# Machine Learning

# Artificial Neural Networks

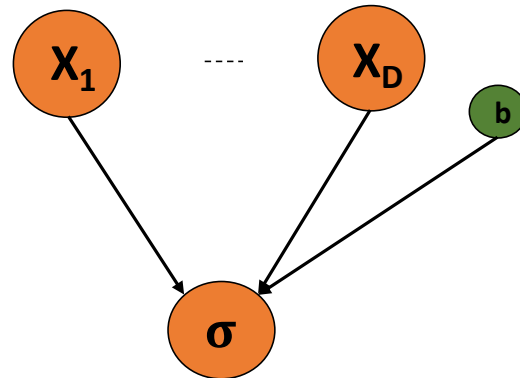| mite | container ship | motor scooter | leopard |
|------|---------------|---------------|---------|
| mite | container ship | motor scooter | leopard |
| black widow | lifeboat | go-kart | jaguar |
| cockroach | amphibian | moped | cheetah |
| tick | fireboat | bumper car | snow leopard |
| starfish | drilling platform | golfcart | Egyptian cat |

**Overview:**

- Artificial Neuron

- Activation Function

- Neuron Capacity

- Single Layer ANN

- Multilayer ANN

- Universal Approximator

- Motivation behind ANN

- Backpropagation

**- Hidden unit pre-activation:**

$$z(x) = \sum_i w_i x_i + b = W^T X + b$$

**- Hidden unit activation:**

$$f(x) = \sigma\left(z(x)\right) = \sigma\left(\sum_i w_i x_i + b\right)$$

**W** are weight matrix connects $i^{th}$ hidden unit with $i^{th}$ input unit

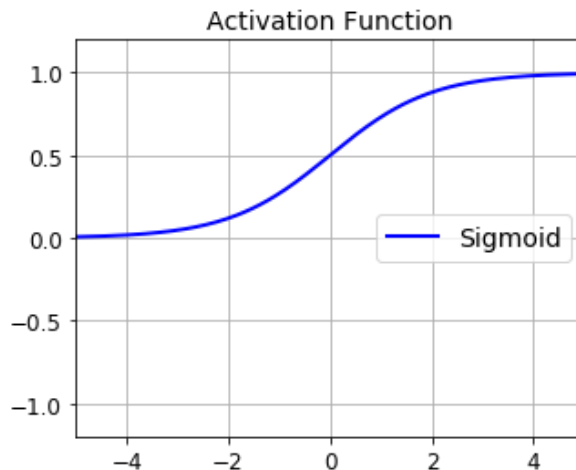**b** are bias vectors

**σ(·)** is the activation function

**Sigmoid:**

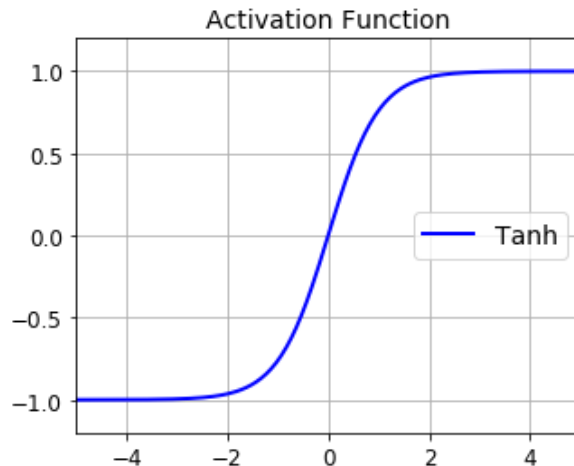$$\sigma(z) = \text{sigm}(z) = \frac{1}{1+\exp(-z)}$$

-   Squashes the hidden unit's pre-activation to between 0 and 1.
-   Always positive.
-   Bounded.
-   Strictly increasing.

**Hyperbolic tangent ("tanh"):**

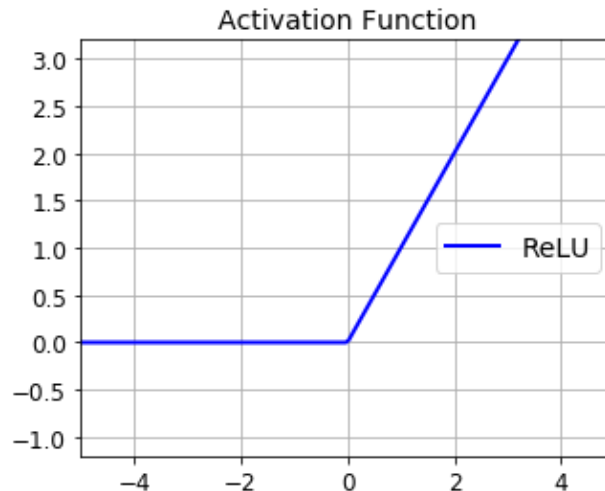$$\sigma(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} = \frac{\exp(2z) - 1}{\exp(2z) + 1}$$

- Squashes the hidden unit's pre-activation to between -1 and 1.
- Can be positive or negative.
- Bounded.
- Strictly increasing.



Activation Function

**Rectified Linear Unit ("ReLu"):**

$\sigma(z) = \text{ReLu}(z) = \max(0, z)$

- Bounded below by 0 (always non-negative).

- Not bounded above.

- Strictly increasing.

