



CS 540 Introduction to Artificial Intelligence Perceptron

Sharon Yixuan Li
University of Wisconsin-Madison

March 4, 2021

Today's outline

- Naive Bayes (cont.)
- Single-layer Neural Network (Perceptron)



Part I: Naïve Bayes (cont.)

Example 1: Play outside or not?

- If weather is sunny, would you likely to play outside?

Posterior probability $p(\text{Yes} | \text{Sun})$ vs. $p(\text{No} | \text{Sun})$

Example 1: Play outside or not?

- If weather is sunny, would you likely to play outside?

Posterior probability $p(\text{Yes} | \text{Sun})$ vs. $p(\text{No} | \text{Sun})$

- Weather = {Sunny, Rainy, Overcast}
- Play = {Yes, No}
- Observed data {Weather, play on day m }, $m=\{1,2,\dots,N\}$

Example 1: Play outside or not?

- If weather is sunny, would you likely to play outside?

Posterior probability $p(\text{Yes} | \text{Sun})$ vs. $p(\text{No} | \text{Sun})$

- Weather = {Sunny, Rainy, Overcast}
- Play = {Yes, No}
- Observed data {Weather, play on day m }, $m=\{1,2,\dots,N\}$

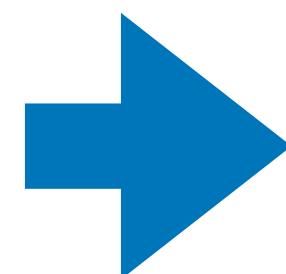
$$p(\text{Play} | \text{Sun}) = \frac{p(\text{Sun} | \text{Play}) p(\text{Play})}{p(\text{Sun})}$$

Bayes rule

Example 1: Play outside or not?

- Step 1: Convert the data to a frequency table of Weather and Play

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No



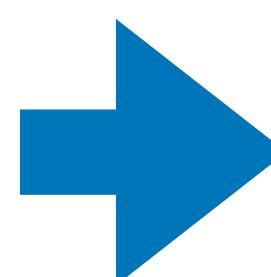
Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Example 1: Play outside or not?

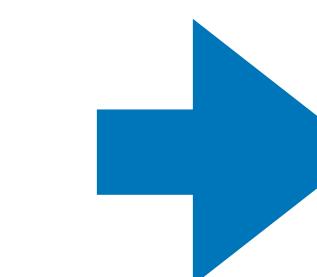
Step 1: Convert the data to a frequency table of Weather and Play

Step 2: Based on the frequency table, calculate **likelihoods** and **priors**

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No



Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9



Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

$$p(\text{Play} = \text{Yes}) = 0.64$$

$$p(\text{Sun} | \text{Yes}) = 3/9 = 0.33$$

Example 1: Play outside or not?

Step 3: Based on the likelihoods and priors, calculate posteriors

$$P(\text{Yes} | \text{Sun}) = P(\text{Sun} | \text{Yes}) * P(\text{Yes}) / P(\text{Sun})$$

$$P(\text{No} | \text{Sun}) = P(\text{Sun} | \text{No}) * P(\text{No}) / P(\text{Sun})$$

Example 1: Play outside or not?

Step 3: Based on the likelihoods and priors, calculate posteriors

$$\begin{aligned} P(\text{Yes} | \text{Sun}) &= P(\text{Sun} | \text{Yes}) * P(\text{Yes}) / P(\text{Sun}) \\ &= 0.33 * 0.64 / 0.36 \\ &= 0.6 \end{aligned}$$

$$\begin{aligned} P(\text{No Sun}) &= P(\text{No Sun} | \text{No}) * P(\text{No}) / P(\text{Sun}) \\ &= 0.4 * 0.36 / 0.36 \\ &= 0.4 \end{aligned}$$

$P(\text{Yes} | \text{Sun}) > P(\text{No} | \text{Sun})$ go outside and play!

Bayesian classification

$$\hat{y} = \arg \max_y p(y | \mathbf{x}) \quad (\text{Posterior})$$

(Prediction)

$$= \arg \max_y \frac{p(\mathbf{x} | y) \cdot p(y)}{p(\mathbf{x})} \quad (\text{by Bayes' rule})$$

$$= \arg \max_y p(\mathbf{x} | y)p(y)$$

Bayesian classification

What if \mathbf{x} has multiple attributes $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\hat{y} = \arg \max_y p(y | X_1, \dots, X_k) \quad (\text{Posterior})$$

(Prediction)

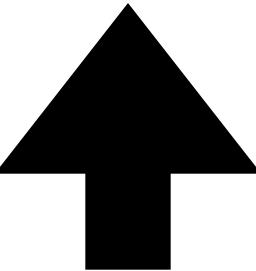
Bayesian classification

What if \mathbf{x} has multiple attributes $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\hat{y} = \arg \max_y p(y | X_1, \dots, X_k) \quad (\text{Posterior})$$

