# COMP 472 Artificial Intelligence Machine Learning Intro to Neural Networks & Perceptrons

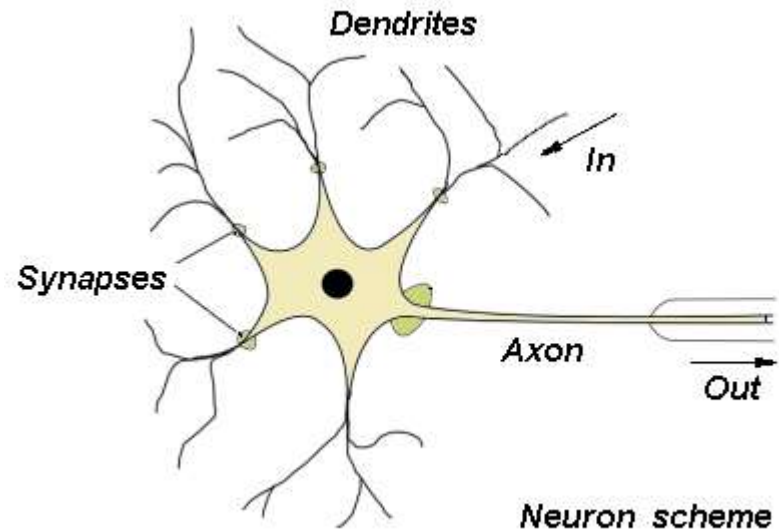- Russell & Norvig: Sections 19.1, 21.1

# Today

# Neural Networks

- Learning approach inspired by biology
    - the neurons in the human brain

- Set of many simple processing units (called neurons) connected together
    - the behavior of each neuron is very simple
    - but a network of neurons can have sophisticated behavior and can be used for complex tasks
    - neurons are connected to each other to form a network
    - the strength of the connection between neurons are determined by weights    connection neurons
    - the network learns by learning the weights between neurons given training data

- Different types of network architectures exist
    - feed forward neural networks (FFNN)
    - recurrent neural networks (RNN)
    - convolutional neural networks (CNN)
    - …

# Biological Neurons

- Human brain =
  - 100 billion neurons
  - each neuron may be connected to 10,000 other neurons
  - passing signals to each other via 1,000 trillion **synapses**

- A neuron is made of:
  - Dendrites: filaments that provide input to the neuron
  - Axon: sends an output signal
  - Synapses: connection with other neurons – releases neurotransmitters to other neurons



Dendrites

In

Synapses

Axon

Out

Neuron scheme

# Behavior of a Neuron

neuron   nabour

- A neuron receives inputs from its neighbors
- If enough inputs are received at the same time:

  input
  - the neuron is activated
  neuron
  - and fires an output to its neighbors

  fire output        neighbor
- Repeated firings across a synapse increases its sensitivity and the future likelihood of its firing

  firing
- If a particular stimulus repeatedly causes activity in a group of neurons, they become strongly associated

# A Perceptron



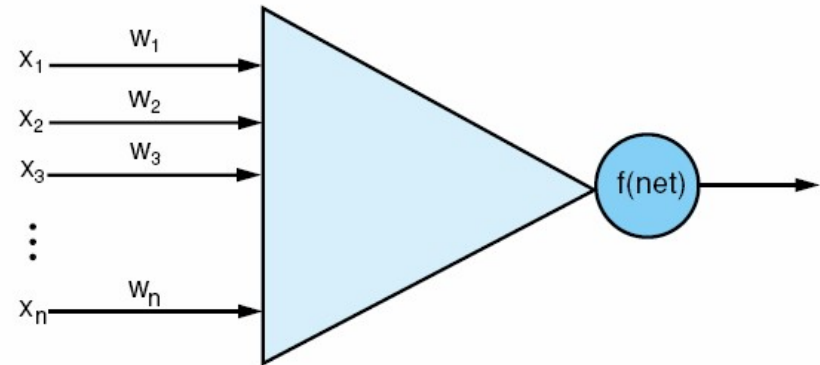- A <u>single</u> computational neuron (no network yet…)

- Input:
  - input signals $x_i$    signal
  - weights $w_i$ for each feature $x_i$    xi
    - represents the strength of the connection with the neighboring neurons
- Output:    weight, fire,output=1
  - if sum of input weights >= some threshold, neuron fires (output=1)
  - otherwise output = 0
    - If $(w_1 x_1 + … + w_n x_n)$ >= $t$
    - Then output = 1
    - Else output = 0

- Learning :    training data    perception    weights
  - use the training data to adjust the weights in the perceptron

source: Luger (2005)

# The Idea

| | Features ($x_i$) | | | | Output |
|---|---|---|---|---|---|
| **Student** | **First last year?** | **Male?** | **Works hard?** | **Drinks?** | **First this year?** |
| Richard | Yes | Yes | No | Yes | **No** |
| Alan | Yes | Yes | Yes | No | **Yes** |
| … | | | | | |

First last year?  1  0.2

Male?  1  0.2

Works hard?  0  0.2  First this year?

