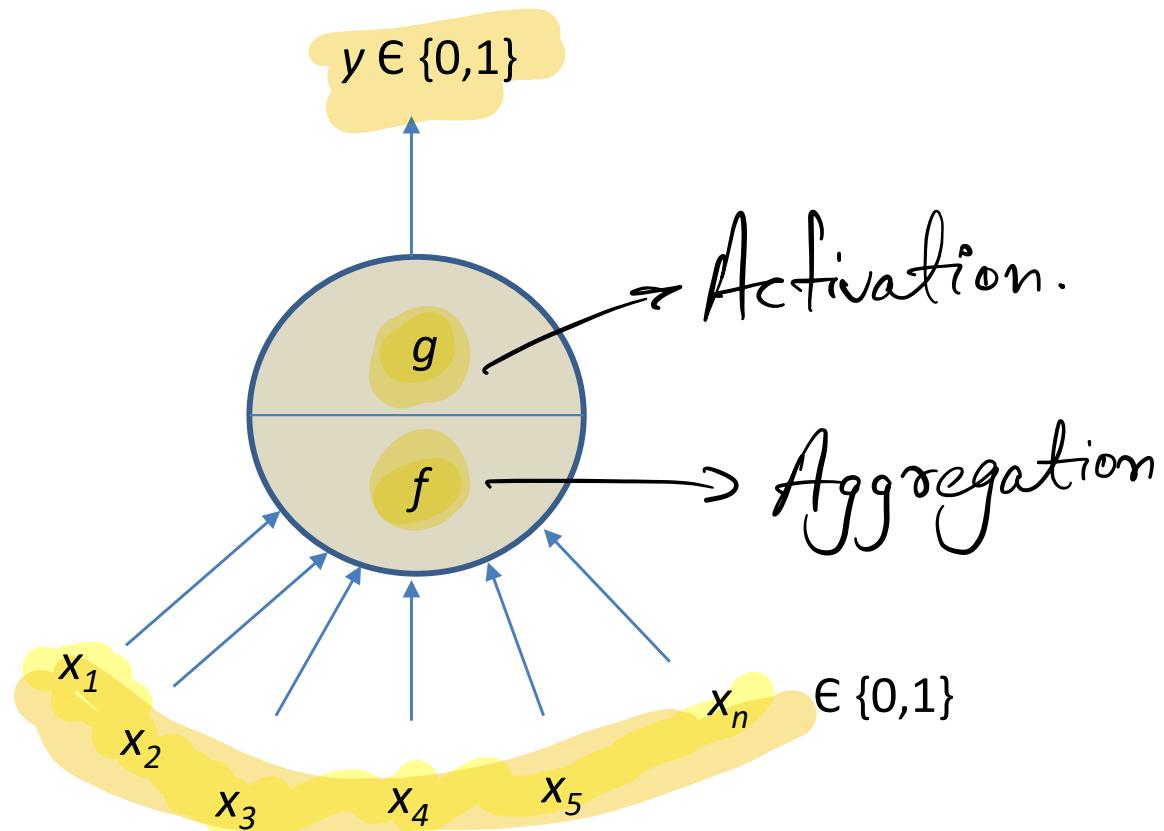


Deep Learning

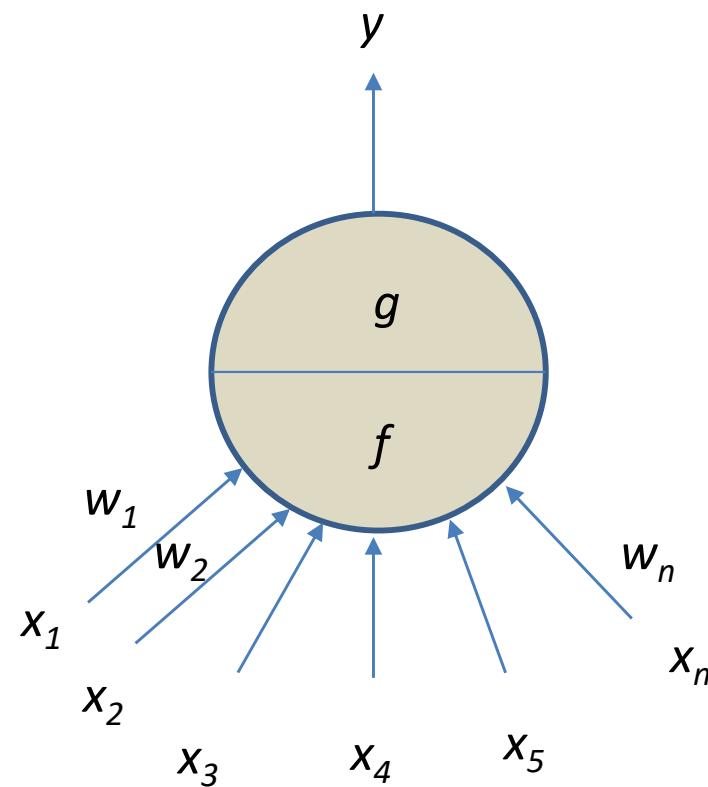
McCulloch Pitts (MP) Neuron

- 1943: McCulloch, neuroscientist, and Pitts, logician, proposed a simplified model of neuron



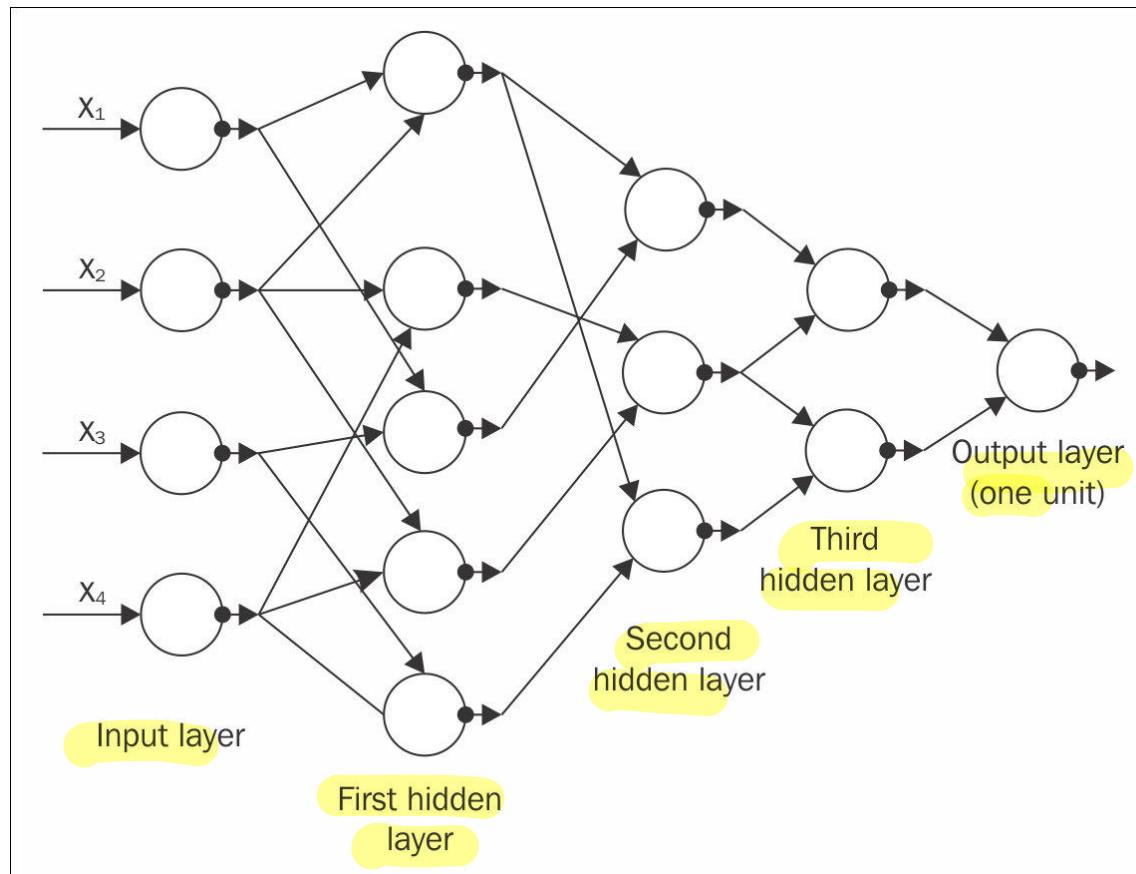
Perceptron

- 1957-58: Frank Rosenblatt, “....perceptron may eventually be able to learn, make decisions, and translate languages..”



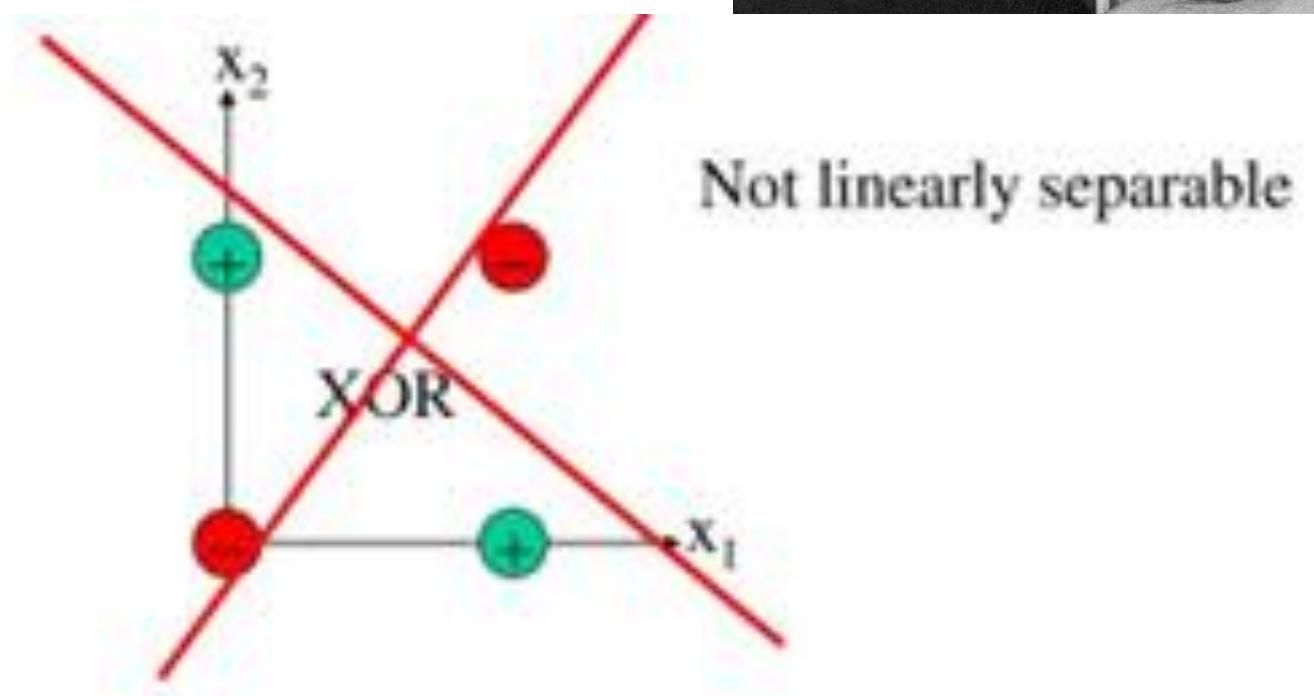
1st Gen Multi-Layer Perceptron

- 1965-68: Ivakhnenko et al.



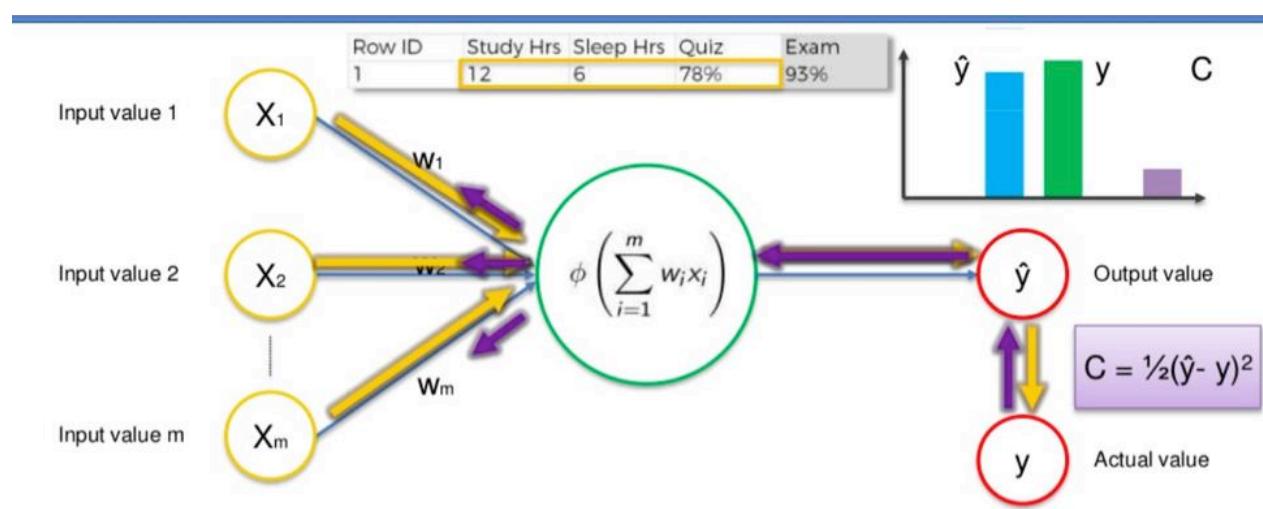
Perceptron Limitations

- 1969: “Perceptrons” by Minsky and Papert



Backpropagation

- 1960-70s: discovered/rediscovered
- 1982: Werbos used it in context of ANN
- 1986: Rumelhart et al. popularized it



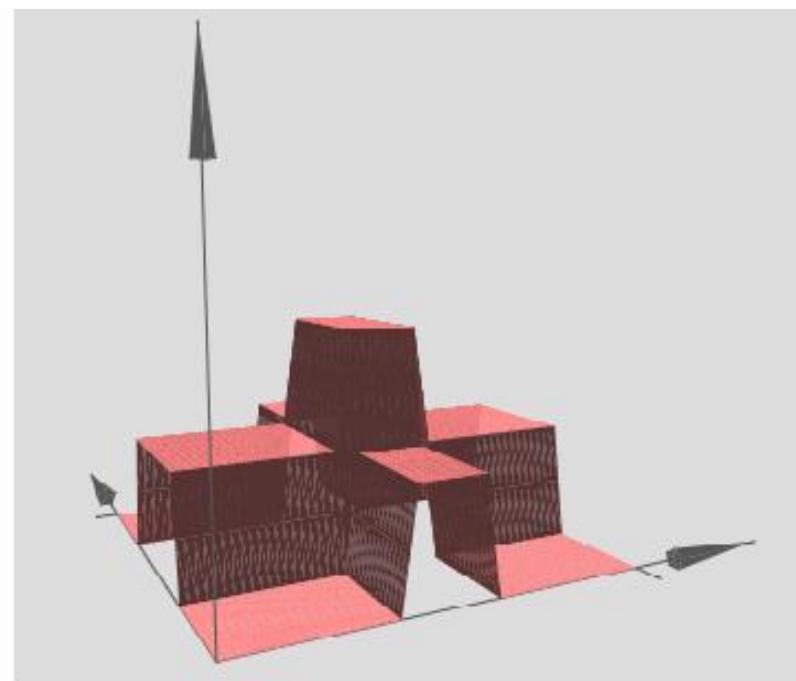
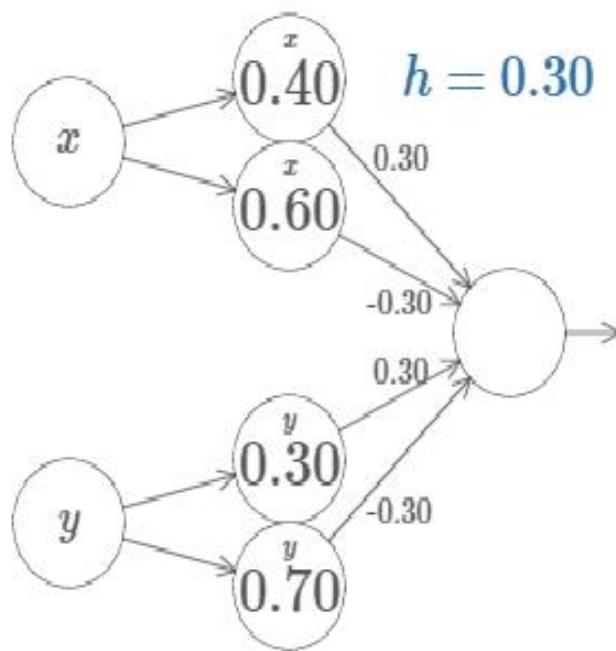
Gradient Descent