**- Hidden layer pre-activation:**

$$z(x)_i = b_i^{(1)} + W_{i,j}^{(1)} x_j$$

*Similarly in Matrix form:*
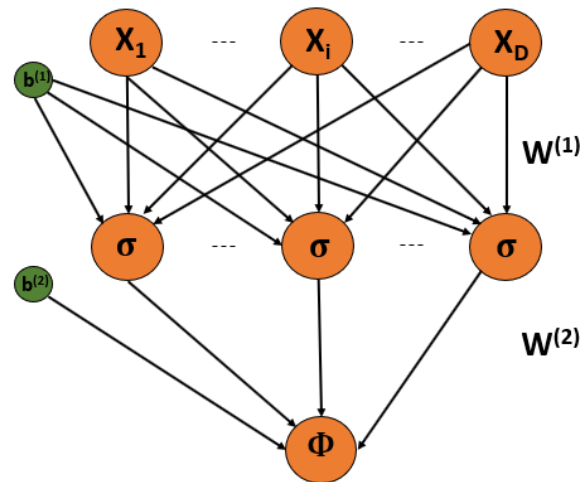
$$z(X) = b^{(1)} + \mathbf{W}^{(1)} X$$

**- Hidden layer activation:**

$$\mathbf{h}(\mathbf{X}) = \boldsymbol{\sigma}(\mathbf{z}(\mathbf{X}))$$

**- Output layer activation "Φ":**

$$F(X) = \Phi(b^{(2)} + \mathbf{W}^{(2)^T} h^{(1)}X)$$



9

**SoftMax Activation Function**

- Multi-class classification:

    - requires multiple outputs i.e. 1 output per class.

    - need to estimate the conditional probability of output belonging to a particular class $c$, $p(y = c | \mathbf{x})$.

- Apply the SoftMax activation function at the output:

    $$\Phi(z) = \text{SoftMax}(z) = [\; \frac{e^{z_1}}{\sum_c e^{z_c}} \;,\; \dots \;,\; \frac{e^{z_1}}{\sum_c e^{z_c}} \;]^\mathsf{T}$$

    - strictly positive

    - sums to one

- Predicted class is the one with highest estimated probability

**Multilayer NN with _L_ hidden layers**
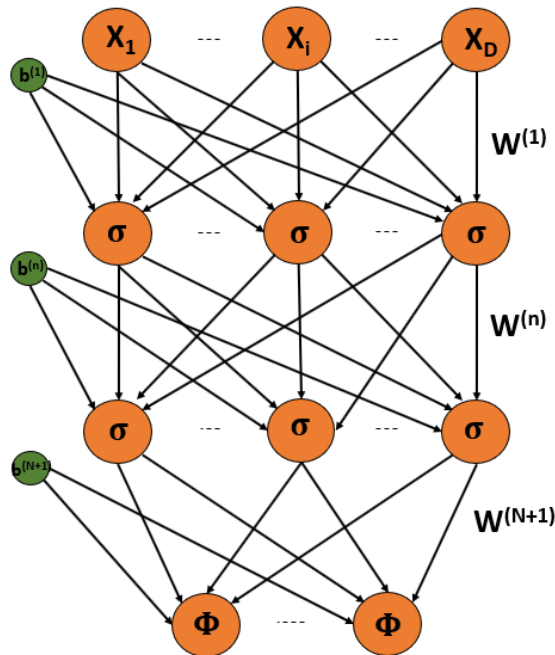
- **Hidden layer pre-activation for $k > 0$ :**

  $\mathrm{h}^{(0)}(\mathrm{X}) = \mathrm{X},$

  $\mathrm{z}^{(k)}(\mathrm{X}) = \mathrm{b}^{(k)} + \mathbf{W}^{(k)}\,\mathrm{h}^{(k-1)}\mathrm{X}$

- **Hidden layer activation (_k_ from 1 to _L_):**

  $\mathbf{h}^{(k)}(\mathrm{X}) = \boldsymbol{\sigma}(\mathbf{z}^{(k)}(\mathrm{X}))$

- **Output layer activation (_k = L+1_):**

  $\mathbf{F}(\mathbf{X}) = \mathbf{h}^{(L+1)}(\mathrm{X}) = \Phi(\mathrm{z}^{(L+1)}(\mathrm{X}))$

**Capacity of Single Hidden Unit**

- **Range of hidden unit determined by $\sigma(.)$**

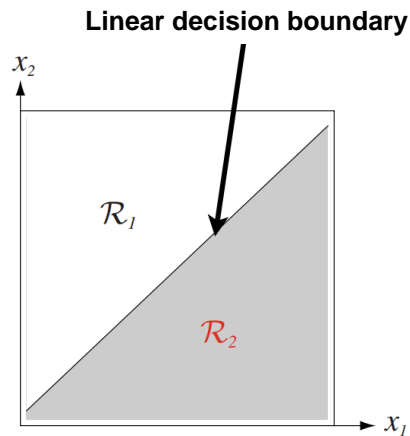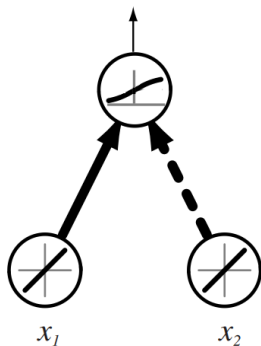- **Bias $b$ changes the position of the riff.**