(Prediction)

$$= \arg \max_y \frac{p(X_1, \dots, X_k | y) \cdot p(y)}{p(X_1, \dots, X_k)} \quad (\text{by Bayes' rule})$$



Independent of y

Bayesian classification

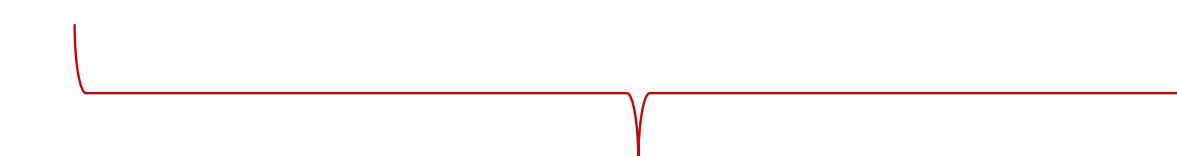
What if \mathbf{x} has multiple attributes $\mathbf{x} = \{X_1, \dots, X_k\}$

$$\hat{y} = \arg \max_y p(y | X_1, \dots, X_k) \quad (\text{Posterior})$$

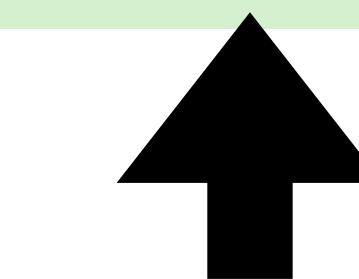
(Prediction)

$$= \arg \max_y \frac{p(X_1, \dots, X_k | y) \cdot p(y)}{p(X_1, \dots, X_k)} \quad (\text{by Bayes' rule})$$

$$= \arg \max_y p(X_1, \dots, X_k | y) p(y)$$



Class conditional
likelihood

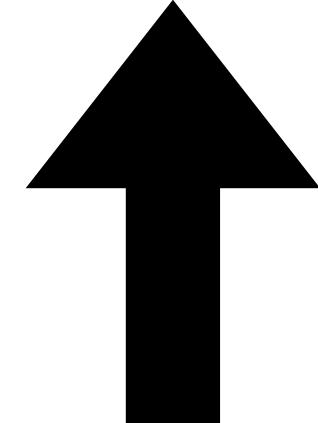


Class prior

Naïve Bayes Assumption

Conditional independence of feature attributes

$$p(X_1, \dots, X_k | y)p(y) = \prod_{i=1}^k p(X_i | y)p(y)$$



Easier to estimate
(using MLE!)