- Discovered by Cauchy in 1847, motivated by need to compute orbit of heavenly bodies
 - Used for backpropagation in neural networks



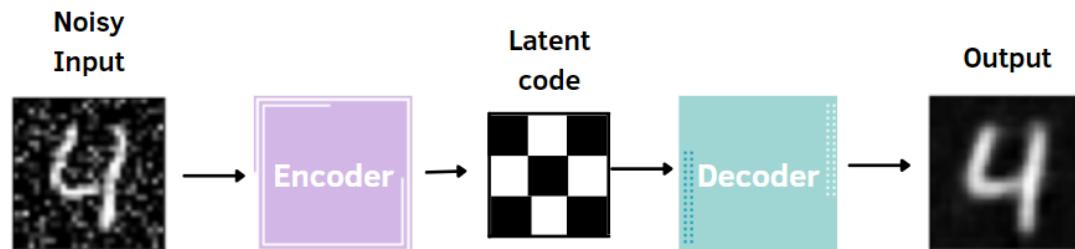
Universal Approximation Theorem

- 1989-90: A multi-layered network of neurons with a single hidden layer can be used to approximate any continuous function to any desired precision



Unsupervised Pre-training

- 2006: Hinton and Salakhutdinov propose method of initializing weights that allows autoencoders to learn low-dimensional representation of data
- 2006-2009: Further explorations

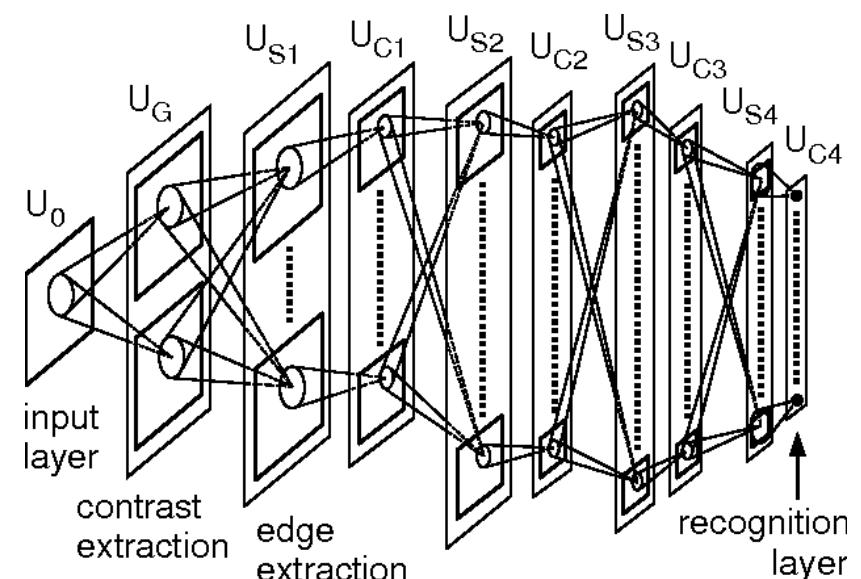
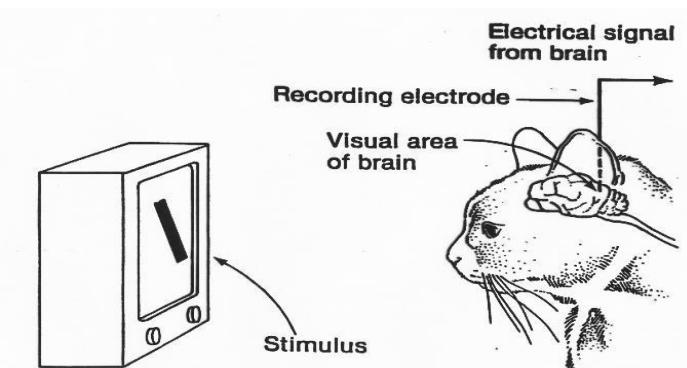


Success Stories

- 2009: Handwriting recognition (Graves et al.)
- 2010: Speech recognition (Dahl et al.)
- 2010: Handwritten digits recognition (MNIST), Ciresan et al.
- 2011: Visual Pattern Recognition (Ciresan et al., IJCNN Traffic Sign Recognition Competition)
- 2012-2016: Visual Recognition Challenges (ImageNet)

Cats to CNN

- 1959: Hubel and Wiesel experiment
 - A neuron fires only in response to a visual stimuli in a specific region in visual space
- 1980: Neocognitron
 - Used for handwritten character and pattern recognition (Fukushima et al.)



CNN

- 1989: handwritten digit recognition using back propagation over CNN (LeCun et al.)
- LeNet-5 introduced MNIST

80322-4129 80306
40004 14310
37878 05153
5502 75216
35460 44209

1611913485726803226414186
6359720299299722510046701
3084111591010615406103631
1064111030475262009979966
8912056708557131427955460
1018730189112993089970984
0109707597331972015519056
1075518255182814358090943
1787541655460554603546055
18255108503047520439401

Figure 1 Examples of original zip codes (top) and normalized digits from the testing set (bottom).

Sequences

- Time series, speech, music, text, video
 - Each unit in sequence interacts with others
 - Models needed to capture interaction
- 1982: Hopfield – Content addressable memory systems for storing and retrieving patterns
- 1986: Jordan network – Output of each time step fed to next time step allowing interactions between time steps
- 1990: Elman network – Hidden state of each time step fed to next time step allowing interactions

Sequences

- 1991-1994: Drawbacks of training RNNs (Bengio et al.)
- 1997: LSTMs – can solve long time lag tasks
- 2014: Attention mechanism – sequence to sequence learning
- RL for Attention: Schmidhuber and Huber – RNNs with reinforcement learning

Better Optimization

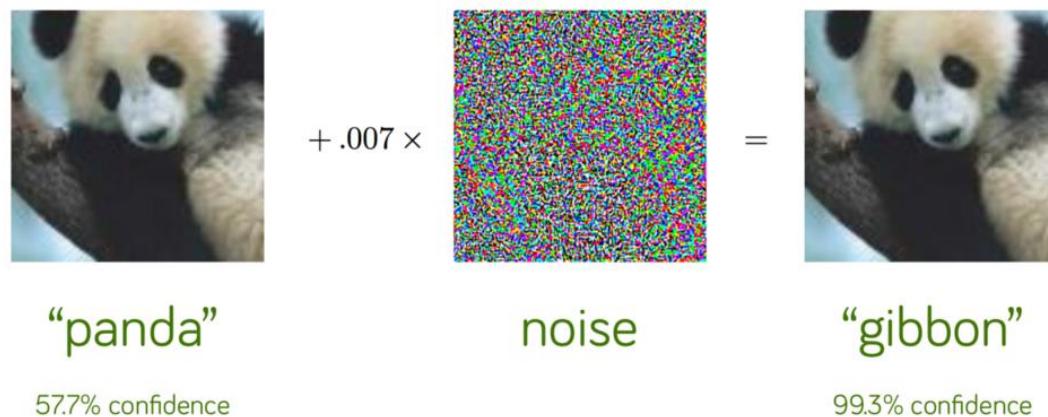
- 2011: Adagrad
- 2012: RMSProp
- 2015: Adam, Batch Normalization
- 2018: Beyond Adam
- Faster convergence, better accuracies

Beating Humans at Playing Games

- 2015: Atari Games
- Alpha Go, GO
- Poker: Deep Stack
- Defense of the Ancients

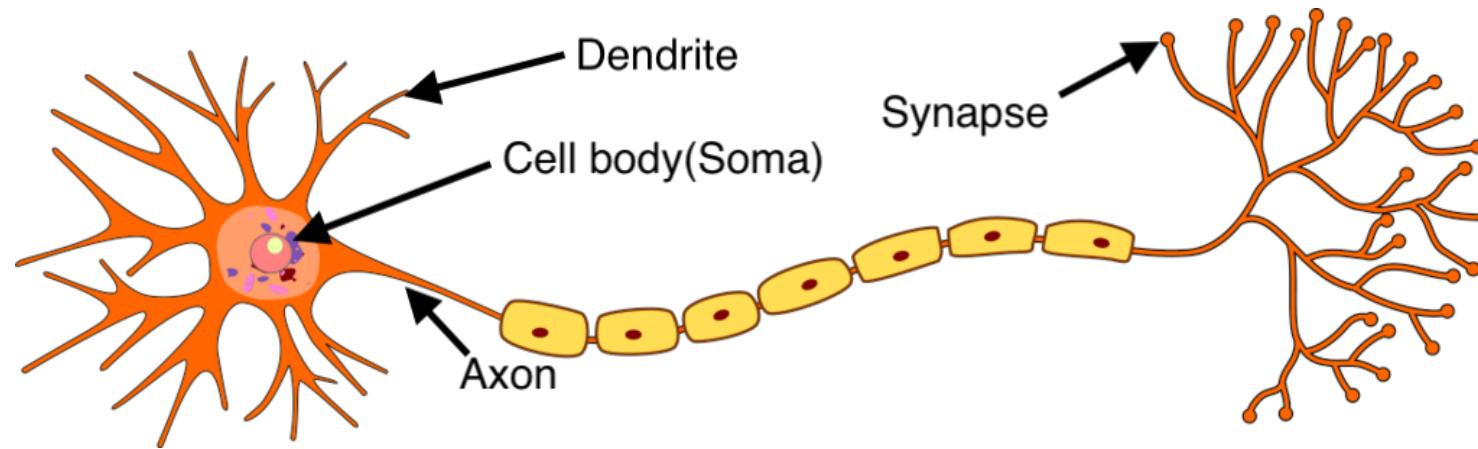
Paradox of Deep Learning

- DL work well despite:
 - High capacity (susceptible to overfitting)
 - Numerical instability (vanishing/exploding gradients)
 - Sharp minima (leads to overfitting)
 - Non-robust