Bias only changes
the position of the
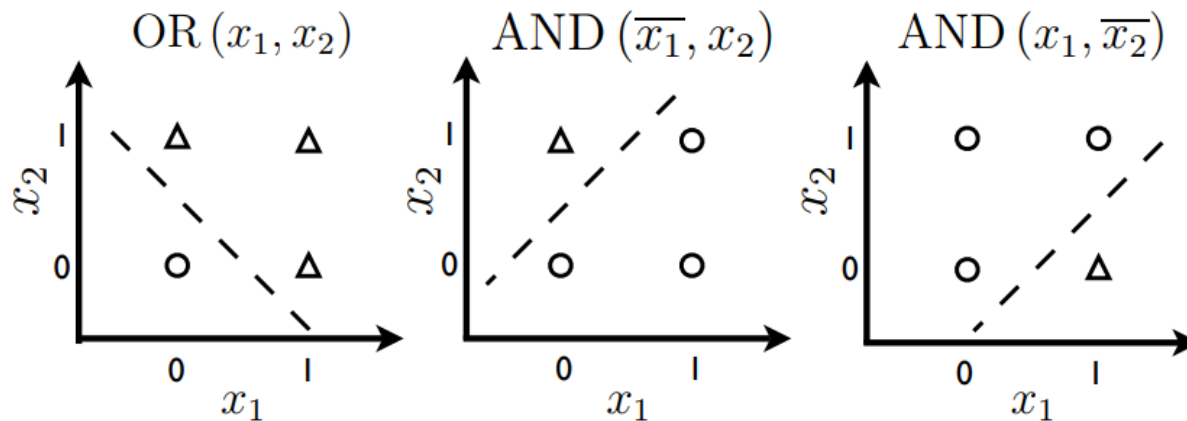riff

Source: Pascal Vincent

**Capacity of Single Hidden Unit**

- Could do binary classification.

- With sigmoid, can interpret neuron as estimating $p(y = 1 | \mathrm{x})$.

- Also known as logistic regression classifier, if greater than 0.5 predict class 1, otherwise predict class 0.
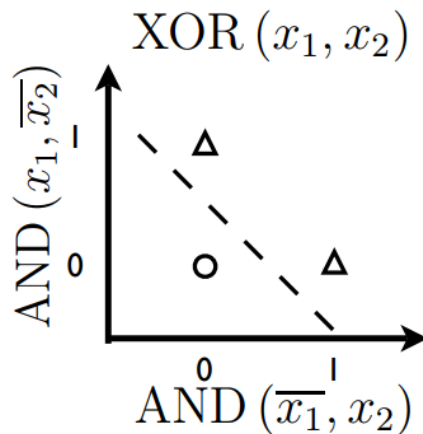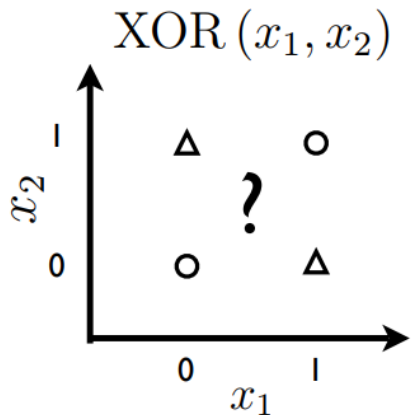


