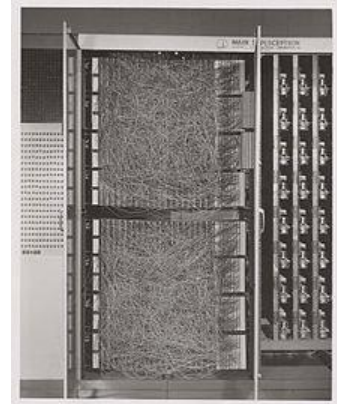# Introduction to Neural Networks

Perceptron, XOR Problem, Multi-layer Perceptron (MLP), Cost Functions, Activation functions, Output units
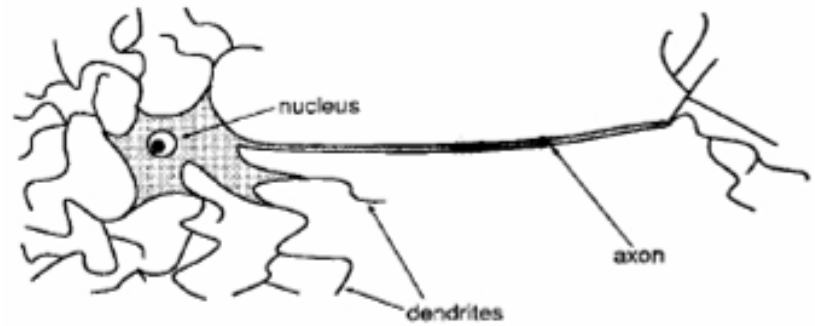
EE 5179
Instructor: Kaushik Mitra

# History

- 1932 - *The Integrative Action of the Nervous System*, **Sir Charles Scott Sherrington**
  - Nervous system - interconnection of individual entities (neuron)
- 1943 - *A Logical Calculus of Ideas Immanent in Nervous Activity*, **Warren McCulloch** and **Walter Pitts**
  - McCulloch and Pitt's model
- 1949 - *The Organization of Behavior,* **Donald Hebb**
  - Hebbian learning
- 1953 - *The Perceptron*, **Rosenblatt**
- 1969 - *Limitation of Perceptrons,* **Minsky** and **Papert**
- 1980's - *Connectionism*
  - *Error Back Propagation,* **Proposed simultaneously by Many**
- 20th century Deep learning
  - CNNs - LeNet, AlexNet, VGGNet, ResNet
  - Deep LSTMs,
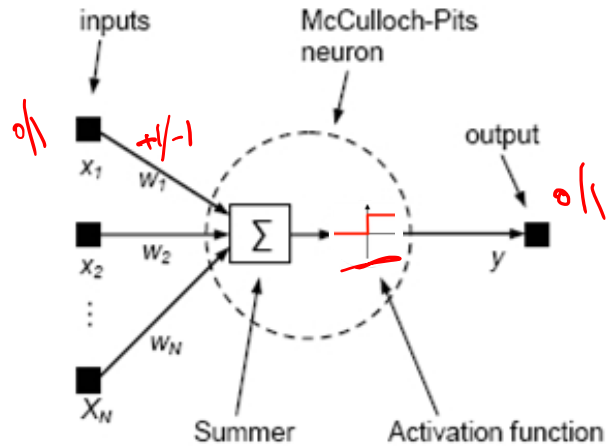- The **deep saga………** what followed is discussed in the last class



The Mark 1 Perceptron
By Rosenblat

# McCulloch - Pits model



Biological neuron

inputs

McCulloch-Pits neuron

output

$x_1$  $w_1$  0/1  +1/-1

$x_2$  $w_2$

$w_N$

$X_N$

Summer

Activation function

$y$  0/1

$$\sum_{i=1}^{n} w_i x_i > \mu$$

$w_1 x_1 + w_2 x_2 + w_N x_N > \mu$

$$\Rightarrow y = 1$$

$w_1 x_1 + \cdots$

$$< \mu$$

$$\Rightarrow y = 0$$

$$y = \textcircled{f}\left(\sum w_i x_i - \mu\right)$$

$> 0$

The input and output are binary and the weights are either excitatory (+1) or inhibitory (-1).