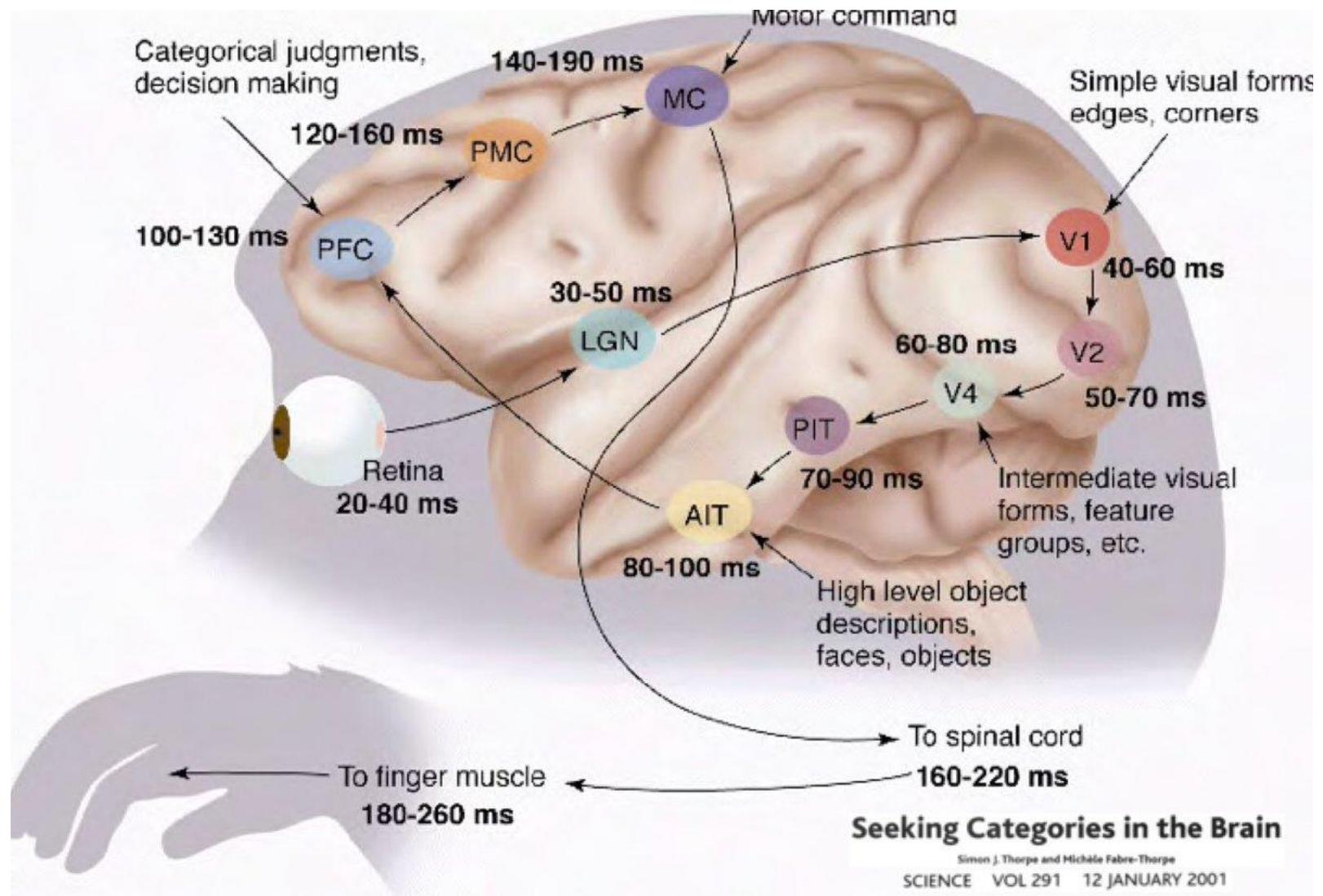


Perceptrons

Biological Neuron

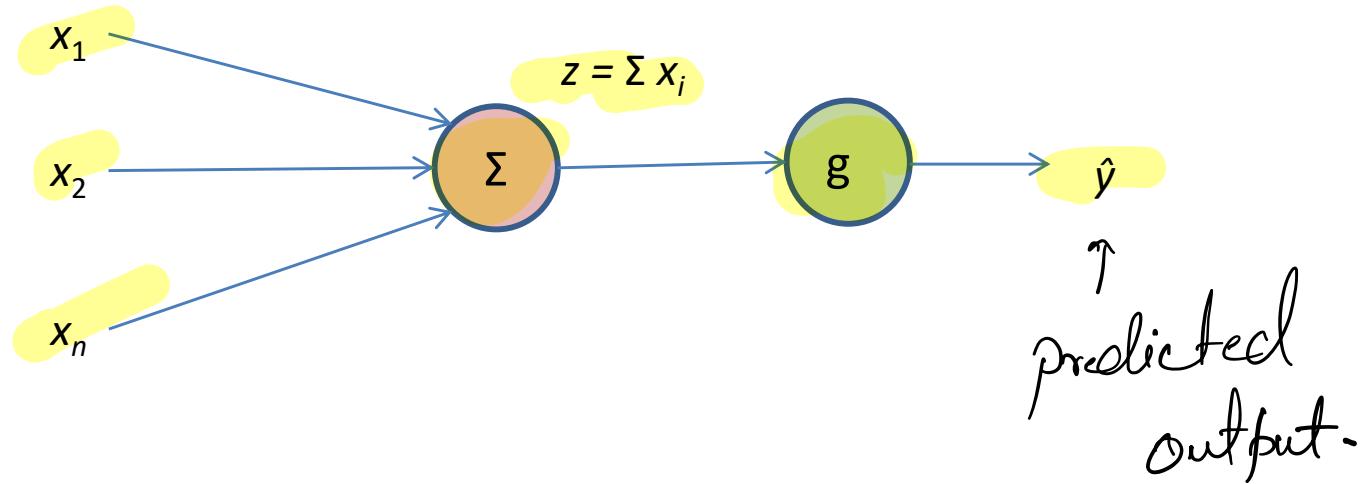
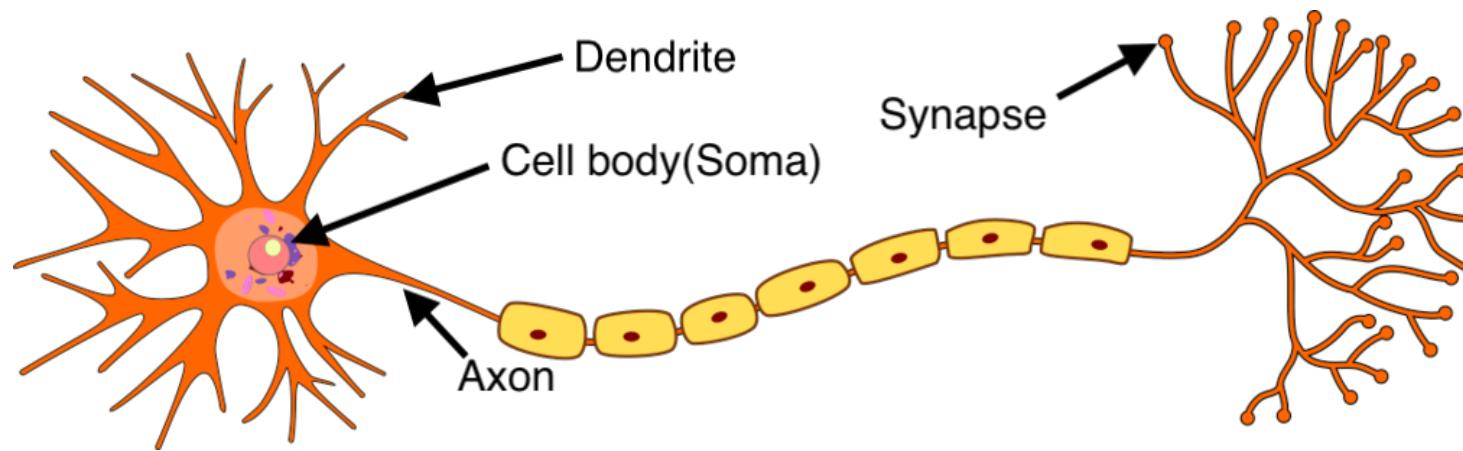


- Massively parallel interconnected network of neurons
- Sense organs relay information to lowest layer of neurons
 - Fired neurons relay information to other connected neurons
 - Division of work – respond to certain stimulus
- Average human brain has around 10^{11} neurons



Picture by Simon Thorpe

Basic Building Block

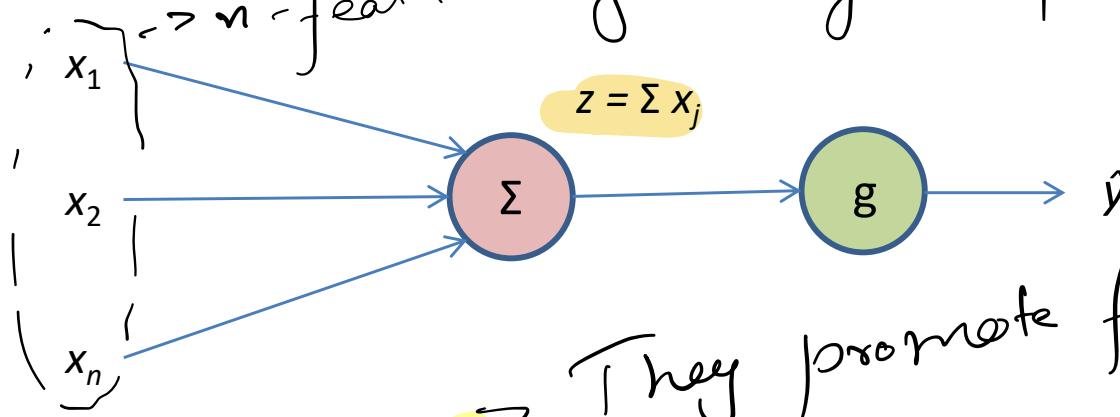


Basic Building Block

(MP Neuron)

Logistic

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix}$$



They promote firing of a neuron

- Inputs can be excitatory or inhibitory → They have an opposite effect.
- Assuming, \mathbf{x} is an n-dimensional vector and $x_i \in \{0,1\}$

$$\hat{y} = g(z) = g(\sum x_i) \quad \forall i = 1 \dots n$$

$$\hat{y} = 1$$

$$\text{if } z \geq \theta$$

$$= 0$$

$$\text{if } z < \theta$$

How to handle?

- What about non-Boolean, real inputs?
- Always need to hard-code threshold?
- Are all inputs equal?? Assign more weights to some inputs?
- Non-linearly separable functions??

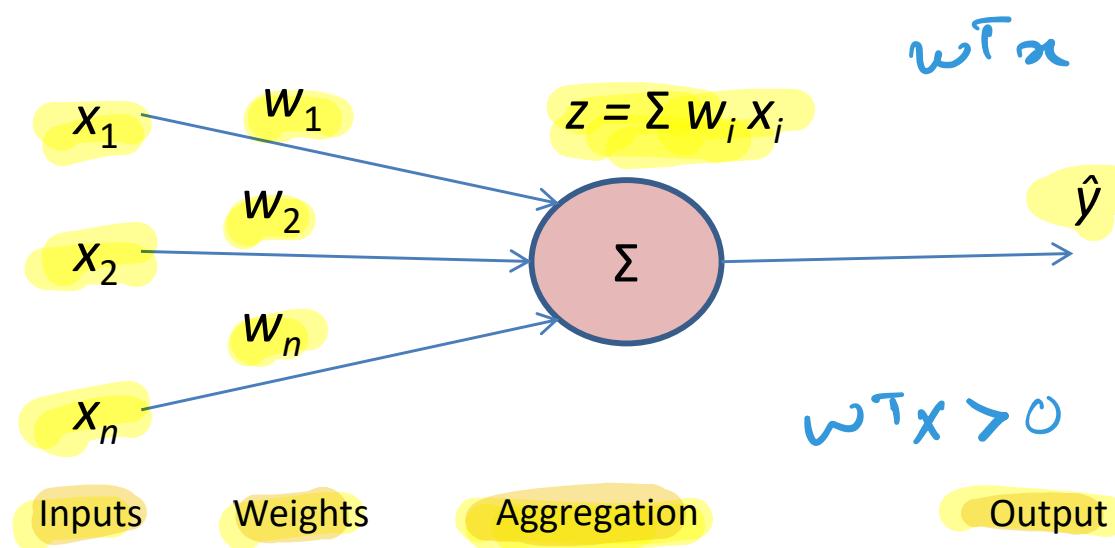
Classical Perceptron

- Frank Rosenblatt proposed classic perceptron (1958)
 - Numerical weights for inputs – mechanism for learning weights
 - Inputs not limited to Boolean values
- Refined by Minsky and Papert – referred to as perceptron model

Classical Perceptron

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_n \end{pmatrix}$$

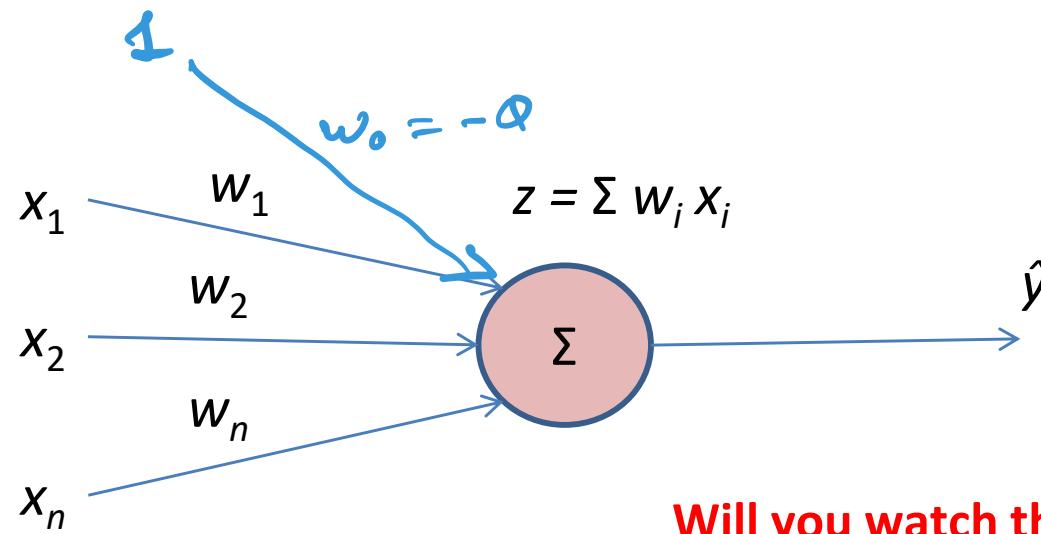


$$\begin{aligned}\hat{y} &= (\sum w_i x_i) \\ &= \mathbf{w}^T \mathbf{x}\end{aligned}$$

$$\forall i = 1 \dots n$$

w_i

Perceptron – Basic Building Block



$$\hat{y} = (\sum w_i x_i)$$

$\forall i = 1 \dots n$

$$\hat{y} = 1 \quad \text{if } \sum w_i x_i \geq \Theta$$

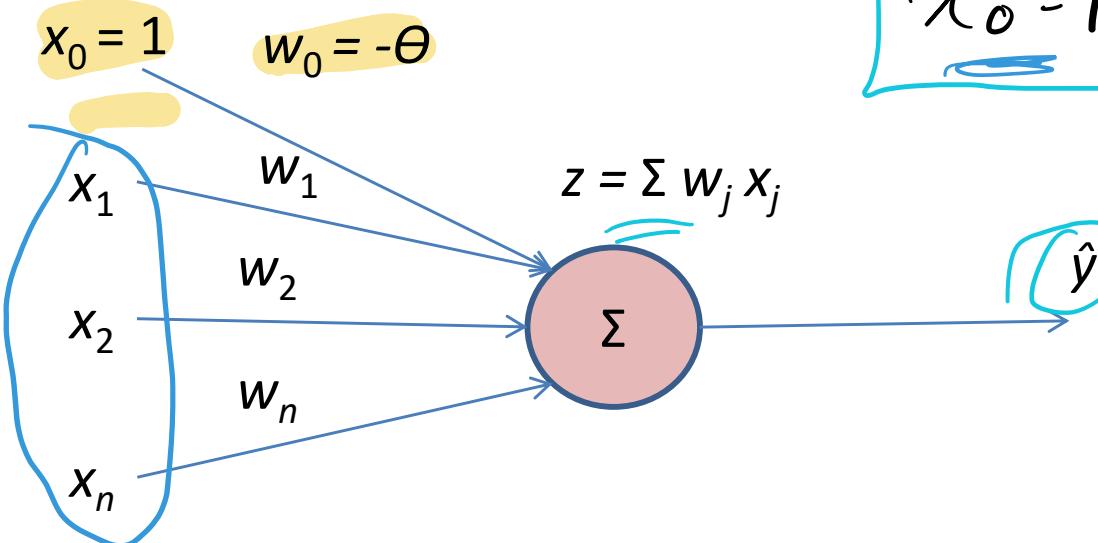
$$= 0 \quad \text{if } \sum w_i x_i < \Theta$$

$$\hat{y} = 1 \quad \text{if } \sum w_i x_i - \Theta \geq 0$$

$$= 0 \quad \text{if } \sum w_i x_i - \Theta < 0$$

Perceptron – Basic Building Block

$$\mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix}$$



$$x_0 = 1, w_0 = -\Theta$$

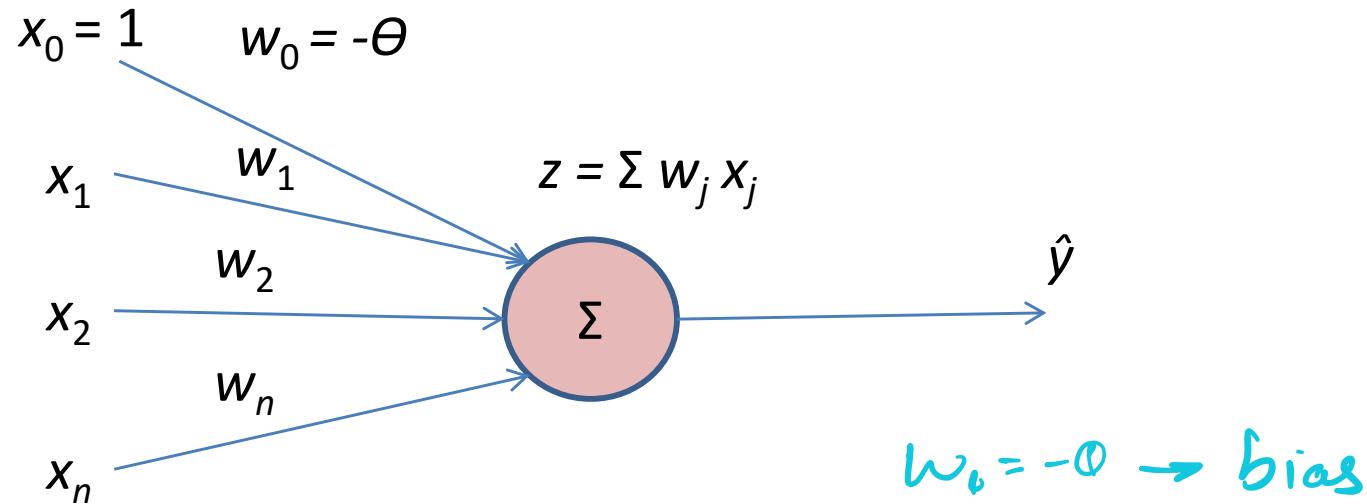
also called
bias-

$$\hat{y} = \begin{cases} 1 & \text{if } \sum w_i x_i - \Theta \geq 0 \\ 0 & \text{if } \sum w_i x_i - \Theta < 0 \end{cases} \quad \forall i = 1 \dots n$$

$$\hat{y} = \begin{cases} 1 & \text{if } \sum w_i x_i \geq 0 \\ 0 & \text{if } \sum w_i x_i < 0 \end{cases} \quad \forall i = 0 \dots n, w_0 = -\Theta, x_0 = 1$$

w_0 / b : bias, represents the prior or prejudice

Perceptron – Basic Building Block



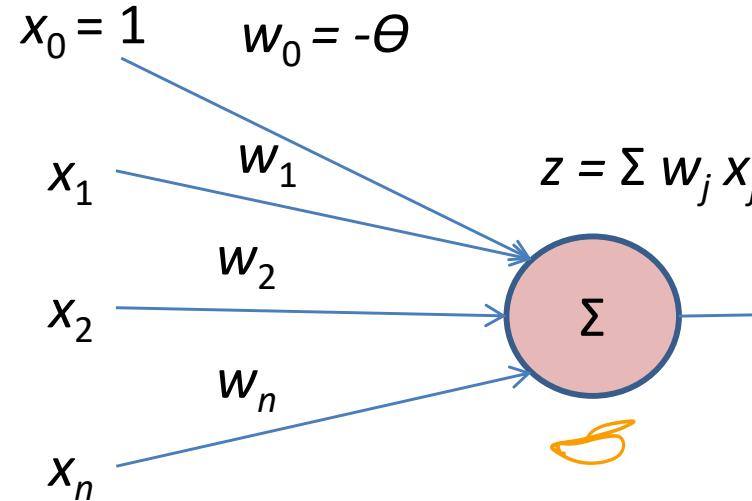
$$\begin{aligned}\hat{y} &= 1 \quad \text{if } \sum w_i x_i \geq 0 \\ &= 0 \quad \text{if } \sum w_i x_i < 0\end{aligned}$$

$$\forall i = 0 \dots n, w_0 = -\Theta, x_0 = 1$$

w_0 : bias

- A movie buff may have a very low threshold and may watch any movie [$\Theta = 0$]
- A selective movie watcher may watch only a thriller, starring Matt Damon and directed by Nolan [$\Theta = 3$]