Drinks a lot?  1  0.2

1. **Step 1:** Set weights to random values
2. **Step 2:** Feed perceptron with an input
3. **Step 3:** Compute the network outputs
4. **Step 4:** Adjust the weights
   1. if output correct → weights stay the same
   2. if output = 0 but it should be 1 →
      - increase weights on <u>active</u> connections (i.e. input $x_i$ =1)
   3. if output = 1 but should be 0 →
      - decrease weights on <u>active</u> connections (i.e. input $x_i$ =1)
5. **Step 5:** Repeat steps 2 to 4 a large number of times until the network converges to the right results for the given training examples

source: Cawsey (1998)

# A Simple Example

- Turn feature values into numerical values
  - yes -> 1   no -> 0
  - e.g. if $x_1$ = 1, then student got an A last year
  - e.g. if $x_1$ = 0, then student did not get an A last year
- Initially, set all weights to random values (all 0.2 here)

First last year?   1     0.2

Male?   1     0.2          First this year?

0.2

Works hard?   0     0.2

Drinks a lot?   1

- Assume:
  0. 55
  - threshold = 0.55
  - constant learning rate = 0.05
    0. 05

# A Simple Example (2)

| Student | Features ($x_i$) | | | | Output |
| | 'A' last year? | Male? | Works hard? | Drinks? | 'A' this year? |
|---------|-----------------|-------|-------------|---------|----------------|
| Richard | 1 | 1 | 0 | 1 | 0 |
| Alan | 1 | 1 | 1 | 0 | 1 |
| Alison | 0 | 0 | 1 | 0 | 0 |
| Jeff | 0 | 1 | 0 | 1 | 0 |
| Gail | 1 | 0 | 1 | 1 | 1 |
| Simon | 0 | 1 | 1 | 1 | 0 |

- **Richard:**
  - $(1 \times 0.2) + (1 \times 0.2) + (0 \times 0.2) + (1 \times 0.2) = 0.6 >= 0.55$ --> output is 1
  - ...but he did not get an A this year
  - So reduce weights of all <u>active</u> connections (with input 1) by 0.05. Do not change the weight of the inactive connections.
  - So we get $w_1 = 0.15$, $w_2 = 0.15$, $w_3 = 0.2$, $w_4 = 0.15$

# A Simple Example (3)

| | Features ($x_i$) | | | | Output |
|---|---|---|---|---|---|
| **Student** | 'A' last year? | Male? | Works hard? | Drinks? | **'A' this year?** |
| ~~Richard~~ | ~~1~~ | ~~1~~ | ~~0~~ | ~~1~~ | ~~0~~ |
| Alan | 1 | 1 | 1 | 0 | **1** |
| Alison | 0 | 0 | 1 | 0 | **0** |
| Jeff | 0 | 1 | 0 | 1 | **0** |
| Gail | 1 | 0 | 1 | 1 | **1** |
| Simon | 0 | 1 | 1 | 1 | **0** |

- Alan:
  - $(1 \times 0.15) + (1 \times 0.15) + (1 \times 0.2) + (0 \times 0.15) = 0.5 < 0.55$ → output is 0
    <br>increase
  - … but expected output is 1
  - So increase all weights of active connections by 0.05
  - So we get $w_1$= 0.2, $w_2$= 0.2, $w_3$= 0.25, $w_4$= 0.15

- Alison…Jeff… Gail… Simon…

# A Simple Example (4)

| | Features ($x_i$) | | | | Output |
|---|---|---|---|---|---|
| **Student** | 'A' last year? | Male? | Works hard? | Drinks? | 'A' this year? |
| ~~Richard~~ | 1 | 1 | 0 | 1 | 0 |
| ~~Alan~~ | 1 | 1 | 1 | 0 | 1 |
| ~~Alison~~ | 0 | 0 | 1 | 0 | 0 |
| ~~Jeff~~ | 0 | 1 | 0 | 1 | 0 |
| ~~Gail~~ | 1 | 0 | 1 | 1 | 1 |
| ~~Simon~~ | 0 | 1 | 1 | 1 | 0 |