Part I: Single-layer Perceptron

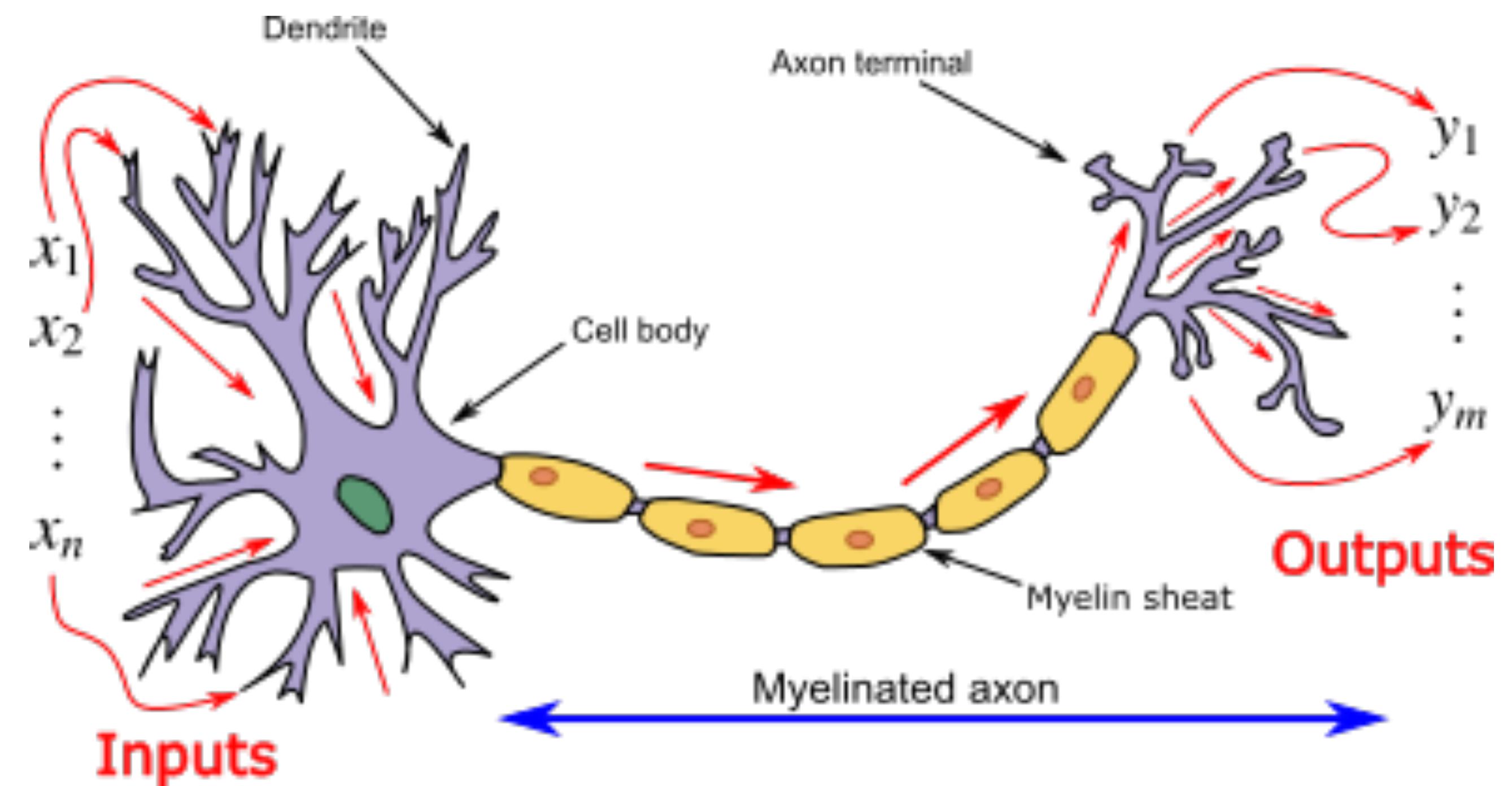
How to classify

Cats vs. dogs?



Inspiration from neuroscience