Perceptron – Basic Building Block

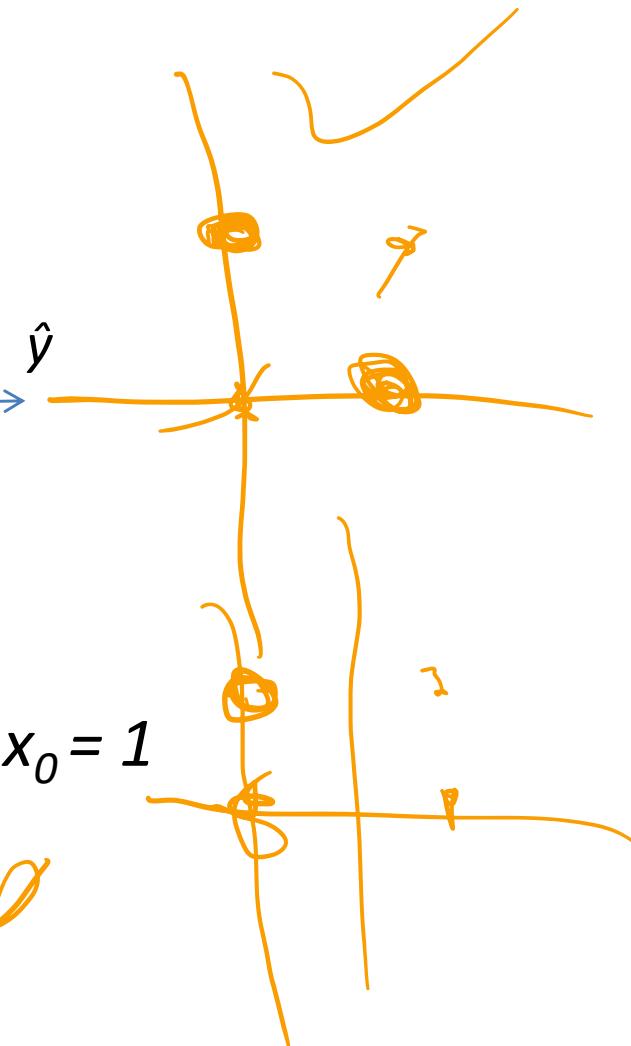


$$\begin{aligned}\hat{y} &= 1 \quad \text{if } \sum w_i x_i \geq 0 \\ &= 0 \quad \text{if } \sum w_i x_i < 0\end{aligned}$$

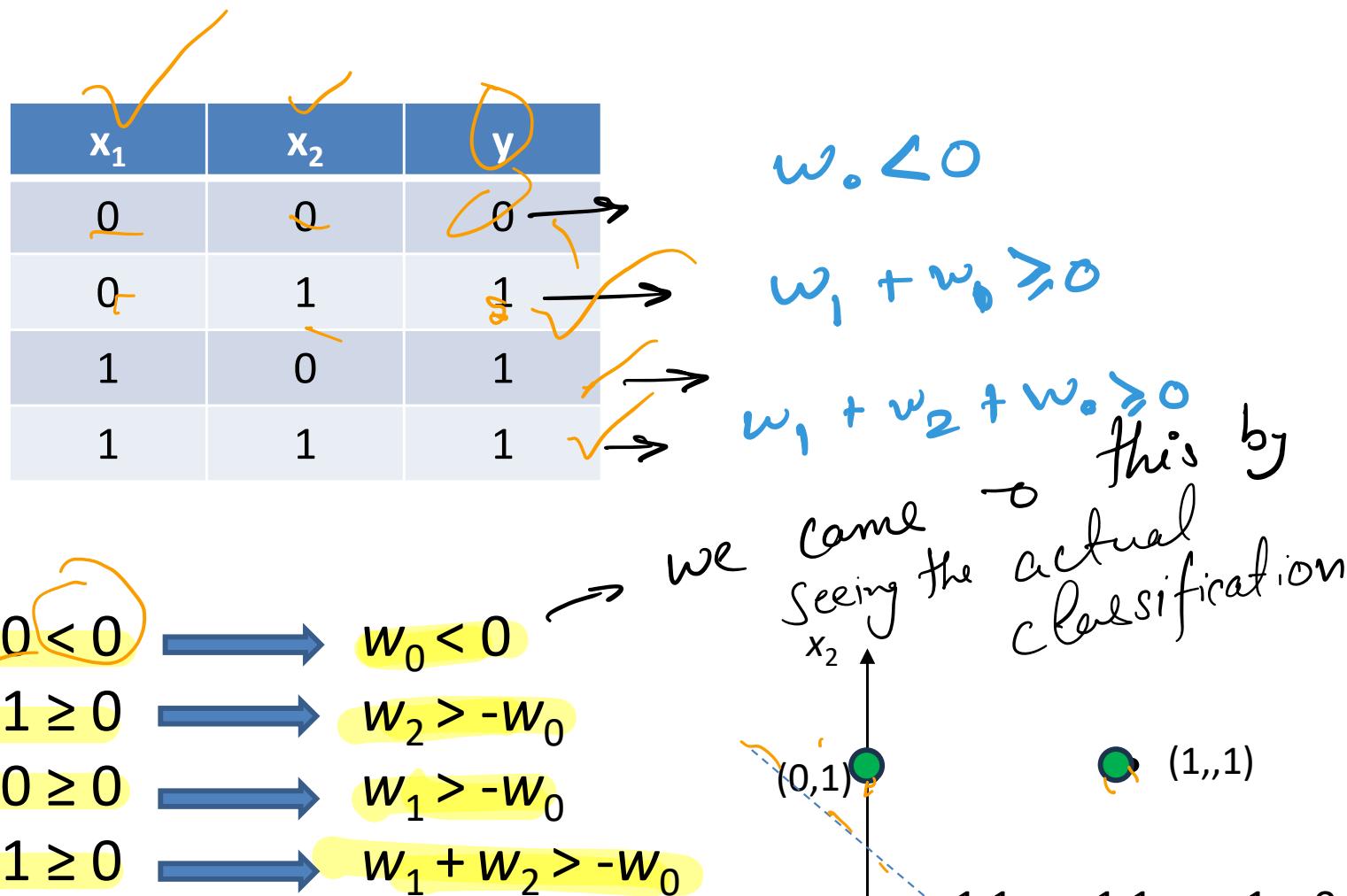
Linearly separable decision boundary: $w^T x = 0$

Perceptron separates input space into two, but:

- Weights (and threshold) can be learned
- Inputs are real-valued



OR function with Perceptron



One possible solution to set of inequalities:

$$w_0 = -1, w_1 = 1.1, w_2 = 1.1$$

$$1.1x_1 + 1.1x_2 - 1 = 0$$

OR function

$$w_0 + w_1 \cdot 0 + w_2 \cdot 0 < 0 \rightarrow w_0 < 0$$

$$w_0 + w_1 \cdot 0 + w_2 \cdot 1 \geq 0 \rightarrow w_2 > -w_0$$

$$w_0 + w_1 \cdot 1 + w_2 \cdot 0 \geq 0 \rightarrow w_1 > -w_0$$

$$w_0 + w_1 \cdot 1 + w_2 \cdot 1 \geq 0 \rightarrow w_1 + w_2 > -w_0$$

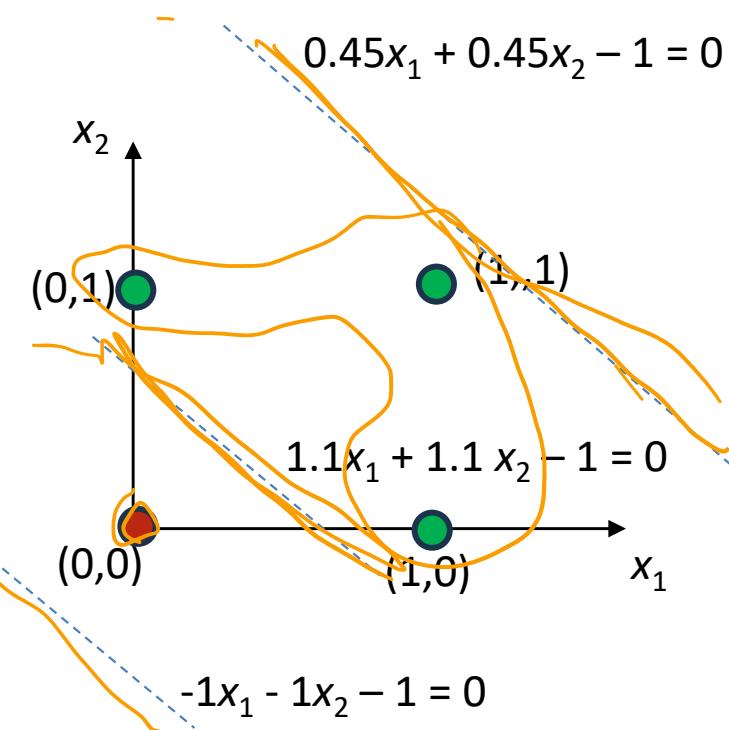
Fix threshold: $w_0 = -1$

Try different values of w_1, w_2

Ex., $w_1 = -1, w_2 = -1 \rightarrow$ one error

Ex., $w_1 = 0.45, w_2 = 0.45 \rightarrow$ 3 errors

Interested in those values of w_0, w_1, w_2 , which gives zero error



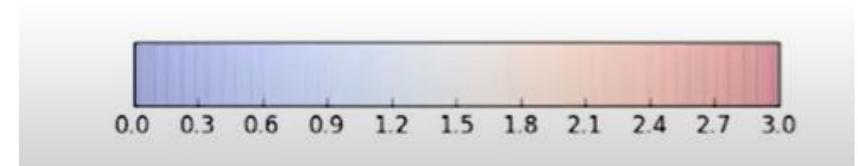
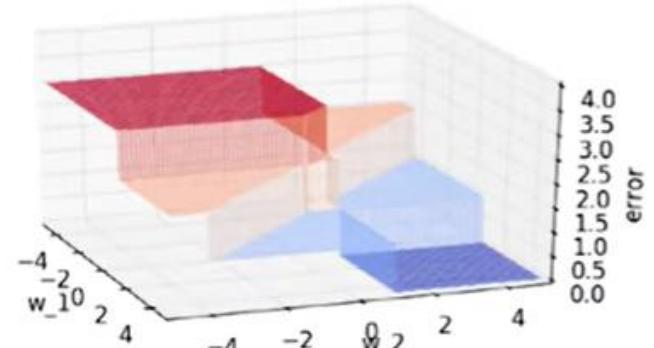
Error Surfaces

Fix threshold: $w_0 = -1$, plot error for different values of w_1, w_2

- For given w_0, w_1, w_2 compute $-w_0 + w_1 x_1 + w_2 x_2$ ($\mathbf{w}^T \mathbf{x}$) for all combinations of (x_1, x_2) and note down errors

For the OR function, error occurs

- if $(x_1, x_2) = (0, 0)$ but $\mathbf{w}^T \mathbf{x} \geq 0$
- or if $(x_1, x_2) \neq (0, 0)$ but $\mathbf{w}^T \mathbf{x} < 0$



Perceptron Learning Algorithm

- How to handle arbitrary functions: $y = f(x)$ instead of Boolean functions

where,

$$x \in \mathbb{R}^n, y \in \mathbb{R}$$

$$[x_1, x_2, x_3, \dots, x_n]^T \rightarrow y$$

Given the n parameters, how much oil will be produced??

- n parameters – some real, some Boolean
- Assuming data is linearly separable, want perceptron to learn how to make decision
- Need a network which can approximately represent these functions

Perceptron Learning Algorithm

$P \leftarrow$ inputs with label 1;

$N \leftarrow$ inputs with label 0;

Initialize $\mathbf{w} = [w_0, w_1, w_2, \dots, w_n]^T$ randomly;

while !convergence **do**

Pick random $\mathbf{x} \in P \cup N$;

if $\mathbf{x} \in P$ and $\sum w_i * x_i < 0$ **then**

$\mathbf{w} = \mathbf{w} + \mathbf{x}$;

→ Adding value to
weights so that $\sum w_i x_i \geq 0$

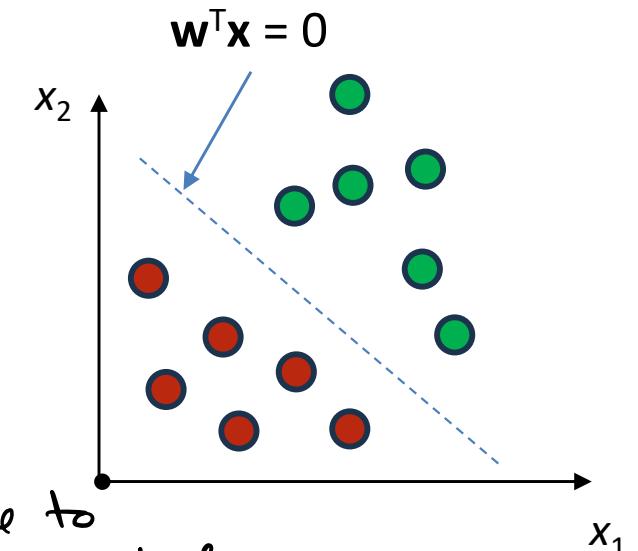
end

if $\mathbf{x} \in N$ and $\sum w_i * x_i \geq 0$ **then**

$\mathbf{w} = \mathbf{w} - \mathbf{x}$;

end

end



Convergence- not making any more errors on training data or predictions are not changing

Perceptron Learning Algorithm

- Consider vectors w and x

$$w = [w_0, w_1, w_2, \dots, w_n]^T$$

$$x = [1, x_1, x_2, \dots, x_n]^T$$

Perceptron rule:

$$\hat{y} = 1$$

$$= 0$$

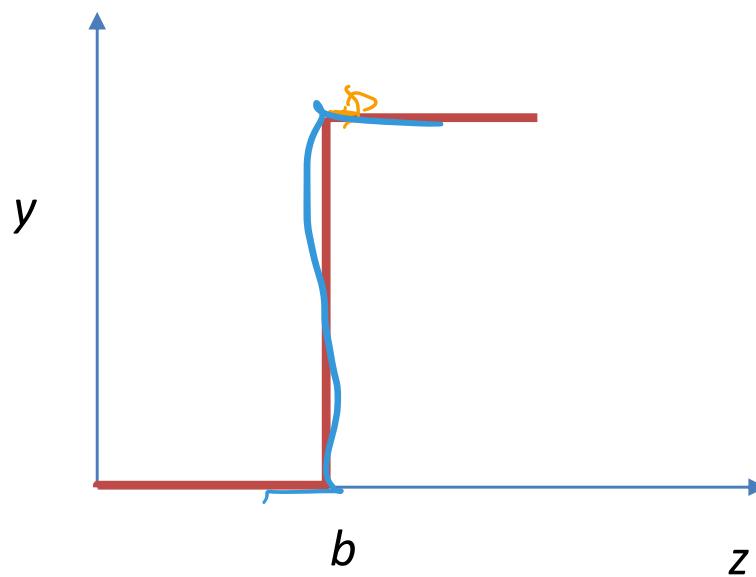
$$\text{if } w^T x \geq 0$$

$$\text{if } w^T x < 0$$

- Find line $w^T x = 0$ which divides input space into two
- Every point x on this line satisfies $w^T x = 0$

Perceptron Logic

- A perceptron will fire if weighted sum of inputs is greater than threshold ($b = -w_0$)