- epoch 1:  Richard, Alan, Alison, Jeff, Gail, Simon
- epoch 2: Richard, Alan, Alison, Jeff, Gail, Simon
- ...
- epoch n... until all training data points are correctly classified
    - $w_1$= 0.25 $w_2$= 0.1 $w_3$= 0.2 $w_4$= 0.1

Output

# A Simple Example (5)

| | Features ($x_i$) | | | | Output |
|---|---|---|---|---|---|
| **Student** | 'A' last year? | Male? | Works hard? | Drinks? | **'A' this year?** |
| Richard | 1 | 1 | 0 | 1 | **0** |
| Alan | 1 | 1 | 1 | 0 | **1** |
| Alison | 0 | 0 | 1 | 0 | **0** |
| Jeff | 0 | 1 | 0 | 1 | **0** |
| Gail | 1 | 0 | 1 | 1 | **1** |
| Simon | 0 | 1 | 1 | 1 | **0** |

- Let's check… ($w_1$= 0.2 $w_2$= 0.1 $w_3$= 0.25 $w_4$= 0.1)
  - Richard: $(1 \times 0.2) + (1 \times 0.1) + (0 \times 0.25) + (1 \times 0.1) = 0.4 < 0.55$ -> output is 0 ✓
  - Alan:    $(1 \times 0.2) + (1 \times 0.1) + (1 \times 0.25) + (0 \times 0.1) = 0.55 \geq 0.55$ -> output is 1 ✓
  - Alison:  $(0 \times 0.2) + (0 \times 0.1) + (1 \times 0.25) + (0 \times 0.1) = 0.25 < 0.55$ -> output is 0 ✓
  - Jeff:    $(0 \times 0.2) + (1 \times 0.1) + (0 \times 0.25) + (1 \times 0.1) = 0.2 < 0.55$ -> output is 0 ✓
  - Gail:    $(1 \times 0.2) + (0 \times 0.1) + (1 \times 0.25) + (1 \times 0.1) = 0.55 \geq 0.55$ -> output is 1 ✓
  - Simon:   $(0 \times 0.2) + (1 \times 0.1) + (1 \times 0.25) + (1 \times 0.1) = 0.45 < 0.55$ -> output is 0 ✓