# The Perceptron - Rosenblatt (1953)



inputs

McCulloch-Pits neuron

output

$x_1$

W

$w_2$

$x_2$

$w_N$

$X_N$

$\Sigma$

$y$

Summer

Activation function

$$\sum_{i=1}^{n} w_i x_i > \mu$$



*Pic courtesy, wikipedia

The Mark 1 Perceptron
By Rosenblat
for digit recognition

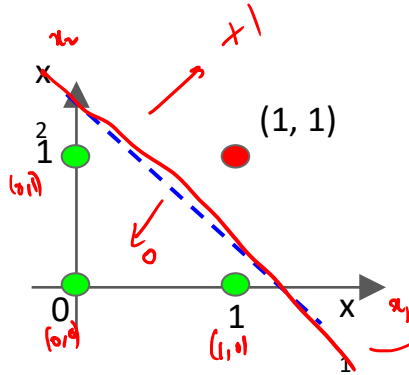Weights can be any real number. He also proposed a learning algorithm for learning the weights from training data.

# Perceptron - geometrical interpretation

$$\sum_{i=1}^{n} w_i x_i > \mu$$

, What does this inequality imply in 2D case?   Half plane

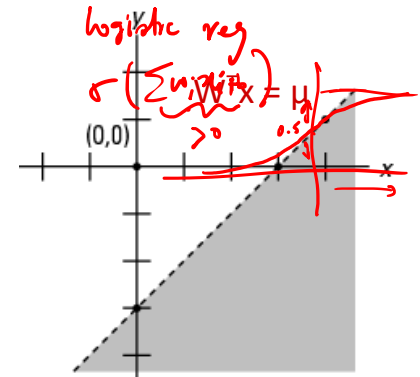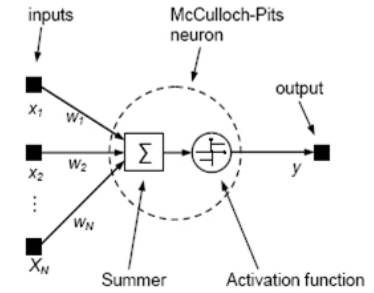| X | AND |
|---|---|
| (0, 0) | 0 |
| (0, 1) | 0 |
| (1, 0) | 0 |
| (1, 1) | 1 |

Solve for W, μ:

$x_1 + x_2 > 1.5$

$w_1 = 1$, $w_2 = 1$ and $\mu = 1.5$

Any function that is linearly separable can be computed by a perceptron

*Pic courtesy, cliffsnotes

# Perceptron - Limitations

**Goal**: learn the XoR function ($f*$)

| $\mathbb{X}$ | $f*$ |
|---|---|
| (0, 0) | 0 |
| (0, 1) | 1 |
| (1, 0) | 1 |
| (1, 1) | 0 |

Original $x$ space

How to tackle this problem?
- Can we use more than one line?
- Yes, but how?

Task is adjust parameters $\vartheta$, such that $f$ is as close as to $f*$

$$y = f(x, \vartheta) \qquad L = \sum_{\{x \in \mathbb{X}\}} \left( f*(x) - f(x, \vartheta) \right)^2$$

Lets use our perceptron for $f$, $\vartheta = \{w, b\}$

$$f(x; w, b) = w^T x + b$$

Solve for $\{w, b\}$

W = 0, b = 0.5; output is 0.5 everywhere

**Why** this linear function can't model XoR**?**

# Perceptron - Limitations

*handwritten:*
$h = A x$
$y = B h$
$y = B A x$
$y = C x$

**Goal**: learn the XoR function ($f^*$)


Original $x$ space

**How** to tackle this problem?

- Add a hidden layer with two units

*handwritten:* $y_{1 \times 1} = f(U h_{2 \times 1} + c_{1 \times 1})$

$y = f^{\cdot(2)}(h; U, c)$

$y = f^{\cdot(2)}( f^{\cdot(1)}(x) )$

$h = f^{\cdot(1)}(x; W, b)$

*handwritten:* $\begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = f\left( W x_{2 \times 1} + b_{2 \times 1} \right)$

What should $f^{\cdot(1)}$ compute?