AND function

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

$$\chi = \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}$$

$$x_0 w_0 + x_1 w_1 + x_2 w_2 \\ w_0 < 0$$

$$x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \quad w = \begin{pmatrix} -1.5 \\ 1 \\ 1 \end{pmatrix}$$

if a
line
exists

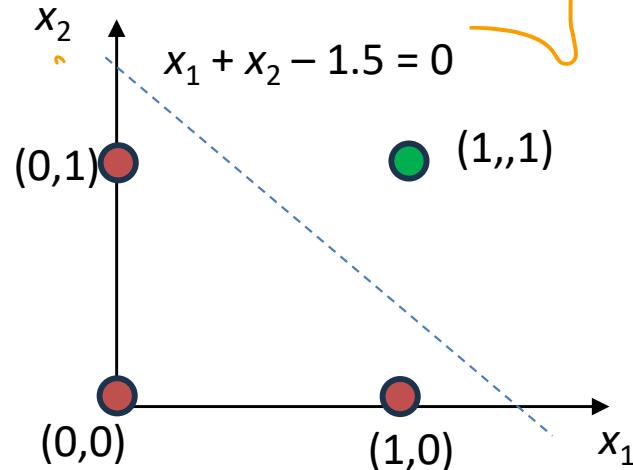
$$X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix} \quad w^T X$$

$w_0 + w_2 < 0$

$w_0 + w_1 < 0$

$w_0 + w_1 + w_2 \geq 0$

Linearly separable



AND function

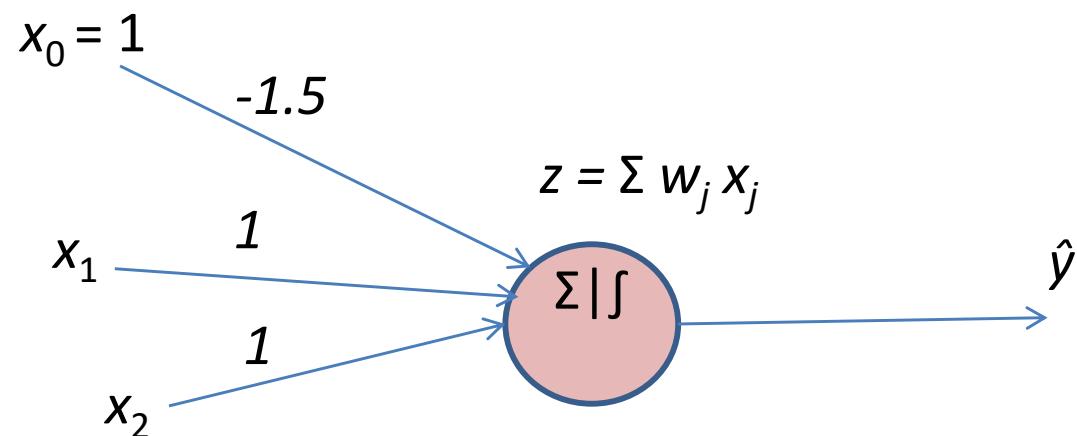
$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix}_{3 \times 4}$$

$$\mathbf{W}^T = [-1.5 \quad 1 \quad 1]_{1 \times 3}$$

$$\mathbf{W}^T \mathbf{X} = [-1.5 \quad -0.5 \quad -0.5 \quad 0.5] \xrightarrow{\text{Nonlinear threshold}} [0 \quad 0 \quad 0 \quad 1]$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
 $<0 \quad <0 \quad <0 \quad >0$

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1



OR function

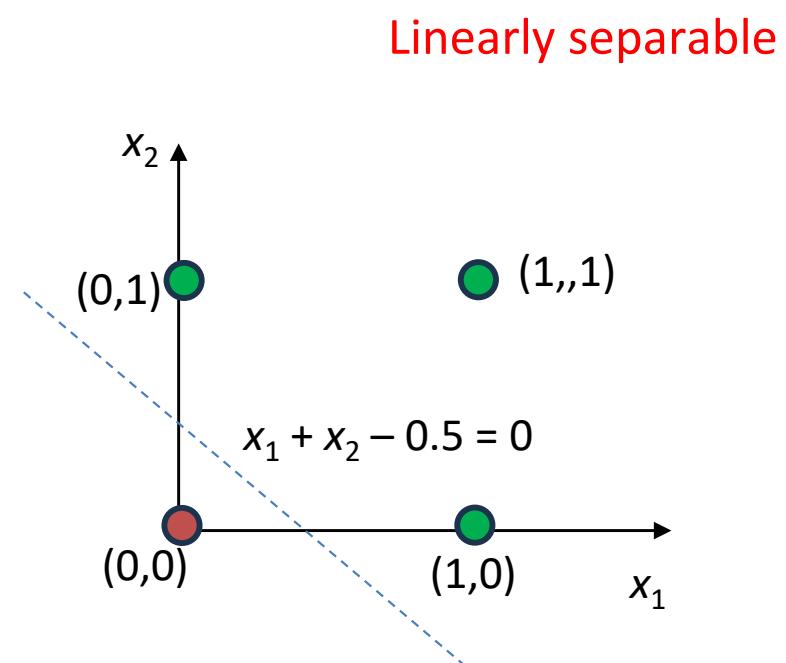
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

$$x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \quad W = \begin{pmatrix} -0.5 \\ 1 \\ 1 \end{pmatrix}$$

\curvearrowright

$$X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix} \quad W^T = \begin{bmatrix} -1.5 & 1 & 1 \end{bmatrix}$$

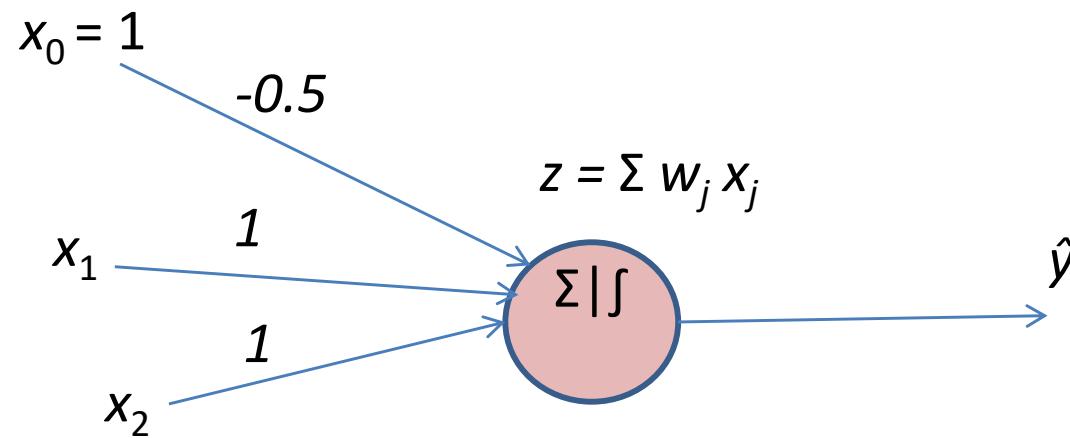
\curvearrowright



OR function

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix} \quad \mathbf{W}^T = [-1.5 \ 1 \ 1] \quad (x_1 + x_2) \cdot ($$

$$\mathbf{W}^T \mathbf{X} = [-0.5 \ 0.5 \ 0.5 \ 0.5] \xrightarrow{\text{Nonlinear threshold}} [0 \ 1 \ 1 \ 1]$$



XOR function

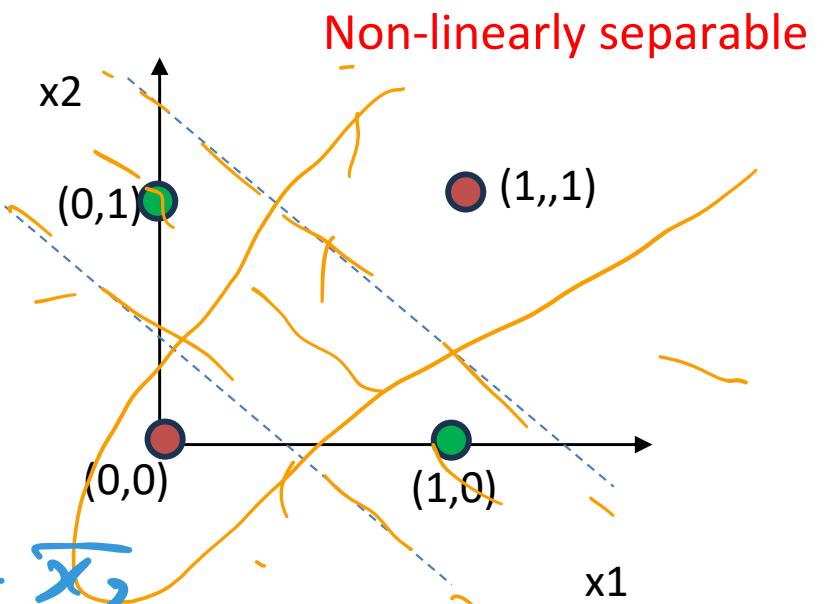
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

$$x_1 \oplus x_2 =$$

$$x_1 +$$



$$\overline{x_1 \cdot x_2} = \overline{x_1} + \overline{x_2}$$



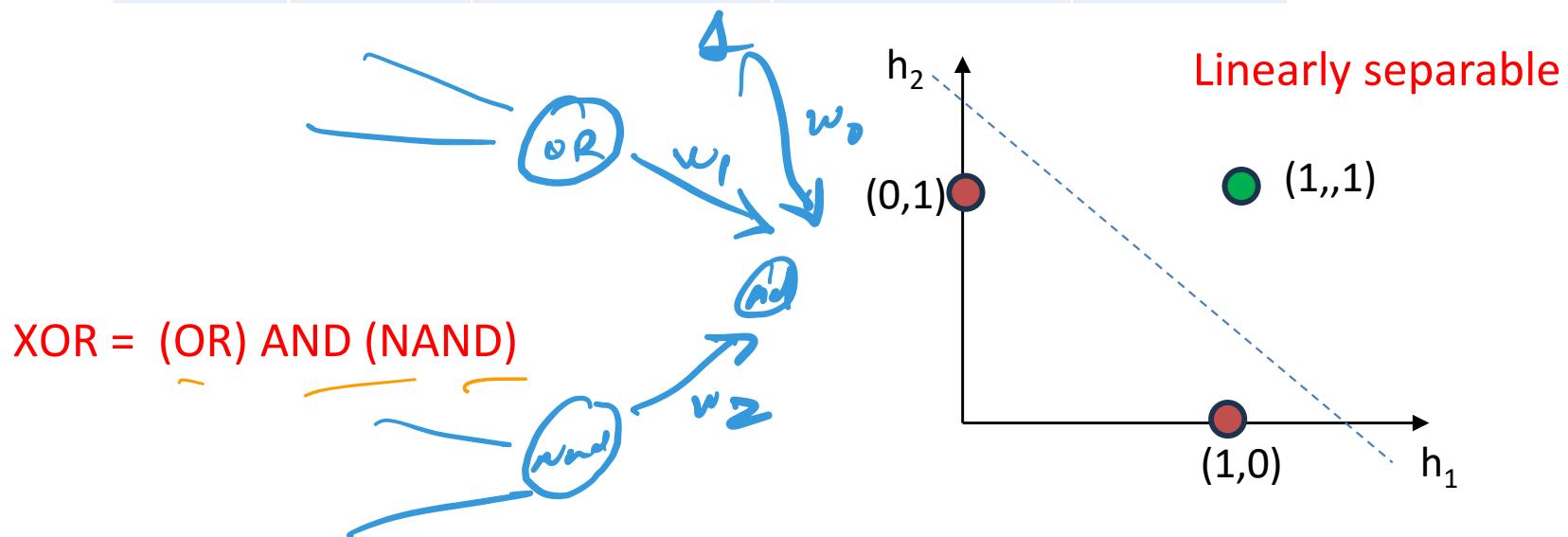
XOR function

$$x_1 \text{ XOR } x_2 = (x_1 + x_2) \cdot (\bar{x}_1 + \bar{x}_2)$$

OR operation

NAND operation

x_1	x_2	$h_1 = x_1 + x_2$	$h_2 = \bar{x}_1 + \bar{x}_2$	y
0	0	0	1	0
0	1	1	1	1
1	0	1	1	1
1	1	1	0	0



XOR function

$$X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix}$$

$$W_1 = \begin{pmatrix} -0.5 & 1.5 \\ 1 & -1 \\ 1 & -1 \end{pmatrix}$$

Weights for OR

Weights for NAND operation (complement of AND weights)

$$W_1^T X = \begin{pmatrix} -0.5 & 0.5 & 0.5 & 1.5 \\ 1.5 & 0.5 & 0.5 & -0.5 \end{pmatrix}$$

Nonlinear threshold

$$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

$$\begin{matrix} h_1 \\ h_2 \end{matrix}$$

$$W_2 = \begin{pmatrix} -1.5 \\ 1 \\ 1 \end{pmatrix}$$

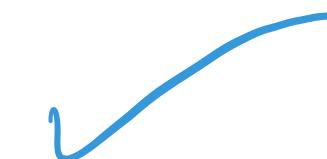
Weights for AND operation

$$W_2^T X' = [0 \ 1 \ 1 \ 0]$$

$$h^T W_2 = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} -1.5 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -0.5 \\ 0.5 \\ 0.5 \\ -0.5 \end{pmatrix}$$

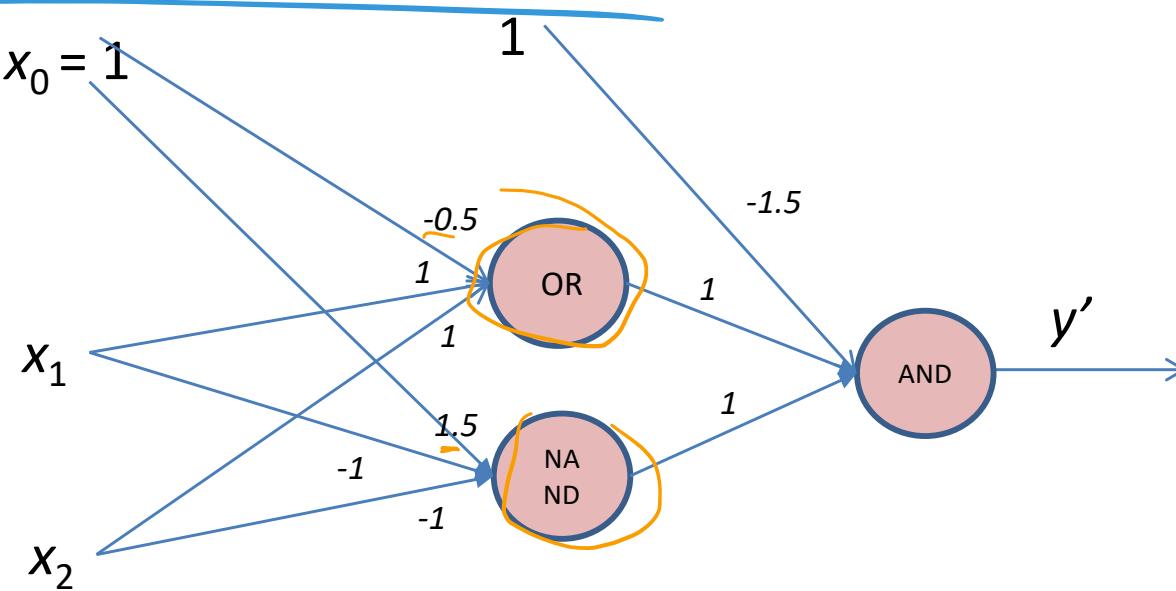
Nonlinear threshold

$$\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$



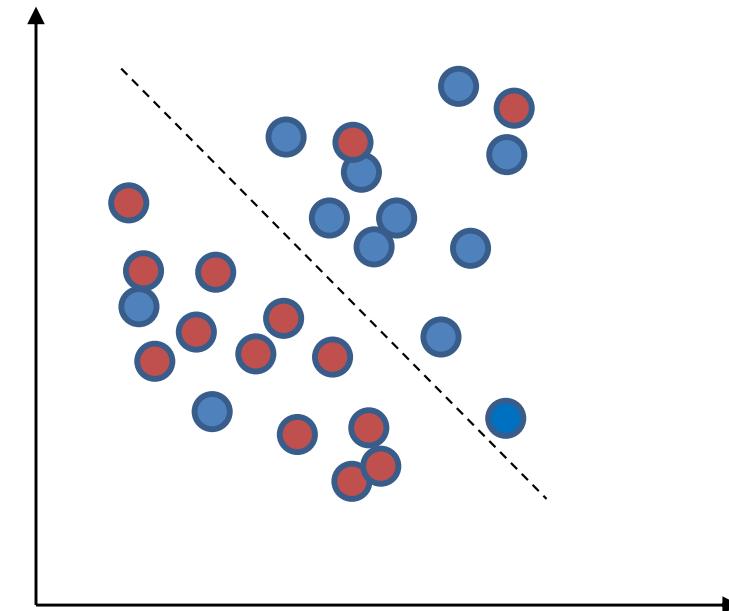
XOR function

- Two layer network
- 1^s layer:
 - 2 neurons: one for OR, one for NAND
- 2nd layer:
 - 1 neuron: AND operation