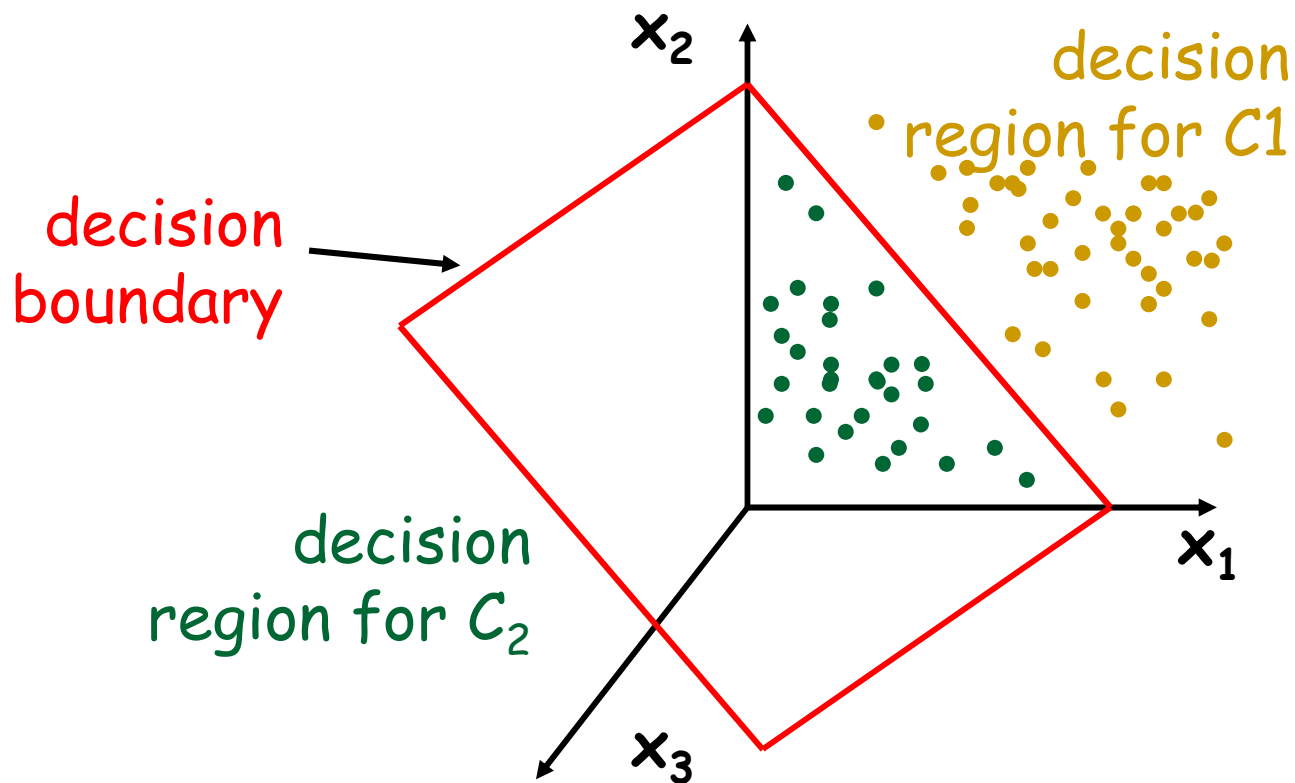# Decision Boundaries of Perceptrons

- So we have just learned the function:
  - If ($0.2x_1 + 0.1x_2 + 0.25x_3 + 0.1x_4 \geq 0.55$) then 1 otherwise 0
  - If ($0.2x_1 + 0.1x_2 + 0.25x_3 + 0.1x_4 - 0.55 \geq 0$) then 1 otherwise 0
- Assume we only had 2 features:
  - If ($w_1x_1 + w_2x_2 - t >= 0$) then 1 otherwise 0              input              t
  - The learned function describes a line in the input space
  - This line is used to separate the two classes C1 and C2
  - t (the threshold, later called 'b') is used to shift the line on the axis

$x_2$

$w_1x_1 + w_2x_2 - t = 0$   decision boundary

$w_1x_1 + w_2x_2 - t \geq 0$

decision region for C1

$C_1$

$C_2$

decision region for $C_2$

$w_1x_1 + w_2x_2 - t < 0$

$x_1$

13

# Decision Boundaries of Perceptrons

- More generally, with n features, the learned function describes a hyperplane in the input space.

hyperplane



decision region for C1

decision boundary

$x_2$

$x_1$

$x_3$

decision region for $C_2$

# Adding a Bias

- We can avoid having to figure out the threshold by using a "bias"

$$b + \sum_i x_i w_i$$

bias          threshold

- A bias is equivalent to a weight on an extra input feature that always has a value of 1.

bias          1          input feature



$b$          $w_1$          $w_2$

$1$          $x_1$          $x_2$

# Perceptron - More Generally

inputs

$x_0 = 1$

bias input set to 1 (*to replace the threshold*)

$X_1$

$w_1$

$w_0$

$w_2$

output

$X_2$

$$\sum_{i=0}^{n} w_i x_i$$

$O$

$w_n$

$X_n$

transfer function

$$f\left( \sum_{i=0}^{n} w_i x_i \right)$$

$$O = f\left( \sum_{i=0}^{n} w_i x_i \right)$$

activation function

final classification

# Common Activation Functions

- ### step

$$O = \begin{cases} if \left( \sum_{i=1}^{n} w_i x_i \right) \geq t \to 1 \\ \\ otherwise \to 0 \end{cases}$$

function step 1
0



(a) Step function

- ### sign

$$O = \begin{cases} if \left( \sum_{i=0}^{n} w_i x_i \right) \geq 0 \to 1 \\ \\ otherwise \to -1 \end{cases}$$

sign 0 1
0 -1



(b) Sign function

$$O = \begin{cases} if \left( \sum_{i=0}^{n} w_i x_i \right) > 0) \to 1 \\ \\ if \left( \sum_{i=0}^{n} w_i x_i \right) = 0) \to 0 \\ \\ otherwise \to -1 \end{cases}$$

[Russell & Norvig, 1995]