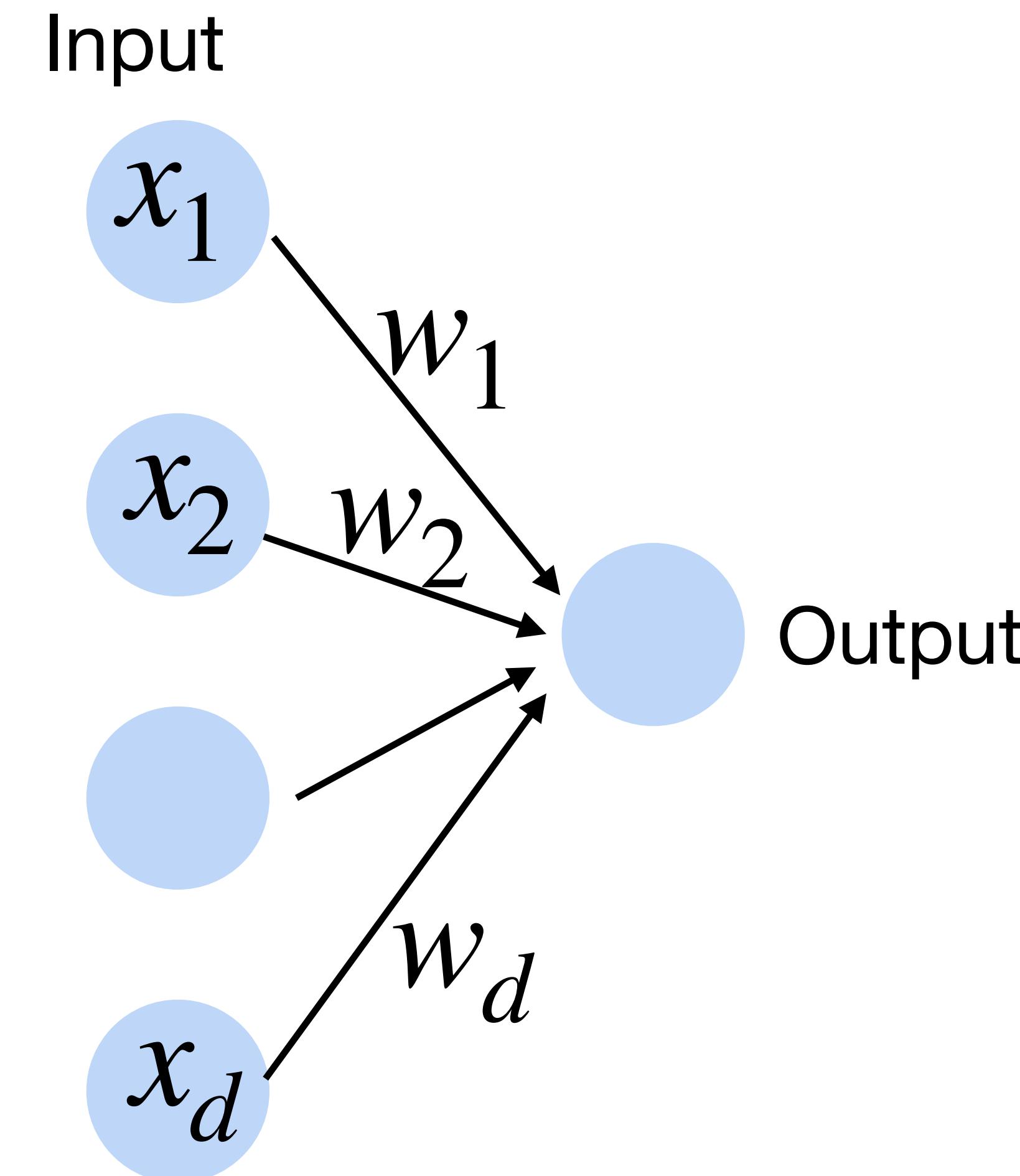
- Inspirations from human brains
- Networks of **simple** and **homogenous** units



(wikipedia)

Perceptron

Cats vs. dogs?