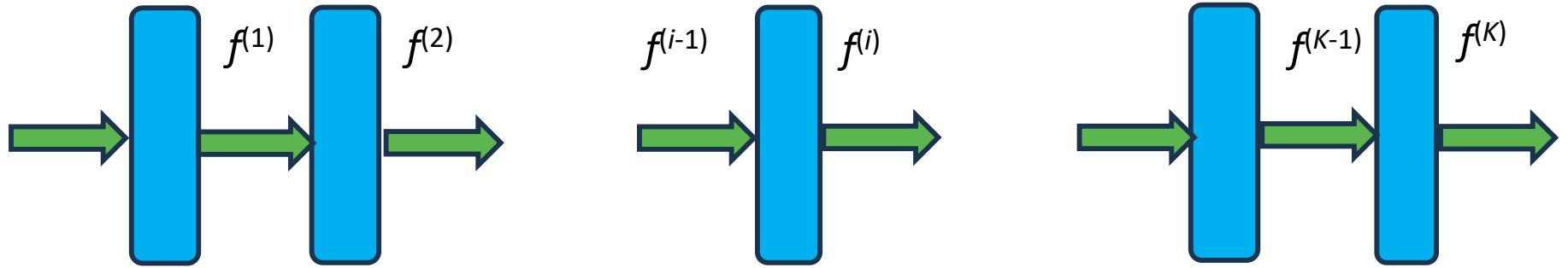


Single perceptron

- Single perceptron cannot deal with data which is not linearly separable
 - Have to be flexible with convergence statement while determining weights...
 - Till almost (say 90%) points are satisfying condition
 - Leads to few errors
 - May not be acceptable in critical real-world applications



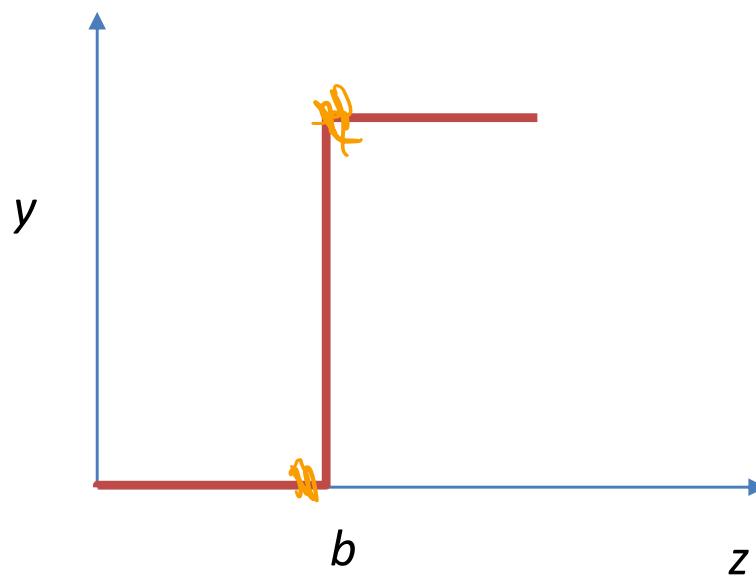
Neural Network Function



$$f^{(K)} \left(f^{(K-1)} \left(\dots f^{(i)} \dots (f^{(2)} (f^{(1)} (x))) \right) \right)$$

Perceptron Logic

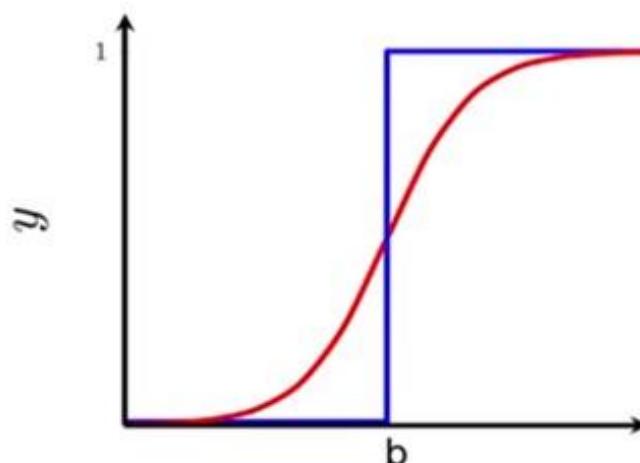
- A perceptron will fire if weighted sum of inputs is greater than threshold ($b = -w_0$)
 - Thresholding logic is harsh



If $x=0.51$, watch movie, if $x=0.49$, do not watch??

For real world problems, we need a smoother decision function

Sigmoid neuron



$$y = \frac{1}{1+e^{-(w^T x + b)}}$$

$$\tau'(z) = \tau(1 - \tau)$$

If $z = (w^T x + b) \rightarrow \infty$, then $y = 1$

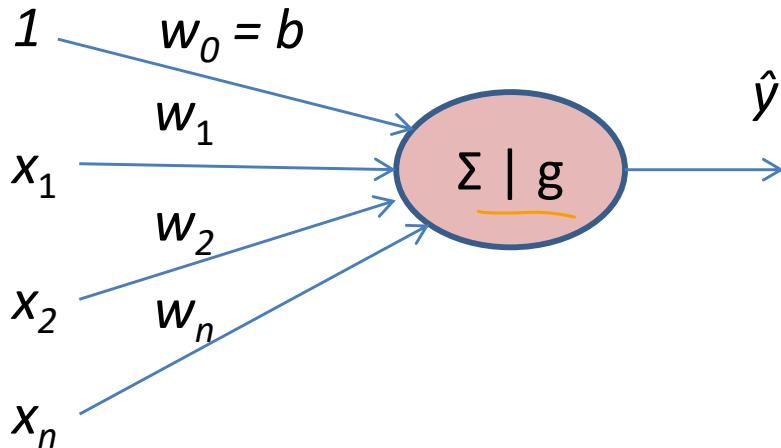
If $z = (w^T x + b) \rightarrow -\infty$, then $y = 0$

If $z = (w^T x + b) = 0$, then $y = 0.5$

Range of sigmoid function: 0 to 1

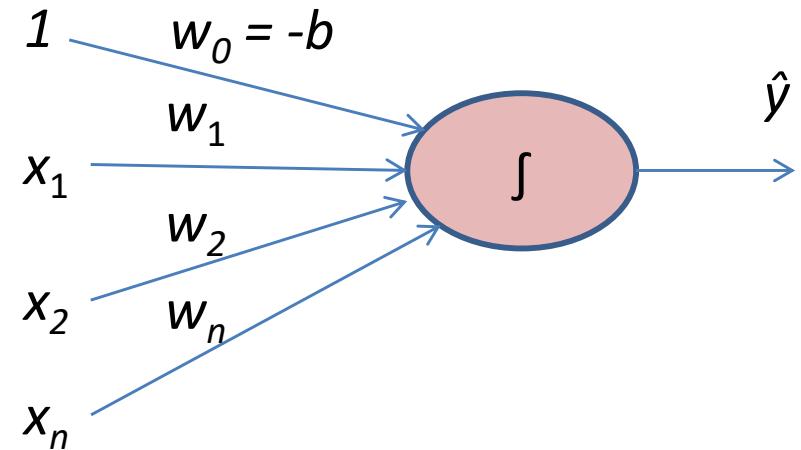
Can be interpreted as probability

Perceptron vs Sigmoid



$$\begin{aligned}\hat{y} &= 1 \quad \text{if } \sum w_i x_i + b \geq 0 \\ &= 0 \quad \text{if } \sum w_i x_i + b < 0\end{aligned}$$

- Not smooth
- Not continuous
- Not differentiable

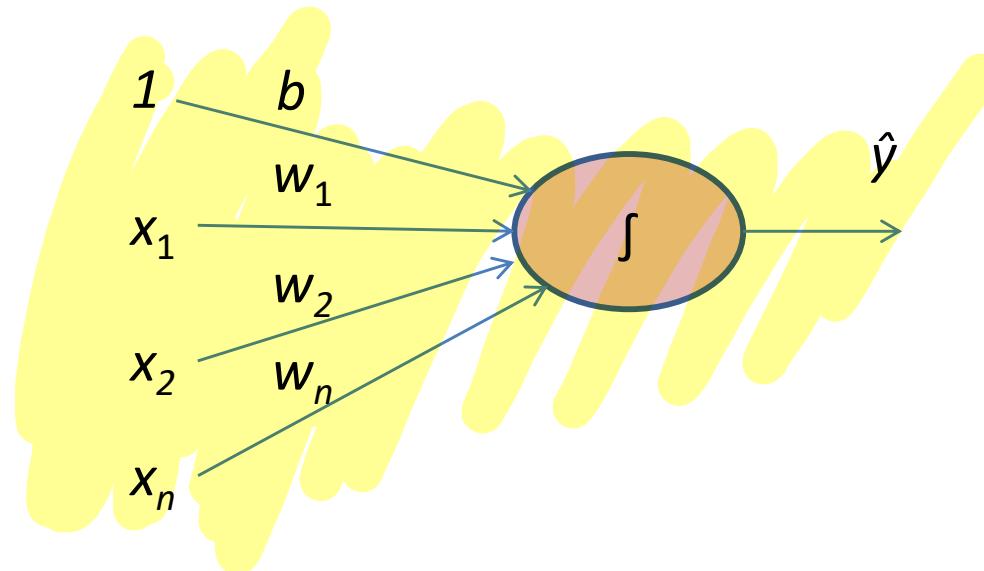


$$\forall i = 1 \dots n$$

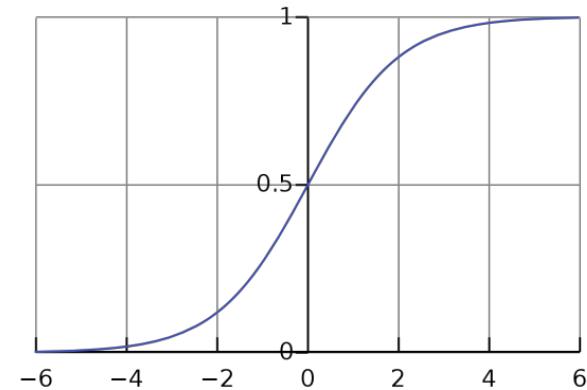
$$y = \frac{1}{1+e^{-(w^T x+b)}}$$

- Smooth
- Continuous
- Differentiable

Perceptron – Forward propagation



$$\begin{aligned}\hat{y} &= g(z) = g(\sum w_i x_i) & \forall i = 0 \dots n \\ &= g(b + \sum w_i x_i) & \forall i = 1 \dots n \\ &= g(b + \mathbf{X}^T \mathbf{W})\end{aligned}$$



$$\text{where, } \mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} \quad \mathbf{W} = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_n \end{pmatrix}$$

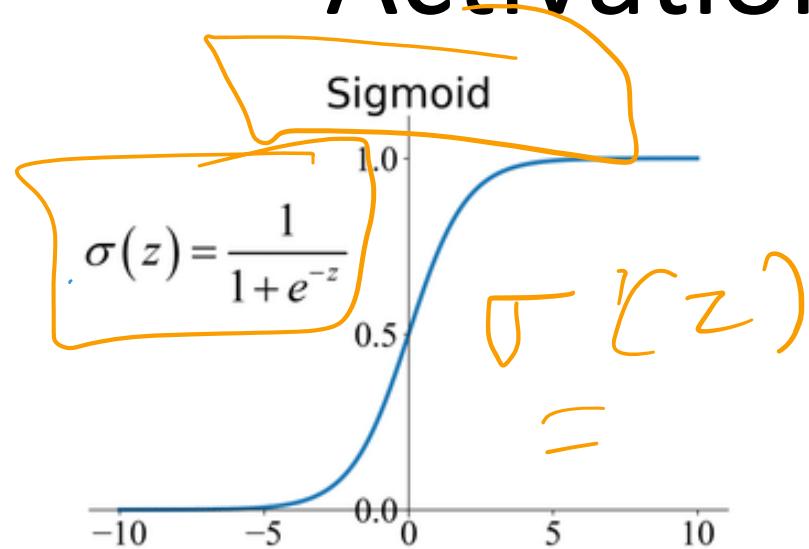
b : bias
 $w^T x + b$

Hand-drawn orange annotations include a bracketed 'z' above the summation node, a bracketed 'w' next to the weight vector, and a bracketed 'x' next to the input vector.

$$\frac{e^z - e^{-z}}{e^z + e^{-z}}$$

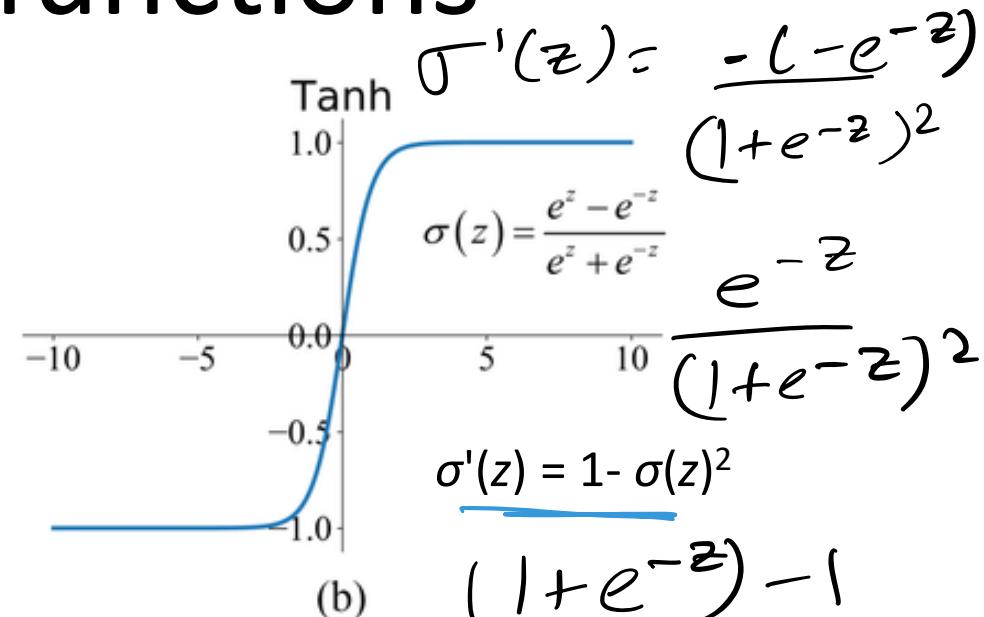
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Activation functions