# Learning Rate

1. Learning rate can be a constant value (as in the previous example)

$$\Delta w = \eta(T - O)$$

learning rate

Error = target output – actual output

learning rate

- So:
  - if T=zero and O=1 (i.e. a false positive) -> decrease w by η
  - if T=1 and O=zero (i.e. a false negative) -> increase w by η
  - if T=O (i.e. no error) -> don't change w

2. Or, a fraction of the input feature $x_i$

$$\Delta w_i = \eta(T - O) \, x_i$$

value of input feature $x_i$

Xi

- So the update is proportional to the value of x
  - if T=zero and O=1 (i.e. a false positive) -> decrease $w_i$ by $\eta x_i$      Xi
  - if T=1 and O=zero (i.e. a false negative) -> increase $w_i$ by $\eta x_i$
  - if T=O (i.e. no error) -> don't change $w_i$

- This is called the delta rule or perceptron learning rule

# Perceptron Convergence Theorem

- **If a solution with zero error exists**
  - i.e. the training data are linearly separable
- **The delta rule will find a solution in finite time**

0error

delta rule

# Example of the Delta Rule

## ■ training data:

| $x_1$ | $x_2$ | Output |
|-------|-------|--------|
| 1.0 | 1.0 | 1 |
| 9.4 | 6.4 | −1 |
| 2.5 | 2.1 | 1 |
| 8.0 | 7.7 | −1 |
| 0.5 | 2.2 | 1 |
| 7.9 | 8.4 | −1 |
| 7.0 | 7.0 | −1 |
| 2.8 | 0.8 | 1 |
| 1.2 | 3.0 | 1 |
| 7.8 | 6.1 | −1 |

## ■ plot of the training data:

## ■ perceptron



source: Luger (2005)

# Example of the Delta Rule

- **assume random initialization**
  - w1 = 0.75
  - w2 = 0.5
  - w3 = -0.6

    initialization
    sign function
    w

- **Assume:**
  - sign function (threshold = 0)
  - learning rate η = 0.2

source: Luger (2005)

# Example of the Delta Rule

- data #1: $f(0.75 \times 1 + 0.5 \times 1 - 0.6 \times 1) = f(0.65) \to 1$   ✓

- data #2: $f(0.75 \times 9.4 + 0.5 \times 6.4 - 0.6 \times 1) = f(9.65) \to 1$ ✗

  - $\to$ error = $(-1 - 1) = -2$
    
    label -   label
    
    $\to w_1 = w_1 - 2 \times 0.2 \times 9.4$   x
    
          $= 0.75 - 3.76 = -3.01$
  
  $\to w_2 = w_2 - 2 \times 0.2 \times 6.4 = -2.06$
  
  $\to w_3 = w_3 - 2 \times 0.2 \times 1 = -1.00$

| $x_1$ | $x_2$ | Output |
|-------|-------|--------|
| 1.0 | 1.0 | 1 |
| 9.4 | 6.4 | −1 |
| 2.5 | 2.1 | 1 |
| 8.0 | 7.7 | −1 |
| 0.5 | 2.2 | 1 |

- data #3: $f(-3.01 \times 2.5 - 2.06 \times 2.1 - 1 \times 1) = f(-12.84) \to -1$ ✗

  - $\to$ error = $(1 - -1) = 2$
    
    $\to w_1 = -3.01 + 2 \times 0.2 \times 2.5 = -2.01$
    
    $\to w_2 = -2.06 + 2 \times 0.2 \times 2.1 = -1.22$
    
    $\to w_3 = -1.00 + 2 \times 0.2 \times 1 = -0.60$

- repeat… over 500 iterations, we converge to:

  $w_1 = -1.3$   $w_2 = -1.1$   $w_3 = 10.9$

source: Luger (2005)

# The Perceptron in 1958



An IBM 704 – a 5-ton computer the size of a room – was fed a series of punch cards.
After 50 trials, the computer taught itself to distinguish cards marked on the left from cards marked on the right.