*handwritten:*
$h_1 = f(W_{11} x_1 + W_{12} x_2 + b_1)$
$h_2 = f(W_{21} x_1 + W_{22} x_2 + b_2)$
$h = f(W$

If its linear again the composition still remains linear

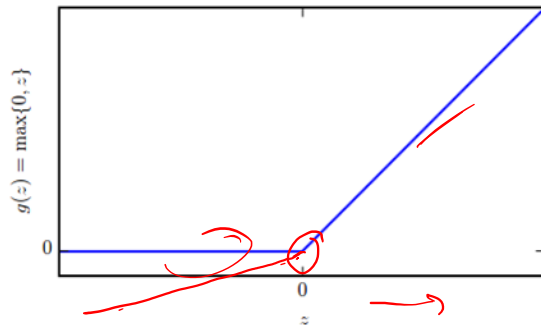$f^{\cdot(2)}(h) = U^T h$ ; since $h = Wx$

$y = U^T W x = W'x$



- $f^{\cdot(1)}$ should be nonlinear to extract useful features

$h = f^{\cdot(1)}(x; W, b) = g(Wx+b)$

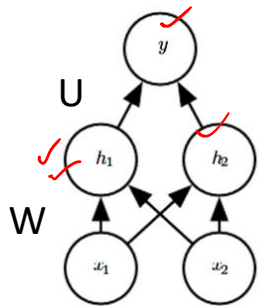- $g$ is referred as **activation** function commonly

- We will use ReLU here
  - Rectified Linear Unit (widely used)
  - $g(z) = \max\{0,z\}$



*Slide courtesy, Ian Goodfellow et al., deep learning book

# Perceptron - Limitations

**How** to tackle this problem?

– Add a hidden layer with two units
– Use ReLU activation in 1st layer

**Goal**: learn the XoR function ($f*$)



$y = U^T h + c$ ; $y = U^T \max\{0, Wx+b\} + c$

ReLU

$h = g(Wx+b)$

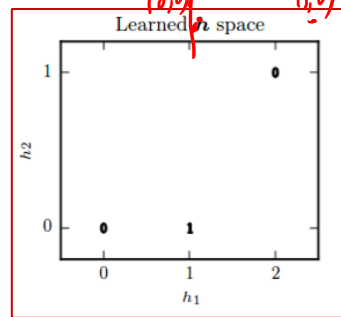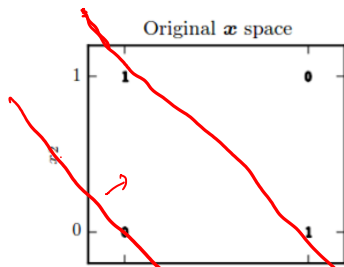$h_1 = {}^{ReLU}(W_1 x_1 + W_{12} x_2 + b_1)$
$h_1 = (x_1 + x_2 + 0)$
$h_2 = {}^{ReLU}(x_1 + x_2 - 1)$

Let,

$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$, $b = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$, $U = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$,

$c = 0$

$y = (h_1 - 2h_2)$
$h_1 - 2h_2 = 0$  $h_2 = \frac{1}{2} h_1$

$X = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$    $WX = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \end{bmatrix}$