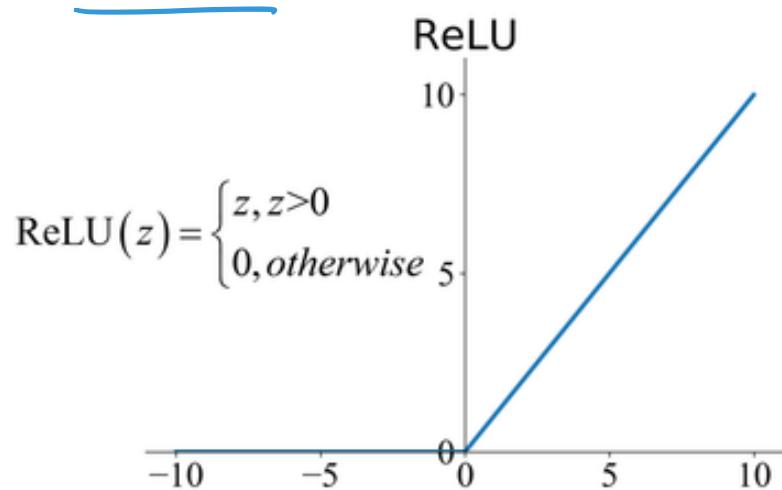


$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

(a)



(b)

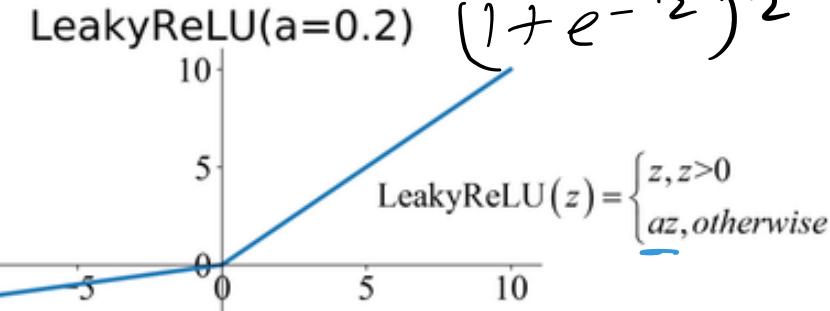


$$\sigma'(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

(c)

$$\sigma(z) - (\sigma(z))^2$$

$$\sigma(z) \left(1 - \frac{\sigma'(z)}{\sigma(z)} \right) = \begin{cases} 1, & z \geq 0 \\ a, & z < 0 \end{cases}$$

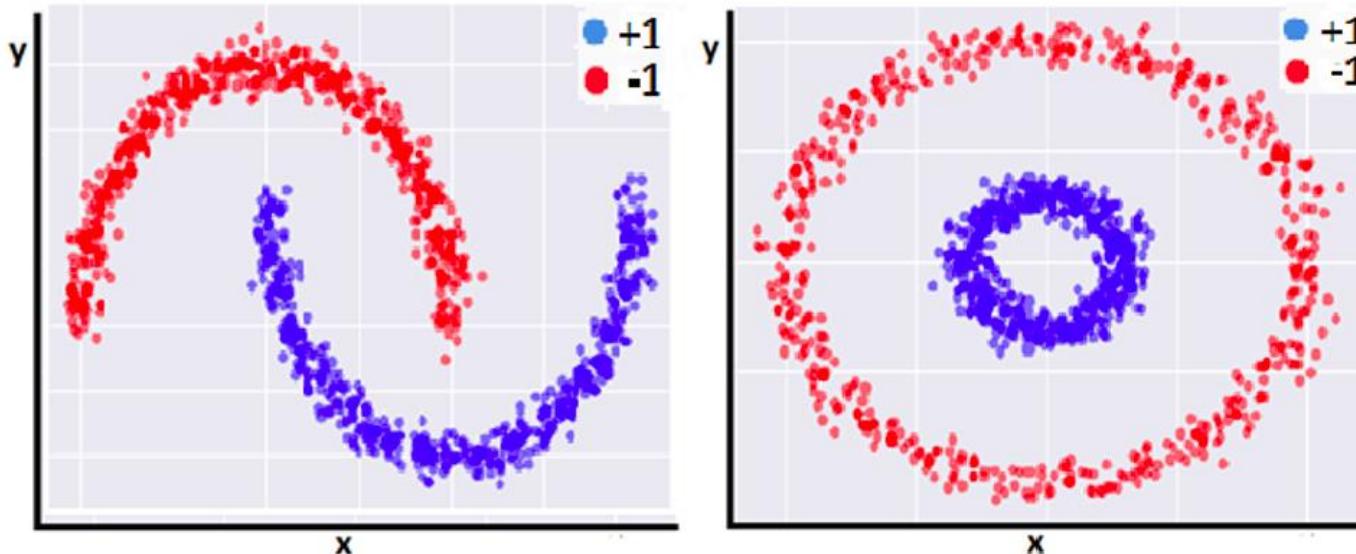


$$\text{LeakyReLU}(z) = \begin{cases} z, & z > 0 \\ az, & \text{otherwise} \end{cases}$$

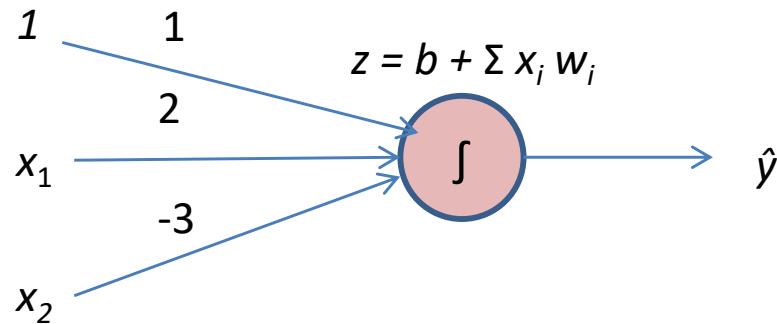
$$\tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Importance - Activation functions

- Introduce non-linearity in the network
 - Allows to deal with non-linear data
 - Allows to approximate complex functions



Perceptron – Example

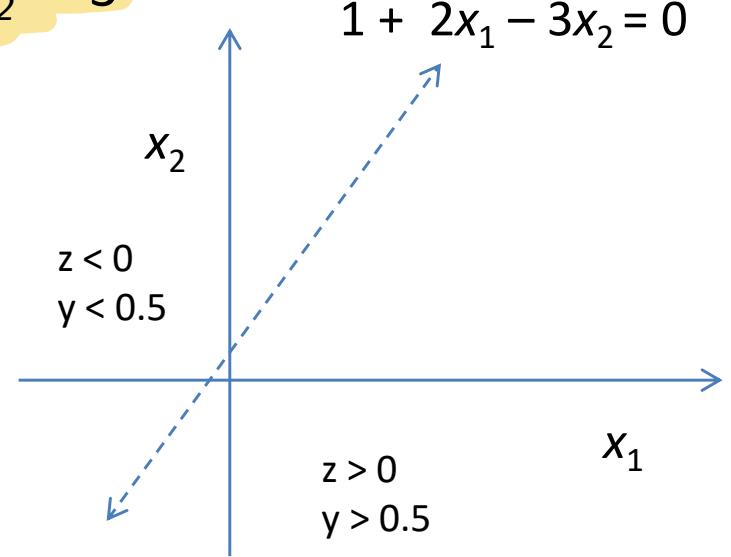


$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad W = \begin{pmatrix} 2 \\ -3 \end{pmatrix} \quad b = 1, w_1 = 2, w_2 = -3$$

$$\hat{y} = g(b + X^T W)$$

$$\hat{y} = g(1 + \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}^T \begin{pmatrix} 2 \\ -3 \end{pmatrix})$$

$$\hat{y} = g(1 + 2x_1 - 3x_2)$$



Typical Machine Learning Problem

x_1	x_2	x_3	...	x_n	y

- Actor, Producer, Genre.....Watch a movie?
- Salinity, Density....Will find oil?
- Salary, Occupation, Tax filed....Commit fraud? $y = f(x)$
- **Model:** Approximation of relation between x and y
$$\hat{y} = f'(x, w)$$
 - w parameter needs to be learned from data
- **Learning algorithm:** Algorithm for learning parameters of model
- **Objective/Loss/Error function:** To guide the learning algorithm – learning algorithm should aim to minimize loss function

Typical Machine Learning Problem

- Actor, Producer, Genre....Watch a movie?
- Input for training:
$$\{x_i, y_i\} \quad i = 1..m \quad m \text{ pairs of } (x, y)$$

- Model: $y = f'(x, w, b)$
- Learning algorithm: Ex. Gradient Descent
- Objective/Loss/Error function:
 - One possibility: $L(w, b) = \sum (\hat{y}_i - y_i)^2$ for $i = 1..m$ *Real \rightarrow Regress.*
- Learning algorithm should find w to minimize the objective function
 - Hope to find sigmoid function such that given x_i, y_i lie on that function
$$-\left(\sum p_i \log(q_i) \right)$$

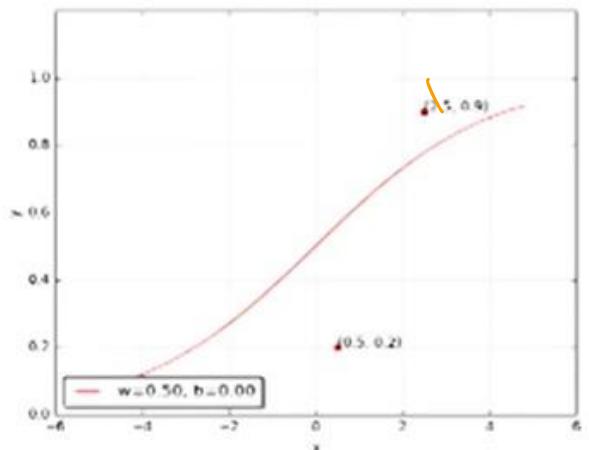
Example

↑
classification

- Ex. data points: $(0.5, 0.2)$, $(2.5, 0.9)$
- At end of training, expect to find w^* , b^* such that
 $f(0.5) \rightarrow 0.2$, $f(2.5) \rightarrow 0.9$
- Loss = $(1/2) * \sum(y - \hat{y})^2$ given: $w = 0.5$, $b = 0$
 $= (1/2) * [(0.2 - f(0.5))^2 + (0.9 - f(2.5))^2]$
 $= 0.073$

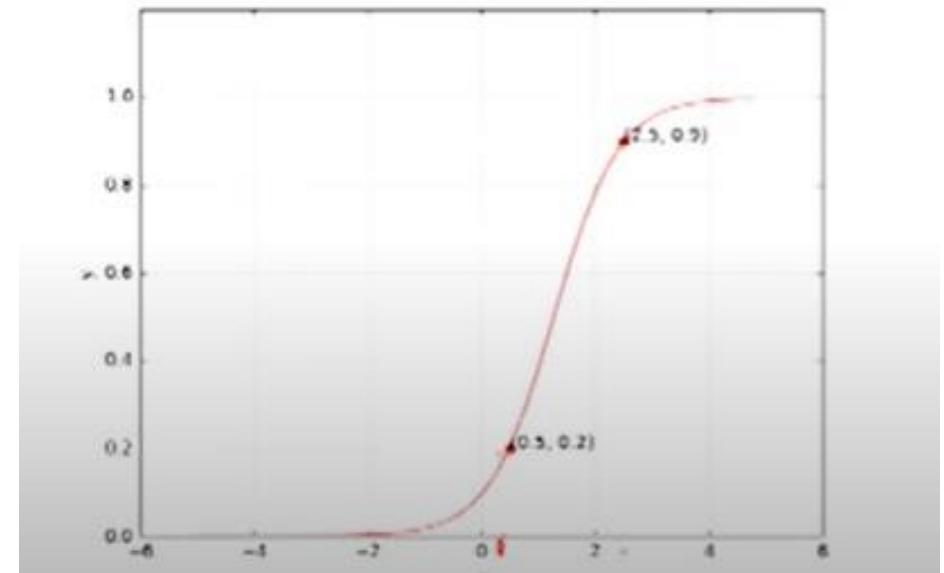
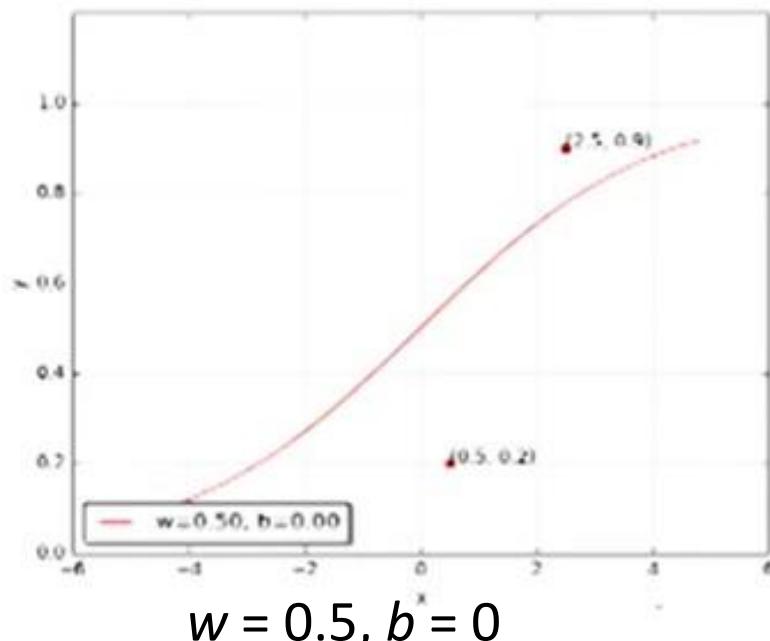
$$y = \frac{1}{1+e^{-(w^T x+b)}}$$

$$w = 0.5, b = 0$$

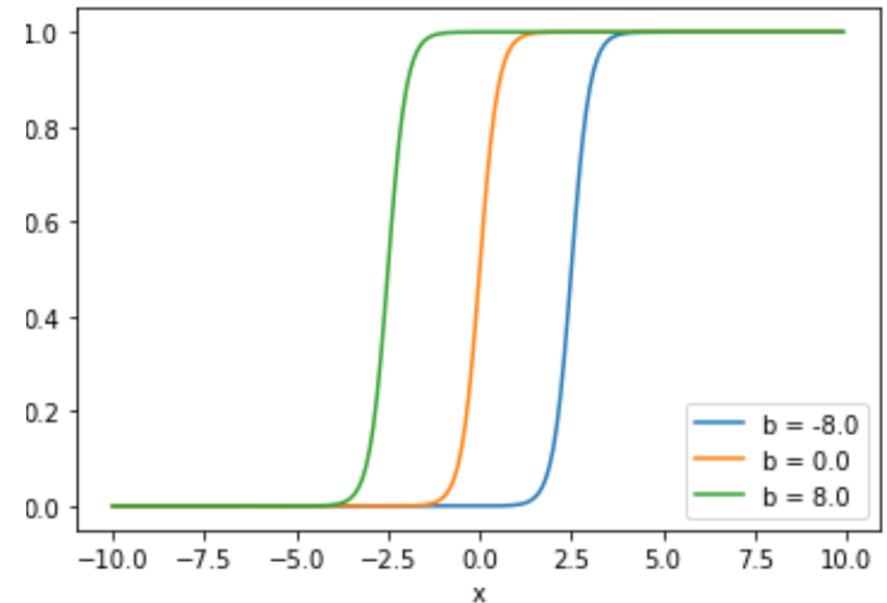
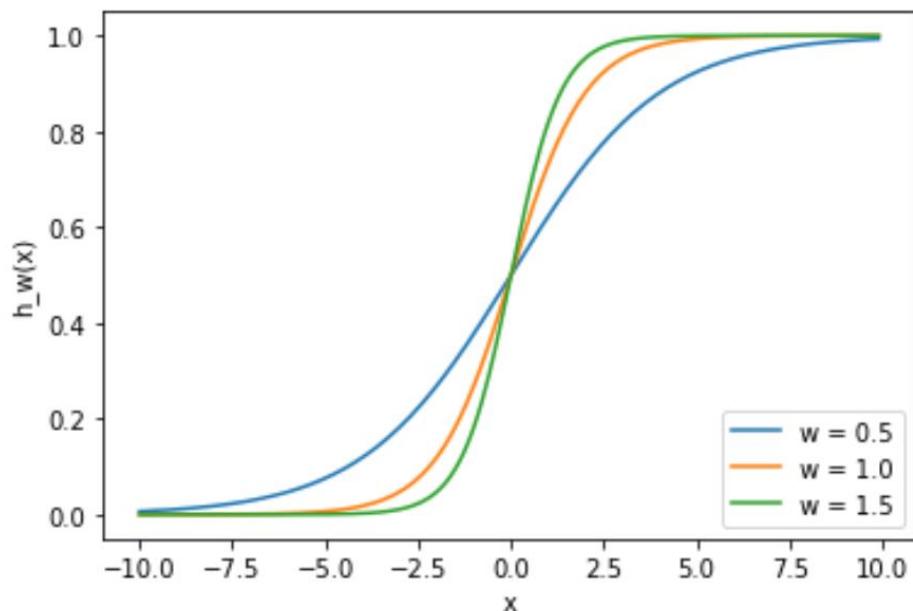