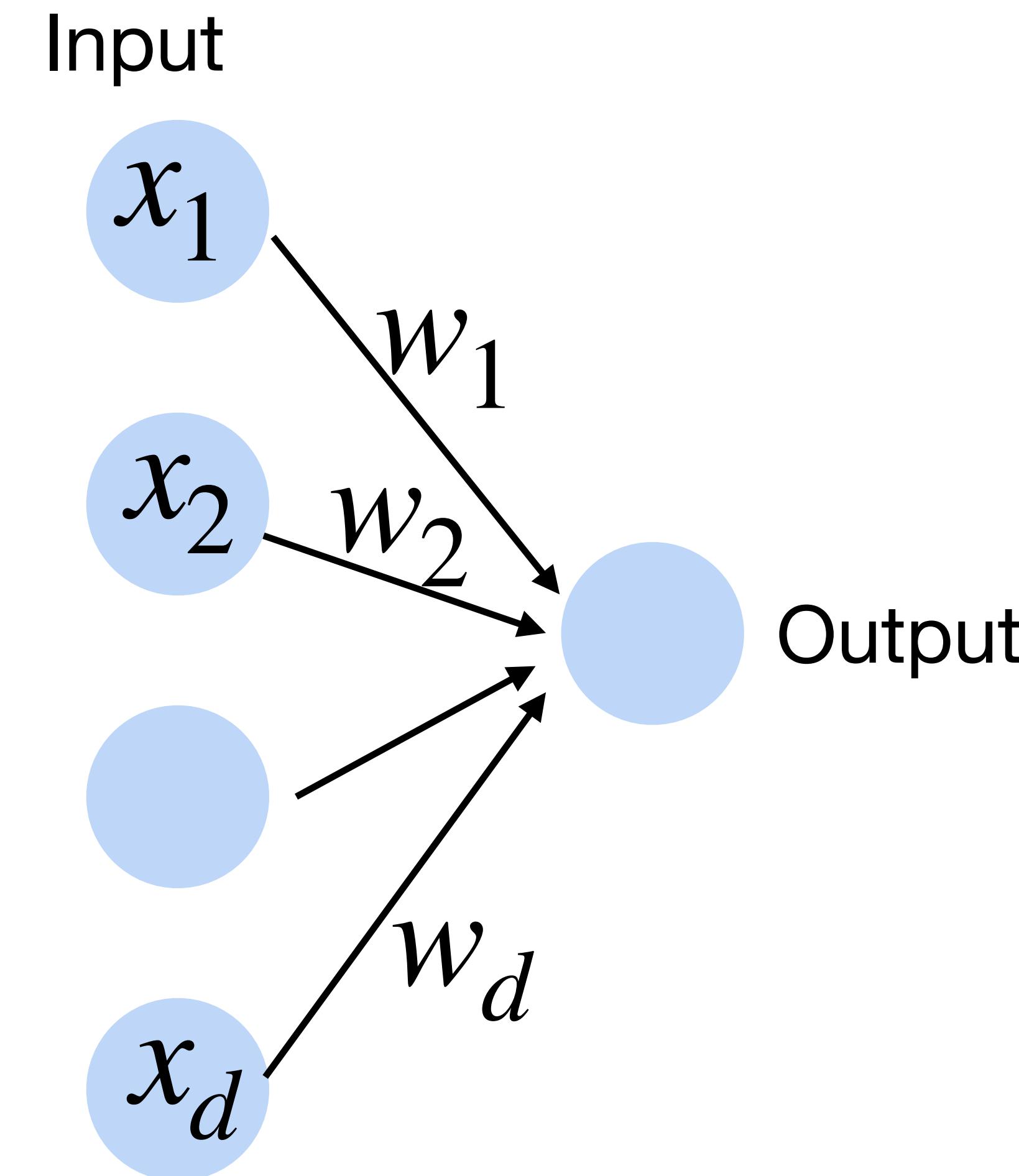


Linear Perceptron

- Given input \mathbf{x} , weight \mathbf{w} and bias b , perceptron outputs:

$$f = \langle \mathbf{w}, \mathbf{x} \rangle + b$$

Cats vs. dogs?