$WX + b = \begin{bmatrix} 0 & 1 & 1 & 2 \\ -1 & 0 & 0 & 1 \end{bmatrix}$    $h = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ **Upon ReLU**
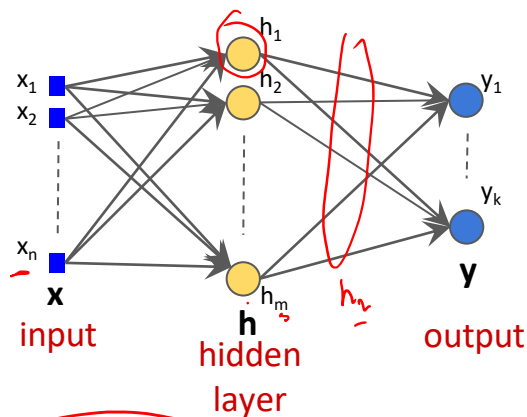
$U^T h + c$

After layer2

$[0 \ 1 \ 1 \ 0]$

| $(x)$ | $h_1$ | $h_2$ | $y$ |
|---|---|---|---|
| $(0,0) \rightarrow$ | 0 | 0 $\rightarrow$ | 0 |
| $(1,0) \rightarrow$ | 1 | 0 $\rightarrow$ | 1 |
| $(0,1) \rightarrow$ | 1 | 0 $\rightarrow$ | 1 |
| $(1,1) \rightarrow$ | 2 | 1 $\rightarrow$ | 0 |

*Slide courtesy, Ian Goodfellow et al., deep learning book

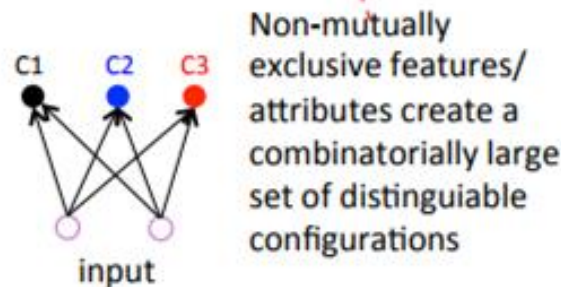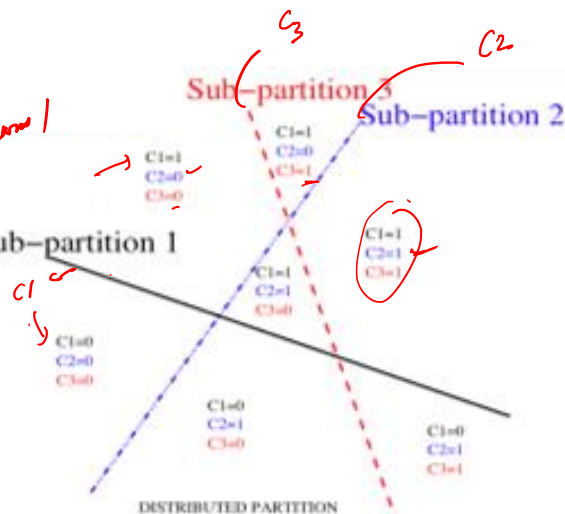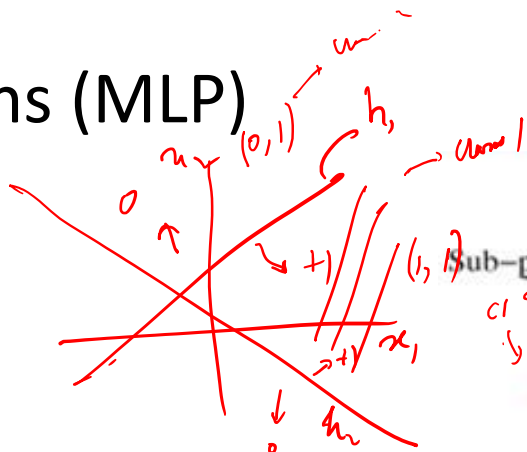# Multi-layer Perceptrons (MLP)

A typical feed forward neural network



$\mathbf{h} = f(\mathrm{W}\mathbf{x} + \mathbf{b_1}); \quad \mathbf{y} = g(\mathrm{U}\mathbf{h} + \mathbf{b_2})$