**Linear decision boundary**

**Capacity of Single Hidden Unit**

- **Can solve linearly separable problems**

**Capacity of Single Hidden Unit**

- **Cannot solve non-linearly separable problems**

- **Unless the input is transformed into a separable representation**

**- Hidden layer pre-activation:**

$$z(x)_i = b_i{}^{(1)} + W_{i,j}^{(1)} x_j$$

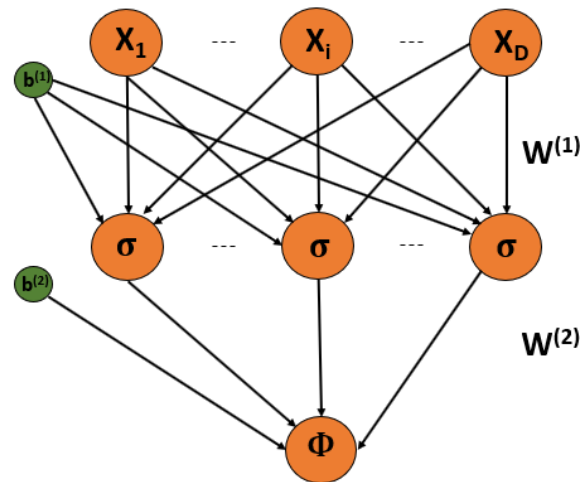*Similarly in Matrix form:*

$$z(X) = b^{(1)} + \mathbf{W}^{(1)} X$$

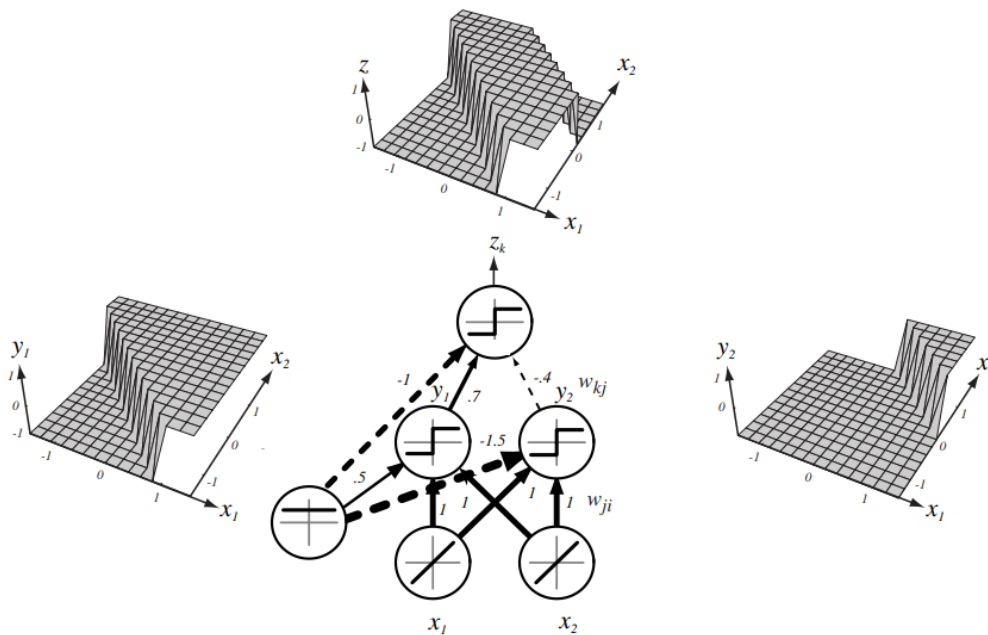**- Hidden layer activation:**

$$\mathbf{h(X)} = \boldsymbol{\sigma}(\mathbf{z(X)})$$

**- Output layer activation "Φ":**

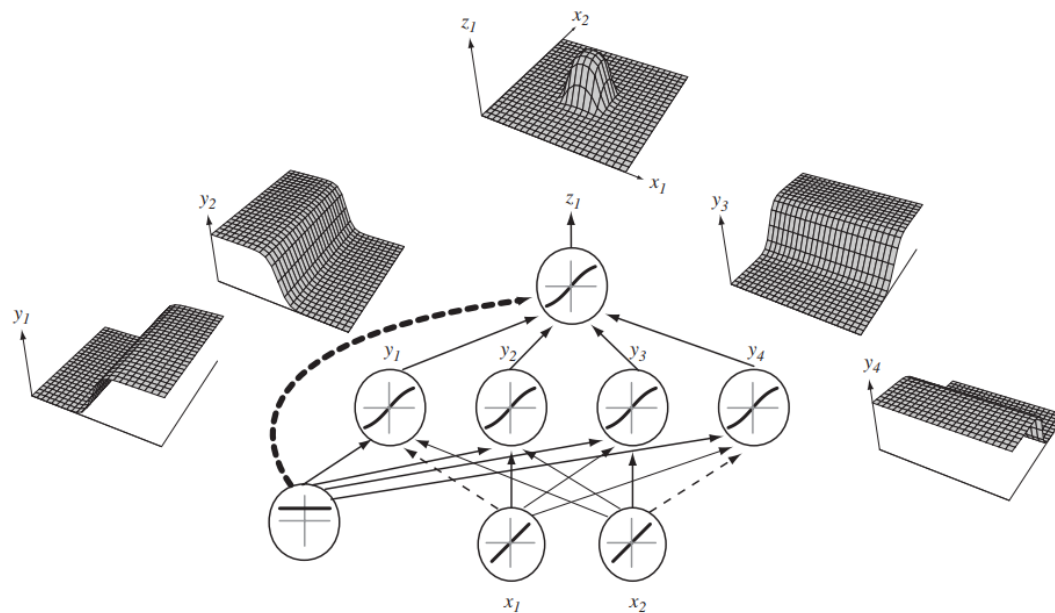$$F(X) = \Phi(b^{(2)} + \mathbf{W}^{(2)^T} h^{(1)}X)$$



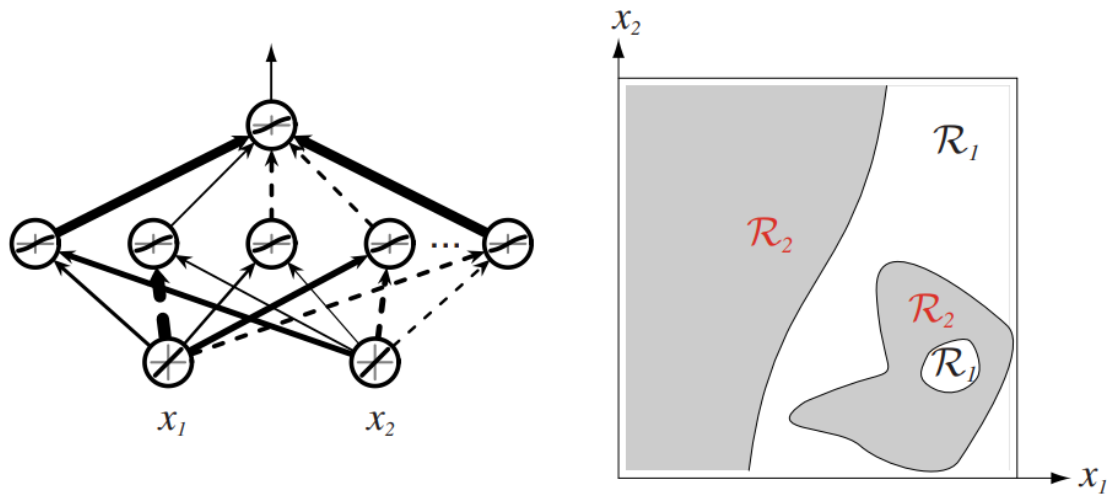16

**Capacity of Single Hidden Layer Neural Network**



Source: Pascal Vincent

**Capacity of Single Hidden Layer Neural Network**



Source: Pascal Vincent

**Capacity of Single Hidden Layer Neural Network**

- **Universal Approximation Theorem (Hornik, 1991)**

  "Single layer feedforward network can approximate any continuous function arbitrarily well if and only if the network's activation function is continuous, non-constant, and bounded."

- Hornik's result applies for sigmoid, tanh and many other hidden layer activation functions.

- However, modern day defacto activation function is ReLu, and it does not satisfy Hornik's theorem as the $\mathrm{ReLu}(z) = \max(0, z)$ is unbounded from above.

- "Multilayer feedforward network can approximate any continuous function arbitrarily well if and only if the network's continuous activation function is not polynomial."