Perceptron

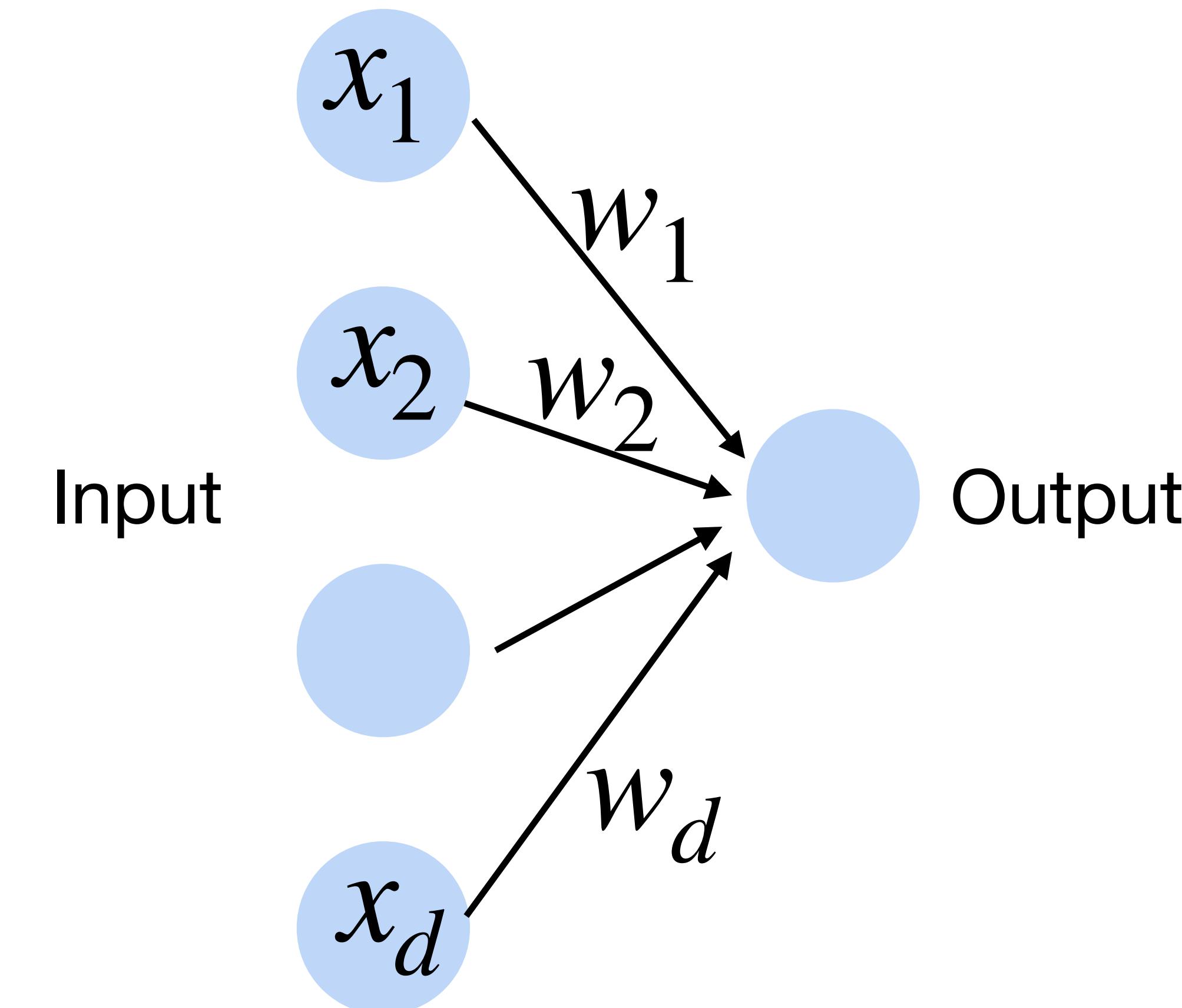
- Given input \mathbf{x} , weight \mathbf{w} and bias b , perceptron outputs:

$$o = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle + b)$$

$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Activation function

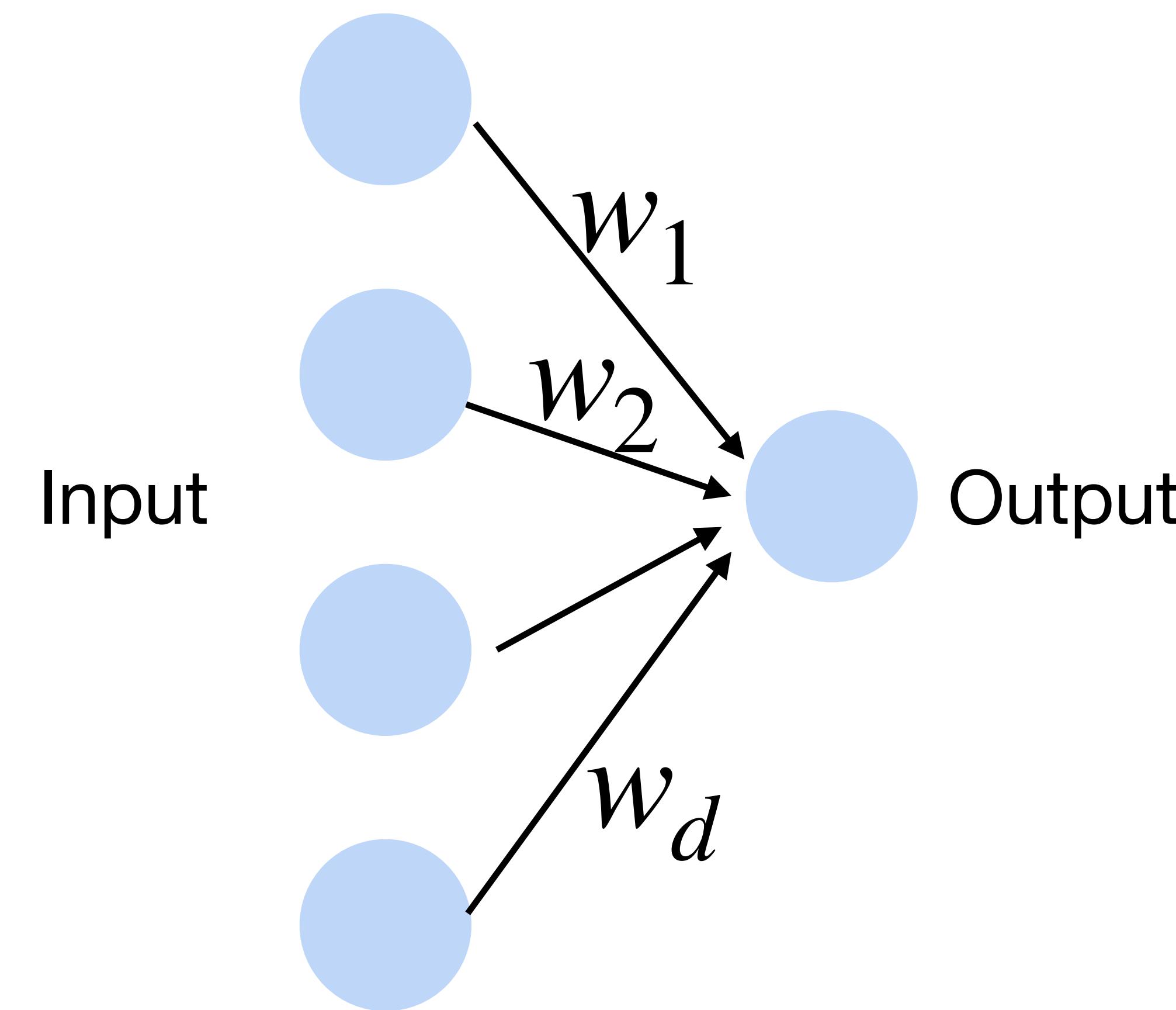
Cats vs. dogs?



Perceptron

- Goal: learn parameters $\mathbf{w} = \{w_1, w_2, \dots, w_d\}$ and b to minimize the classification error

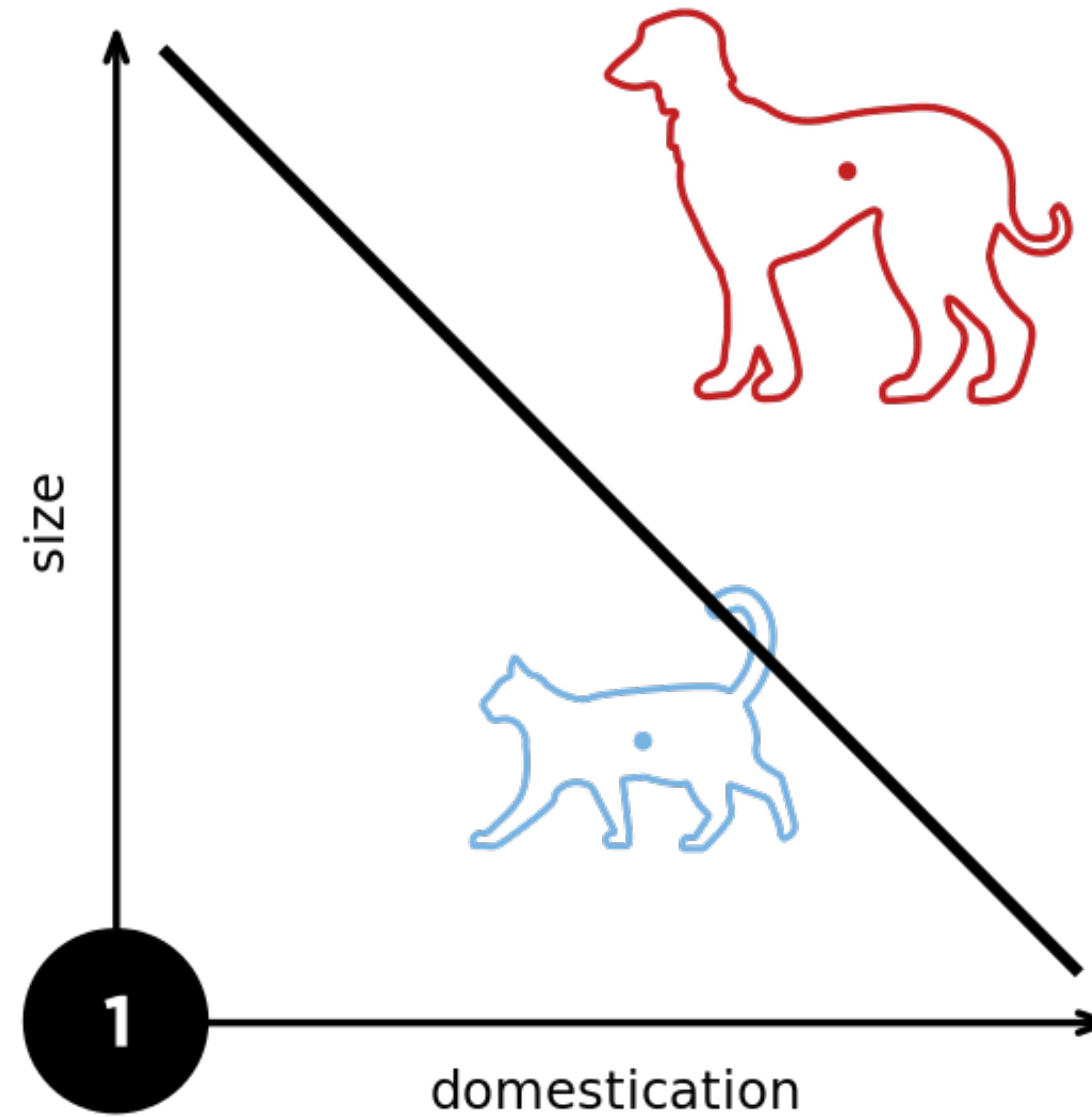
Cats vs. dogs?