With more hidden units network is more expressible



DISTRIBUTED PARTITION

**Non-mutually exclusive features/ attributes create a combinatorially large set of distinguiable configurations**

# Feedforward Neural Networks - Cost functions

$$x \rightarrow y = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ N \end{pmatrix} \begin{smallmatrix} class \, 1 \\ \\ 2 \end{smallmatrix}$$

**For regression,**

$$J(\theta) = \frac{1}{2}\mathbb{E}_{\mathbf{x},\mathbf{y}\sim\hat{p}_{\text{data}}}\|\boldsymbol{y} - f(\boldsymbol{x};\boldsymbol{\theta})\|^2$$

$$\frac{1}{2}\sum_{\{x_i,y_i\}}\|y_i - f(x_i,\theta)\|^2$$

training data

$(x_i, y_i)$

$\min J(\theta)$

$\theta = \{w, b\}$

Cross entropy

entropy

$= p(x)$

$\sum_x p(x) \ln \dfrac{1}{q(x)}$

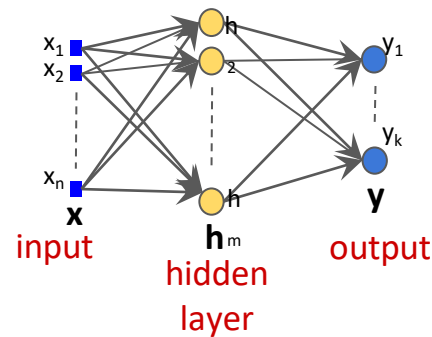ent $(p) = \sum_x p(x)\left(\ln \dfrac{1}{p(x)}\right)$

**For classification,**

$q(x) = f(x, \theta)$

- Typically outputs a probability vector $q(c = k|x) \ \forall \ k$

- How do you compare two distributions?
  - KL divergence, KL($p\|q$)

$$D_{KL}(p(x)\|q(x)) = \sum_{x\in X} p(x)\ln\frac{p(x)}{q(x)} = \sum p(x)\ln\frac{1}{q(x)} - p(x)\ln\frac{1}{p(x)}$$

Cross-entropy

$$= \sum p(x)\ln p(x) - p(x)\ln q(x)$$

$H(p)$

$$= -H(p) + H(p,q)$$

$w, b$

$\underbrace{\text{entropy}}\underbrace{\text{cross-entropy}}$

$a \rightarrow 0.5 \rightarrow p(a) \quad q(a)$ of y

$b \rightarrow 0.2 \rightarrow p(b)$

$c \rightarrow 0.2 \rightarrow p(c)$

$d \rightarrow 0.1 \rightarrow p(d)$

$$\min_{\theta=\{w,b\}} J(\theta) = \sum_{x_i,y_i} H(p(x_i), q(x_i))$$

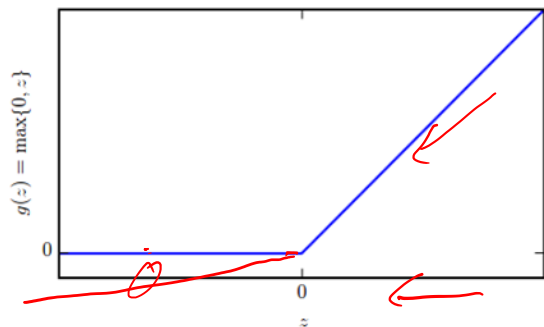# Activation functions
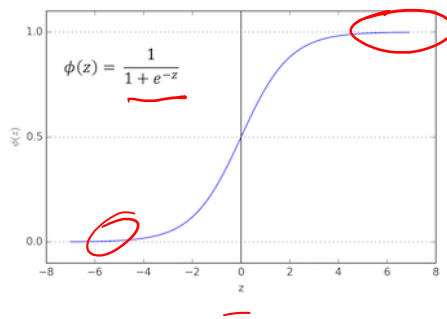