Example

- Ex. data points: (0.5,0.2), (2.5, 0.9)
- Hope to find a sigmoid function that the data points lie on the function



Sigmoid function



Changes w.r.t w

↳ weights change
slope

Changes w.r.t b

↓
bias change

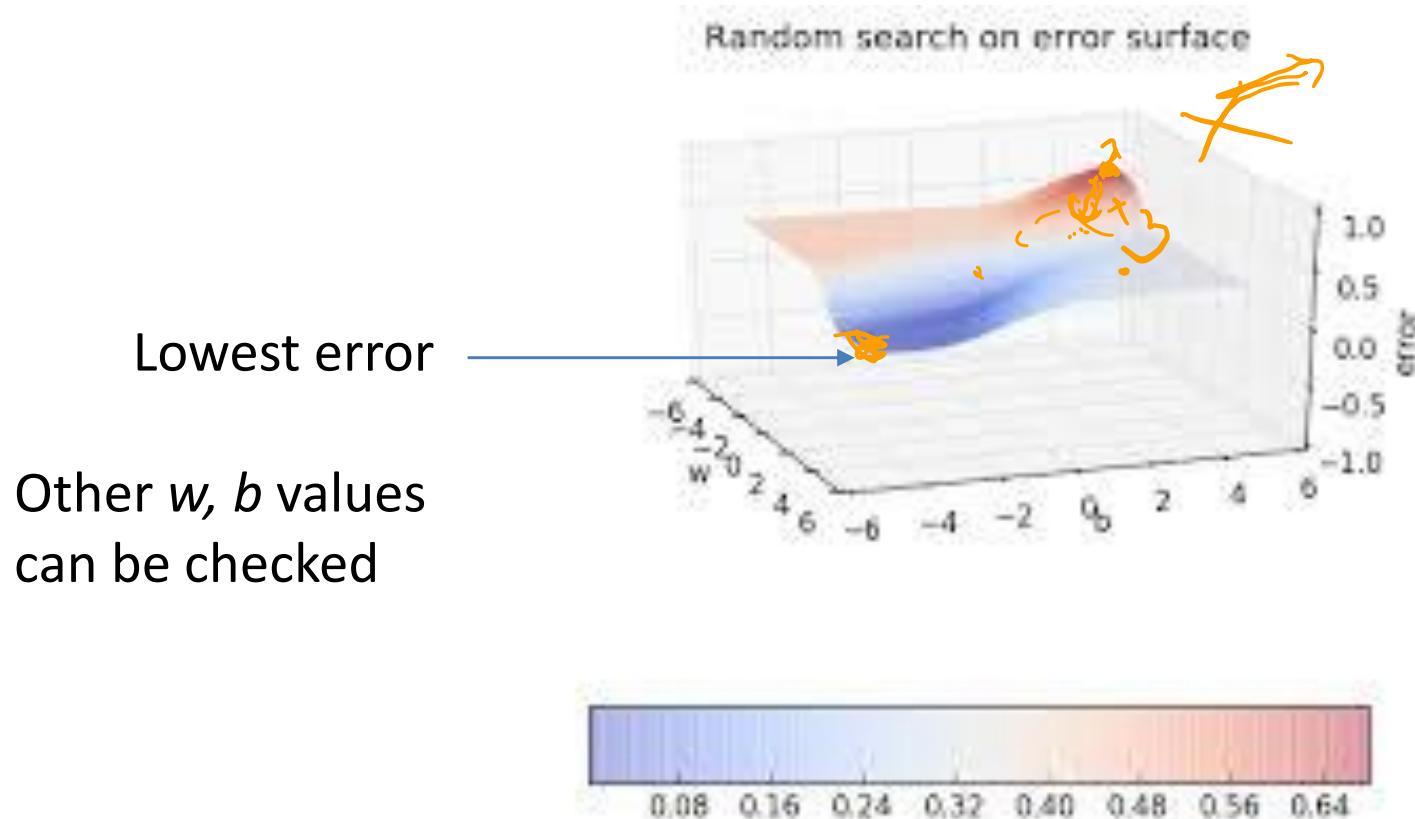
position.

Guesses!!

w	b	Loss
0.5	0	0.073
-0.1	0	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

Error surface

- 2 data points, 2 parameters (w, b)



- How to handle more data points and parameters??

Gradient Descent

- Find a way to traverse error surface to reach minimum value quickly
- Parameters: $\theta = [w, b]$
- Change in parameters: $\Delta\theta = [\Delta w, \Delta b]$
- New parameters: $\theta_{\text{new}} = \theta + \Delta\theta = [w_{\text{new}}, b_{\text{new}}]$
 - Change with a small stride η
- How to choose Δw and Δb ?

$$w_{t+1} = w_t - \eta$$

$$\Delta w = \frac{\partial L}{\partial w}$$

$$\Delta b = \frac{\partial L}{\partial b}$$

$$w_{\text{new}} = w + \eta \cdot \Delta w$$

$$b_{\text{new}} = b + \eta \cdot \Delta b$$

$\nabla \rightarrow \text{Gradient}$

$\Delta \rightarrow \text{step}$

$$(\nabla L(\theta))^T (\nabla_w L(\theta)) < 0$$

Gradient Descent

- According to Taylor series:

$$\begin{aligned}L(\theta + \eta \Delta \theta) &= L(\theta) + \eta * (\Delta \theta)^T \nabla_{\theta} L(\theta) + (\eta^2 / 2!) * (\Delta \theta)^T \nabla_{\theta}^2 L(\theta) (\Delta \theta) + \dots \\&= L(\theta) + \eta (\Delta \theta)^T \nabla_{\theta} L(\theta) \quad (\text{if } \eta \text{ is small})\end{aligned}$$

- New loss should be less than old loss
- $(\eta \Delta \theta)$ would be favourable only if:

$$L(\theta + \eta \Delta \theta) < L(\theta)$$

$$L(\theta + \eta \Delta \theta) - L(\theta) < 0$$

- Implies: $(\Delta \theta)^T \nabla_{\theta} L(\theta) < 0$ (η is positive constant)

$$\nabla_{\theta} L(\theta) = \frac{\partial L(\theta)}{\partial \theta}$$

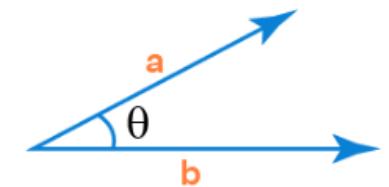
Gradient Descent

- Desired: $(\Delta\theta)^T \nabla_{\theta} L(\theta) < 0$ $\frac{\partial L(\theta)}{\partial \theta}$
 - Want $\Delta\theta$ to be as negative as possible
 - Dot product of two vectors: product of their magnitudes multiplied by the cosine of the angle between them

- Let β be the angle between $\Delta\theta$ and $\nabla_{\theta} L(\theta)$

- Then:

$$-1 \leq \cos(\beta) = \frac{(\Delta\theta)^T \nabla_{\theta} L(\theta)}{\|\Delta\theta\| * \|\nabla_{\theta} L(\theta)\|} \leq 1$$



$$\cos \theta = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$$

Most negative: $\cos(\beta) = -1$ when $\beta=180^\circ$

$\rightarrow \Delta\theta$ should be such that it is at 180° to the gradient $\nabla L(\theta)$

Gradient Descent

- $(\Delta\theta)^\top \nabla L(\theta) < 0$
- Move in direction opposite to gradient (180° w.r.t. the gradient)

$$(\Delta\theta)^\top (\nabla L(\theta)) < 0$$

$$w_{t+1} = w_t - \eta \nabla w_t$$

$$b_{t+1} = b_t - \eta \nabla b_t$$

where, $\nabla w_t = \frac{\partial L(w, b)}{\partial w}$

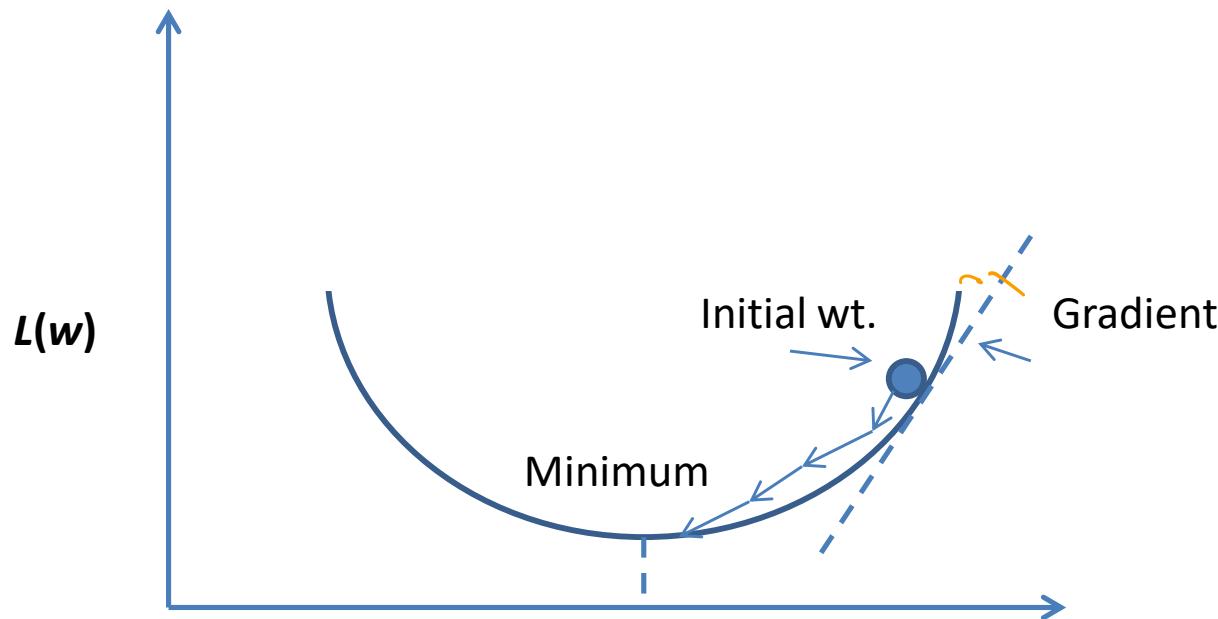
$$\nabla b_t = \frac{\partial L(w, b)}{\partial b}$$

- Repeat till convergence

at $w = w_t$ and $b = b_t$

at $w = w_t$ and $b = b_t$

Gradient Descent



Handwritten notes:

\mathbf{w}

$\mathbf{w} := \mathbf{w} - \eta * \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$

$\mathbf{P}^2 = \begin{bmatrix} -SSE \\ SST \end{bmatrix}$

Randomly pick (w_0, w_1)

Compute L

Compute gradient
 $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$ - gives ascent at that point

Take small step in opposite direction of gradient

Repeat until convergence

$$SSE = \sum (\hat{y} - y)^2$$

$$SST = \sum (y - \bar{y})^2$$

1

$$SSR = \sum (\hat{y} - \bar{y})^2$$

Gradient Descent

Assuming, $L(w, b) = (1/2) * \sum (\hat{y} - y)^2$

$$\hat{y} = f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Assuming one data point only:

$$\frac{\partial L(w, b)}{\partial w} = (1/2) * [2 * (f(x) - y) * \frac{\partial (f(x) - y)}{\partial w}]$$

$$= (f(x) - y) * \frac{\partial}{\partial w} \left[\frac{1}{1 + e^{-(wx+b)}} \right]$$

$$\frac{\partial L(w, b)}{\partial w} = (f(x) - y) * f(x) * (1 - f(x)) * x = (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x$$

$$\frac{\partial L(w, b)}{\partial b} = (f(x) - y) * \frac{\partial}{\partial b} \left[\frac{1}{1 + e^{-(wx+b)}} \right]$$

$$= (f(x) - y) * f(x) * (1 - f(x)) = (\hat{y} - y) * \hat{y} * (1 - \hat{y})$$

Gradient Descent

$$\nabla w = \sum [(\hat{y}_i - y_i) * \hat{y}_i * (1 - \hat{y}_i) * x_i]$$

$$\nabla b = \sum [(\hat{y}_i - y_i) * \hat{y}_i * (1 - \hat{y}_i)]$$

for all points

for all points

- Algorithm:

1. Initialize weights randomly
2. Loop until convergence:

1. Compute gradient $\frac{\partial L(w, b)}{\partial w}, \frac{\partial L(w, b)}{\partial b}$

2. Update weights $w := w - \eta * \frac{\partial L(w, b)}{\partial w}$

$$b := b - \eta * \frac{\partial L(w, b)}{\partial b}$$

Return weights

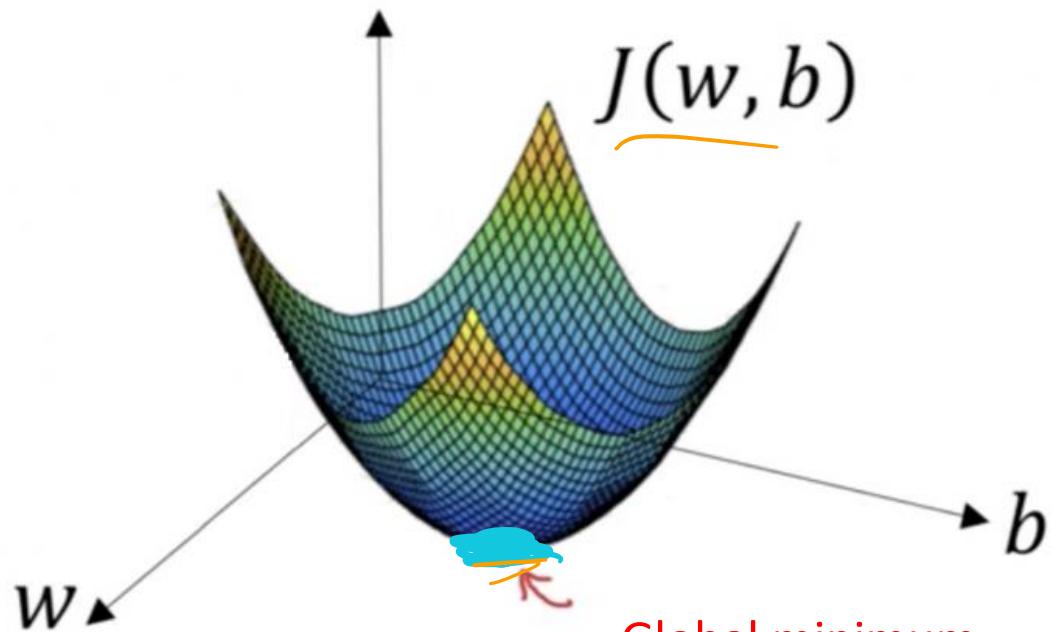
$$\frac{\partial f(x)}{\partial w}$$

$$f(\theta) = \frac{1}{1 + e^{\frac{w_0x + b}{z}}}$$

$$\frac{\partial f(x)}{\partial z} \times \frac{\partial z}{\partial w}$$

$$f(z)(1-\sigma) \cancel{f(z)}$$

Gradient Descent



Global minimum

Assuming $J(W)$ to be a convex function