**Definition**

A set $F$ of functions in $L^\infty_{loc}(R^n)$ is dense in $C(R^n)$ if for every function $g \in C(R^n)$ and for every compact set $K \subset R^n$, there exists a sequence of functions $f_j \in F$ such that

$$\lim_{f \to \infty} \|g - f_j\|_{L^\infty(K)} = 0.$$

**Theorem**

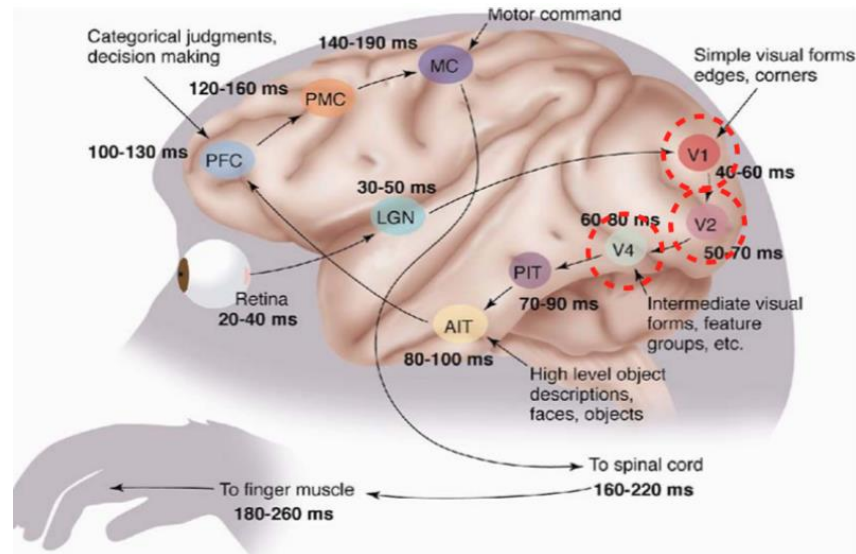(Leshno et al., 1993) Let $\sigma \in M$, where $M$ denotes the set of functions which are in $L^\infty_{loc}(\Omega)$.

$$\Sigma_n = span\{\sigma(w \cdot x + b) : w \in R^n, b \in R\}$$

Then $\Sigma_n$ is dense in $C(R^n)$ if and only if $\sigma$ is not an algebraic polynomial (a.e.).
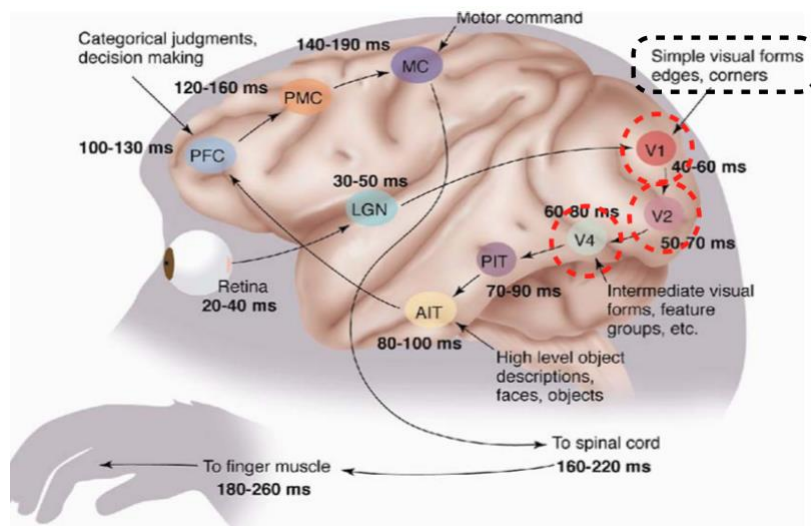
**Biological Inspiration**
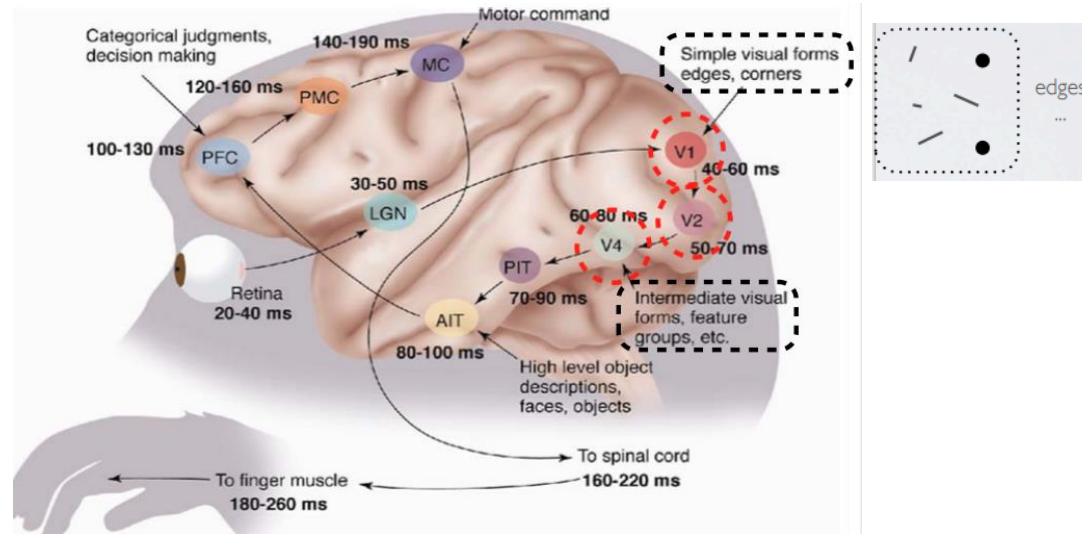
- **Parallel with the visual vortex**



