rectified Linear Unit has become very popular recently

$$f(x) = \max(0, x)$$

- Activation is simply thresholded at zero
- It was found to greatly accelerate (Krizhevsky et al.) the convergence of stochastic gradient descent compared to the sigmoid/tanh functions
- Compared to tanh/sigmoid neurons that involve expensive operations (exponentials, etc.), the ReLU is simply thresholding

Microsoft

# Dropout

- A form of ensemble learning
- Avoids overfitting (by preventing inter-dependencies from emerging between nodes)
- Dropout - an extreme version of bagging
- At each training step, the dropout procedure creates a different network by removing some neurons randomly



Microsoft

# Fully Connected Layer (FC layer)