Frank Rosenblatt

https://news.cornell.edu/stories/2019/09/professors-perceptron-paved-way-ai-60-years-too-soon

# Remember this slide?

## History of AI

- Reality hits (late 60s - early 70s)
  - 1966: the ALPAC report kills work in machine translation (and NLP in general)
  - People realized that scaling up from micro-worlds (toy-worlds) to reality is not just a manner of faster machines and larger memories...
  - Minsky & Papert's paper on the limits of perceptrons (cannot learn just any function...) kills work in neural networks
  - in 1971, the British government stops funding research in AI due to no significant results
  - it's the first major AI Winter...

https://www.vectorstock.com/royalty-free-vector/freezing-snowman-vector-689086

32

24

# Limits of the Perceptron

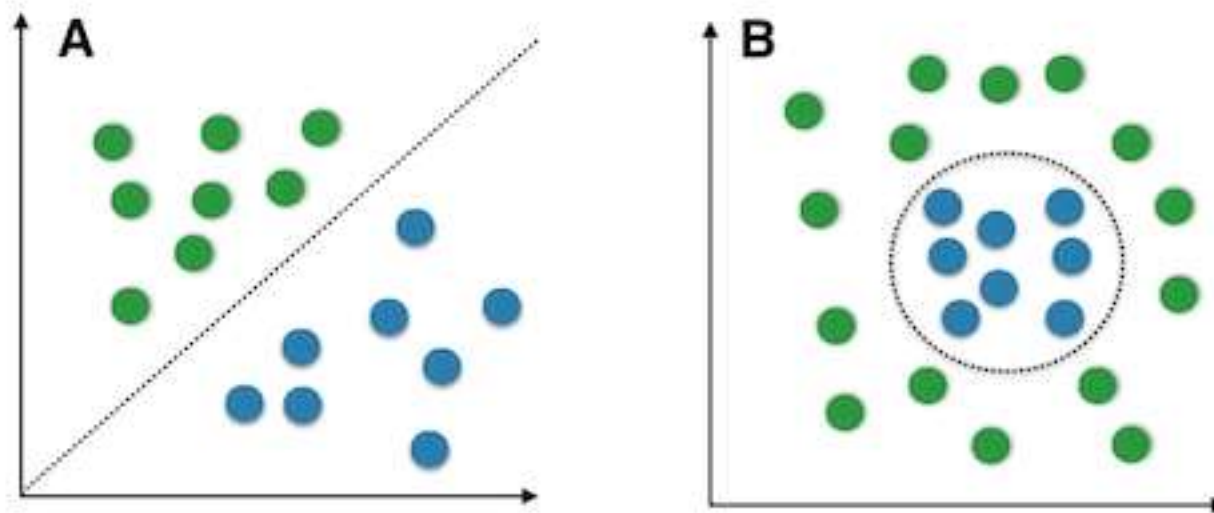- In 1969, Minsky and Papert showed formally what functions could and could not be represented by perceptrons

- Only linearly separable functions can be represented by a perceptron



(a) $I_1$ **and** $I_2$

(b) $I_1$ **or** $I_2$

(c) $I_1$ **xor** $I_2$

source: Luger (2005)

# AND and OR Perceptrons



| x | y | $x + y - 2$ | Output |
|---|---|---|---|
| 1 | 1 | 0 | 1 |
| 1 | 0 | -1 | -1 |
| 0 | 1 | -1 | -1 |
| 0 | 0 | -2 | -1 |

source: Luger (2005)

# The XOR Function - Visually

- In a 2-dimentional space (2 features for the X)
- No straight line in two-dimensions can separate
  - (0, 1) and (1, 0) from
  - (0, 0) and (1, 1).

$x_1$

$(0, 1)$    $(1, 1)$

$(0, 0)$    $(1, 0)$    $x_2$

source: Luger (2005)

# A Perceptron Network



$I_j$     $W_j$     $O$

Input Units     Output Unit

**Single Perceptron**

- A perceptron is a binary classifier (i.e. 2 classes)
- if the output needs to learn more than a binary decision
- we can have a network of perceptrons

- Eg: learning to recognize digit --> 10 possible outputs --> need a perceptron network



$I_j$     $W_{j,i}$     $O_i$

Input Units     Output Units

**Perceptron Network**

# Non-Linearly Separable Functions

- Real-world problems cannot always be represented by linearly-separable functions...



Linear vs. nonlinear problems

- This caused a decrease in interest in neural networks in the 1970's

http://sebastianraschka.com/Articles/2014_kernel_pca.html

# Today

1. Introduction to ML
2. Naïve Bayes Classification
    a. Application to Spam Filtering
3. Decision Trees
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
    a. Perceptrons
    b. Multi Layered Neural Networks

# Up Next

1. Introduction to ML
2. Naïve Bayes Classification
   a. Application to Spam Filtering
3. Decision Trees
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
   a. Perceptrons
   b. **Multi Layered Neural Networks**