Source: Simon Thorpe

**Biological Inspiration**

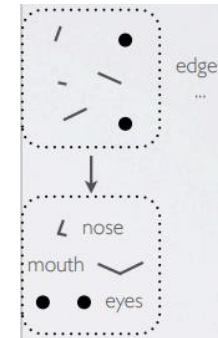**- Parallel with the visual vortex**



Source: Simon Thorpe

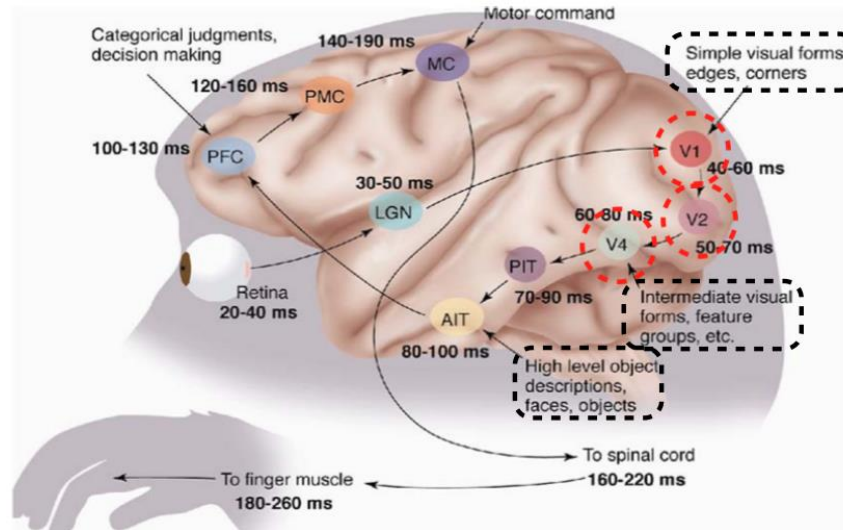**Biological Inspiration**

**- Parallel with the visual vortex**



Source: Simon Thorpe

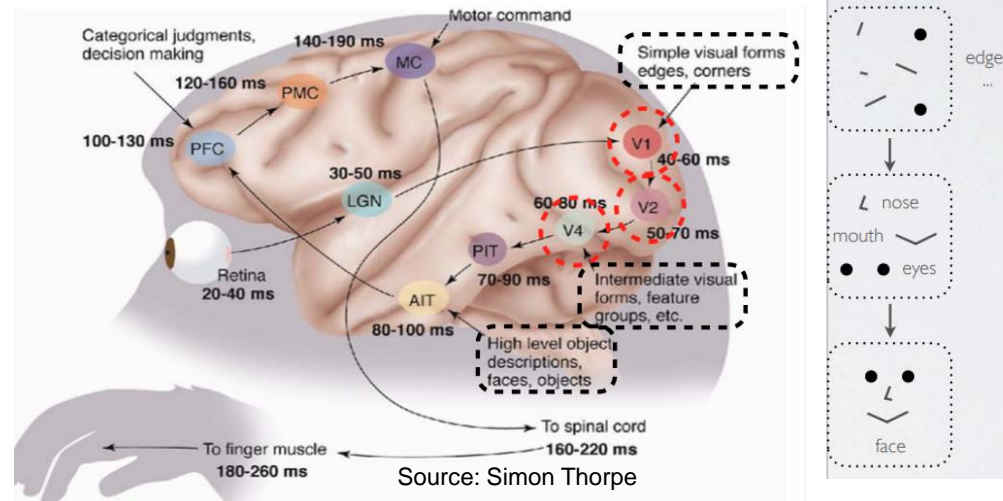NYU

**Biological Inspiration**

- **Parallel with the visual vortex**

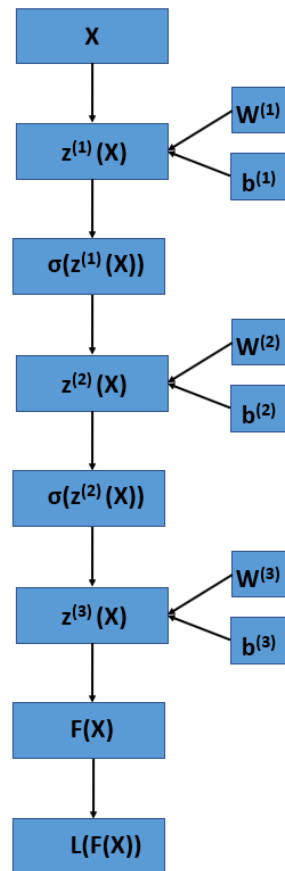- **Edges**



Source: Simon Thorpe

**Biological Inspiration**

- **Parallel with the visual vortex**

- **Edges**

- **Higher level features such as nose, mouth, eyes**

- **Face**



Source: Simon Thorpe

- **Firing rates of different input neurons combine to influence the firing rate of other neurons:**

  - depending on the dendrite and axon, a neuron can either work to increase (excite) or decrease (inhibit) the firing rate of another neuron

- **This is what artificial neurons approximate:**

  - the activation corresponds to a "sort of" firing rate
  - the weights between neurons model whether neurons excite or inhibit each other
  - the activation function and bias model the threshold behavior of action potentials