$h = g(Wx+b)$;    Affine transformation followed by activation function, $g$
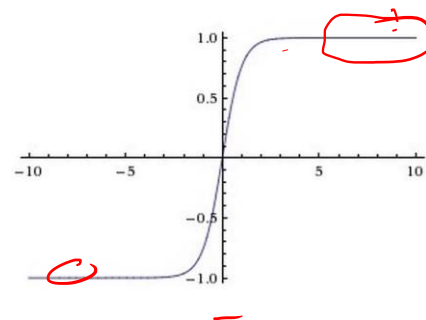
– Very important factor in learning features
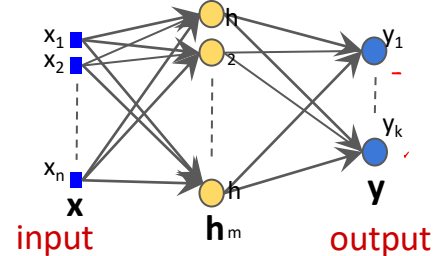
$g(z) = \max\{0,z\}$, ReLU          $g(z) = \sigma(z)$, sigmoid          $tanh(z) = 2\sigma(2z) - 1$

# Output units


input     $\mathbf{h}_m$     output

– Linear units for real valued outputs
  ❑ Activation function is left to be linear
  ❑ Given features h,

$$y' = Wh+b$$

  ❑ Most commonly used with regression tasks

– Say you want to do binary classification
  ❑ What kind of distribution describes output?
  **Bernouli**
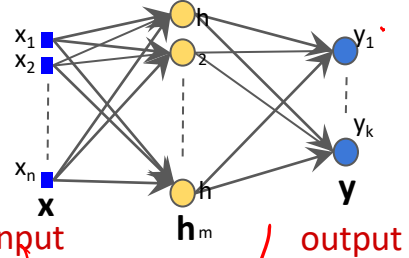  ❑ How to constrain the output - valid probability? Can you use linear activation?

$$P(y = 1 \mid \boldsymbol{x}) = \max\left\{0, \min\left\{1, \boldsymbol{w}^\top \boldsymbol{h} + b\right\}\right\}. \quad \in [0, 1]$$

  ❑ What is the problem?    *Not amenable for gradient based learning*

  ❑ Instead, use sigmoid unit - output ∈ [0,1]

$$\hat{y} = \sigma\left(\boldsymbol{w}^\top \boldsymbol{h} + b\right)$$

# Output units



input  $\mathbf{h}_m$  output

- Now, say we want to do multi-class classification (K classes)
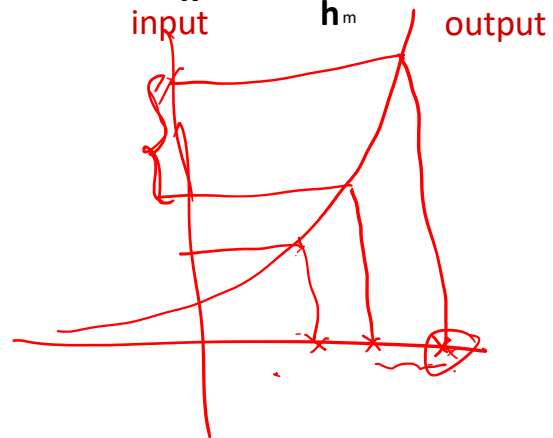  - ❏ Output should be K probabilities,
    $p_k = p(class = k | x) \; \forall \; k = 1$ to K

  - ❏ Can we use K sigmoid units?
    Won't be sufficient, since probabilities are not constrained to sum to 1
    $$\sum_k p_k = 1$$

  - ❏ We will look at softmax unit for this
    Idea is to convert a vector of real values to valid probabilities,
    How? Make all the elements positive
    - ❏ Normalize the values

- Let, $\mathbf{z} = [z_1, \ldots, z_K]^T$;     $\mathbf{z} = W\mathbf{h} + \mathbf{b}$

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}.$$