## Forward propagation

- Randomly Initialize $\Theta$

- $\Theta = \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, ..., \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$
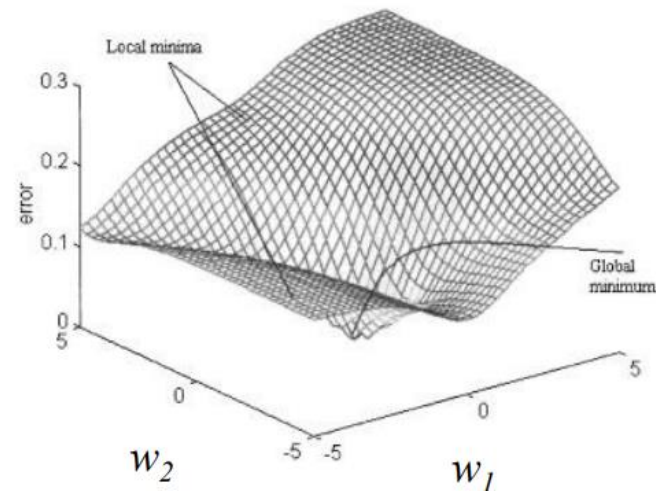
NYU

- Objective Function for Multi-class Classification

$$\underset{\boldsymbol{\theta}}{\text{argmin}} \quad \frac{1}{m} \sum_{i=1}^{m} l(\mathbf{F}(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)}) + \lambda\, \Omega(\boldsymbol{\theta})$$
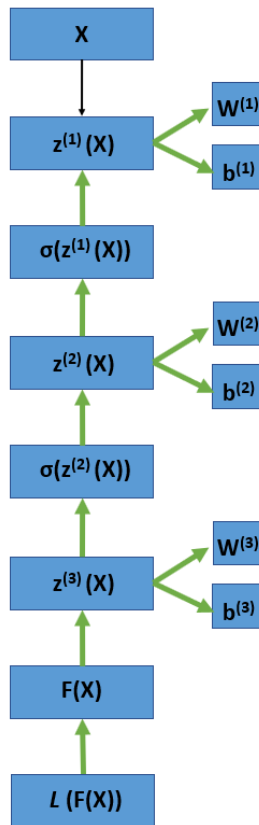
- Loss Function: Negative log likelihood

$$l(\mathbf{F}(\mathbf{x}, \boldsymbol{\theta}), y) = -\sum_c \mathbf{1}_{(y=c)}\, log\, \mathbf{F}(\mathbf{x})_c$$



Cho & Chow, Neurocomputing 1999

29

## Backpropagation

- Optimization Algorithm e.g. stochastic gradient descent
- Objective Function
- Gradient of Output Layer
- Gradient of Hidden Layer
- Gradient of Activation Function
- Gradient of Parameters

**Optimization Algorithm**
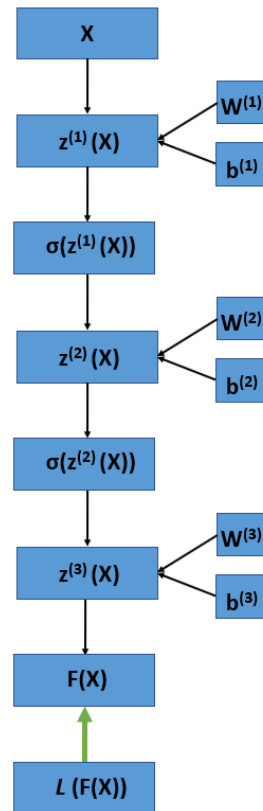
Stochastic (incremental) gradient descent

For each training example ( $\mathrm{x}^{(i)}, y^{(i)}$ )

    For N iterations

- Compute $\Delta = - \nabla_{\boldsymbol{\theta}} \, \boldsymbol{l}(\mathbf{F}(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)}) - \nabla_{\boldsymbol{\theta}} \lambda \boldsymbol{\Omega}(\boldsymbol{\theta})$

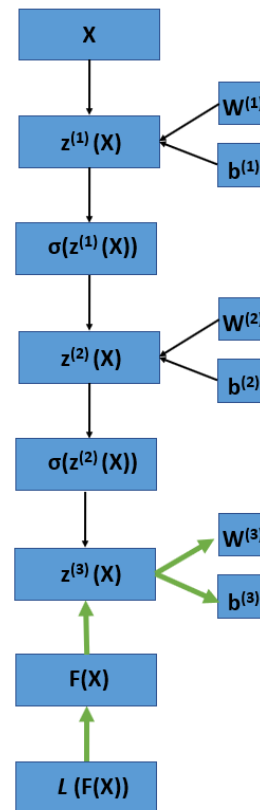- Update $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \, \Delta$

Loss Gradient at Output Layer

- $\nabla_{\mathbf{F(x)}} - \boldsymbol{log}\mathbf{F(x)}_y = \dfrac{-\mathbf{e}(\boldsymbol{y})}{\mathbf{F(x)}_y}$
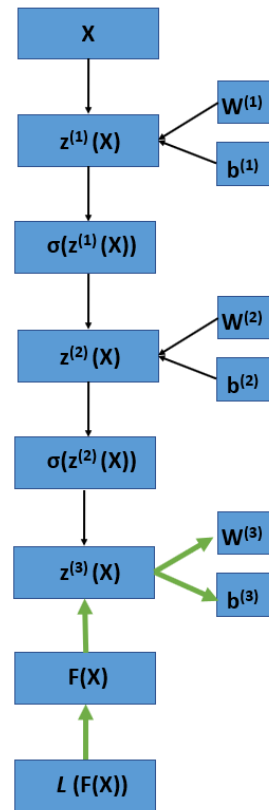
Loss Gradient at Output Pre-activation

- $\nabla_{\mathbf{z}^{(L+1)}(\mathbf{x})} - log\mathbf{F}(\mathbf{x})_y = -(\mathbf{e}(y) - \mathbf{F}(\mathbf{x}))$
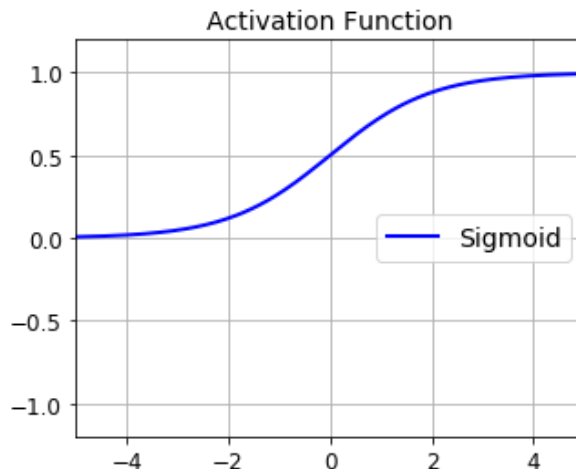
NYU

## Loss Gradient at Hidden Layer

- $\nabla_{z^{(k)}(x)} - log\mathbf{F}(x)_y$

  $= \nabla_{\mathbf{h}^{(k)}(\mathbf{x})} - log\mathbf{F}(\mathbf{x})_y \odot [\dots, \sigma'(z^{(k)}(\mathbf{x})), \dots]$
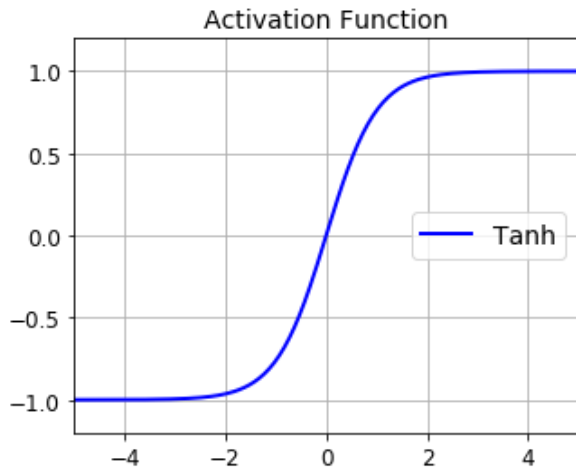


34

**Sigmoid Function Derivative:**

- $\sigma(z) = \text{sigm}(z) = \frac{1}{1+\exp(-z)}$

- $\sigma'(z) = \sigma(z)\,(1 - \sigma(z))$



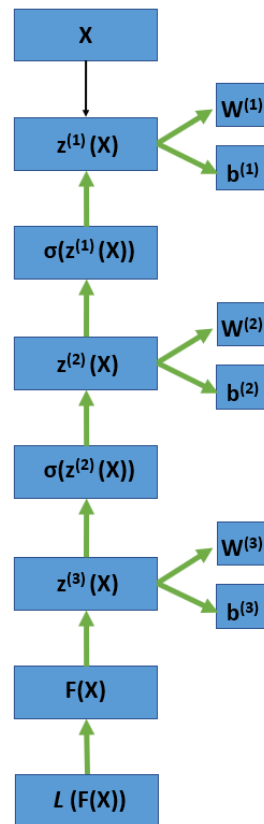Activation Function

**Hyperbolic tangent ("tanh") Derivative:**

- $\sigma(z) = \tanh(z) = \dfrac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} = \dfrac{\exp(2z) - 1}{\exp(2z) + 1}$

- $\sigma'(z) = 1 - \sigma(z)^2$



Activation Function

Gradient of Weights

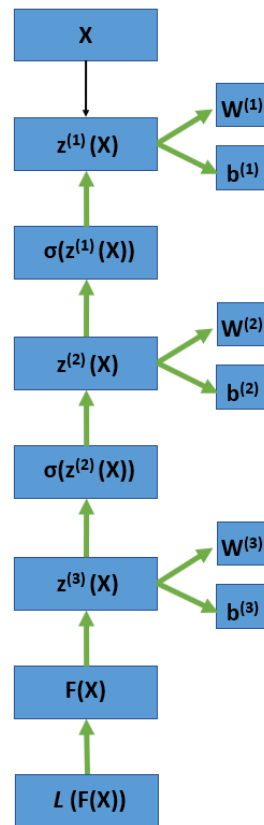- $\nabla_{\mathbf{w}^{(k)}} - log\mathbf{F}(\mathbf{x})_y$

$$= (\nabla_{\mathbf{a}^{(k)}(\mathbf{x})} - log\mathbf{F}(\mathbf{x})_y)\, \mathbf{h}^{(k-1)}(\mathbf{x})^\mathsf{T}$$

Gradient of Biases

- $\nabla_{\mathbf{b}^{(k)}} - log\mathbf{F}(\mathbf{x})_y$

$= (\nabla_{\mathbf{a}^{(k)}(\mathbf{x})} - log\mathbf{F}(\mathbf{x})_y)$

**NYU**

## Books

- **Hands-On Machine Learning with Scikit-Learn and TensorFlow, Aurélien Géron**
- **Pattern recognition and machine learning, Christopher M. Bishop**
- **The Elements of Statistical Learning, Hastie, Tibshirani, and Friedman**

## Blogs, code snippets, lecture notes, etc.

- **https://github.com/stephencwelch/Neural-Networks-Demystified This tutorial uses numpy and python**
- **http://www.wildml.com/2015/09/implementing-a-neural-network-from-scratch This tutorial implements a classifier in python from scratch**
- **http://deeplearning.net/tutorial This set of tutorials uses the Theano package, which is pretty tricky to learn (but very powerful because it calculates gradients for you automatically, among other things)**
- **http://www.cs.stir.ac.uk/courses/ITNP4B/lectures/kms/1-Intro.pdf For some of the history and possible connections to neuroscience**
- **http://www.dmi.usherb.ca/~larocheh/index_en.html Modified from Hugo Larochelle**