Training the Perceptron

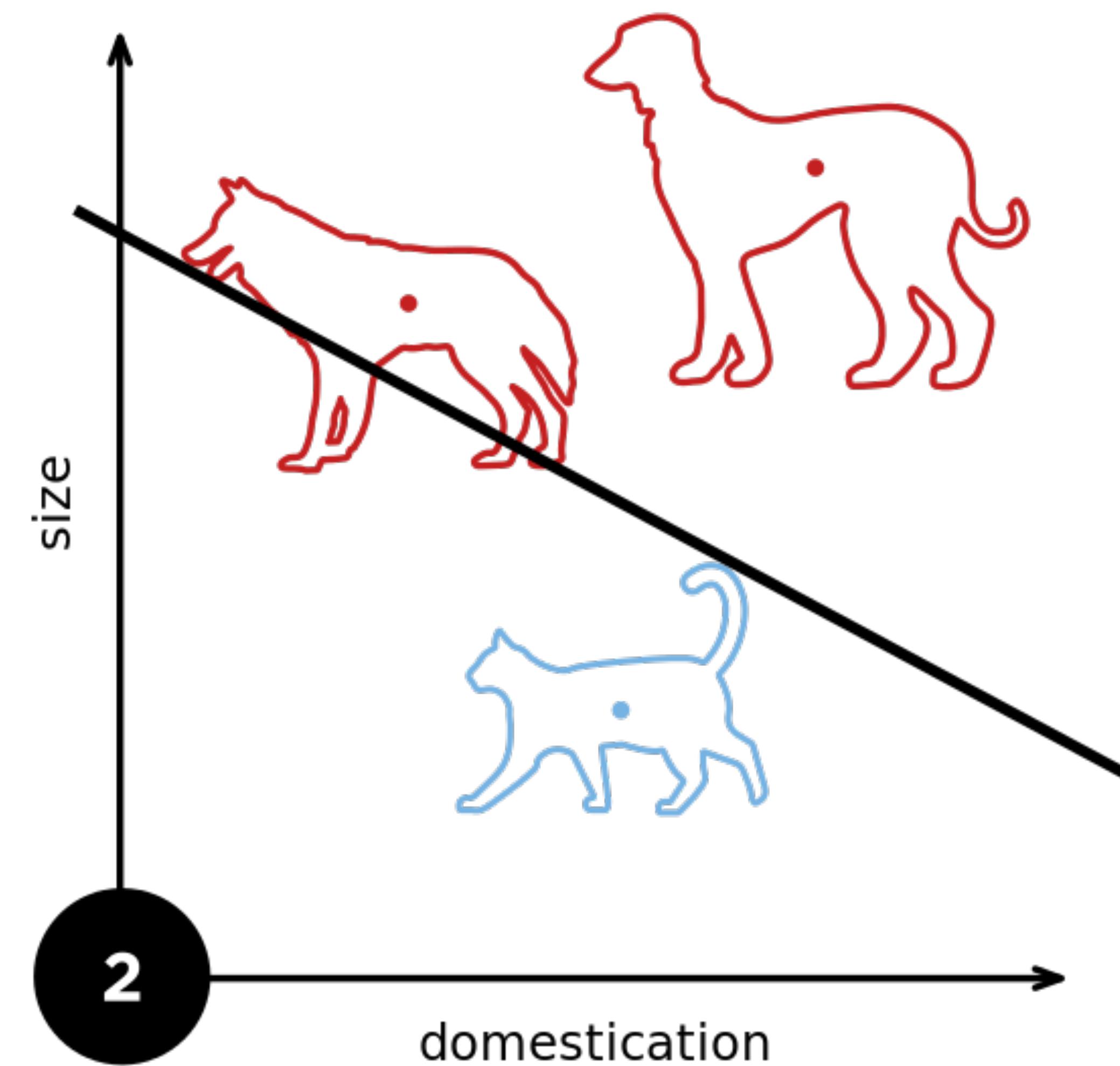
Perceptron Algorithm

Perceptron



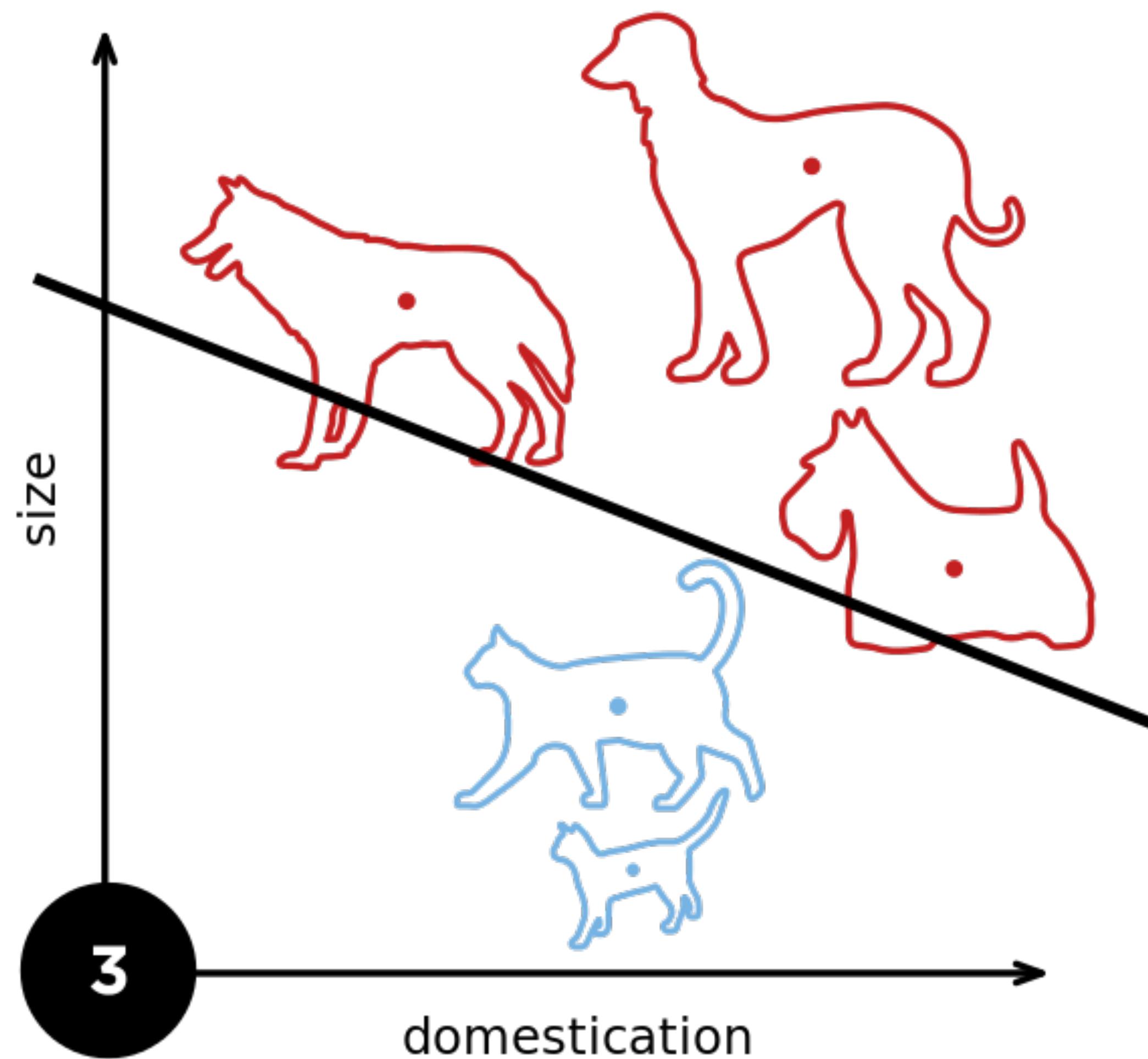
From wikipedia

Perceptron



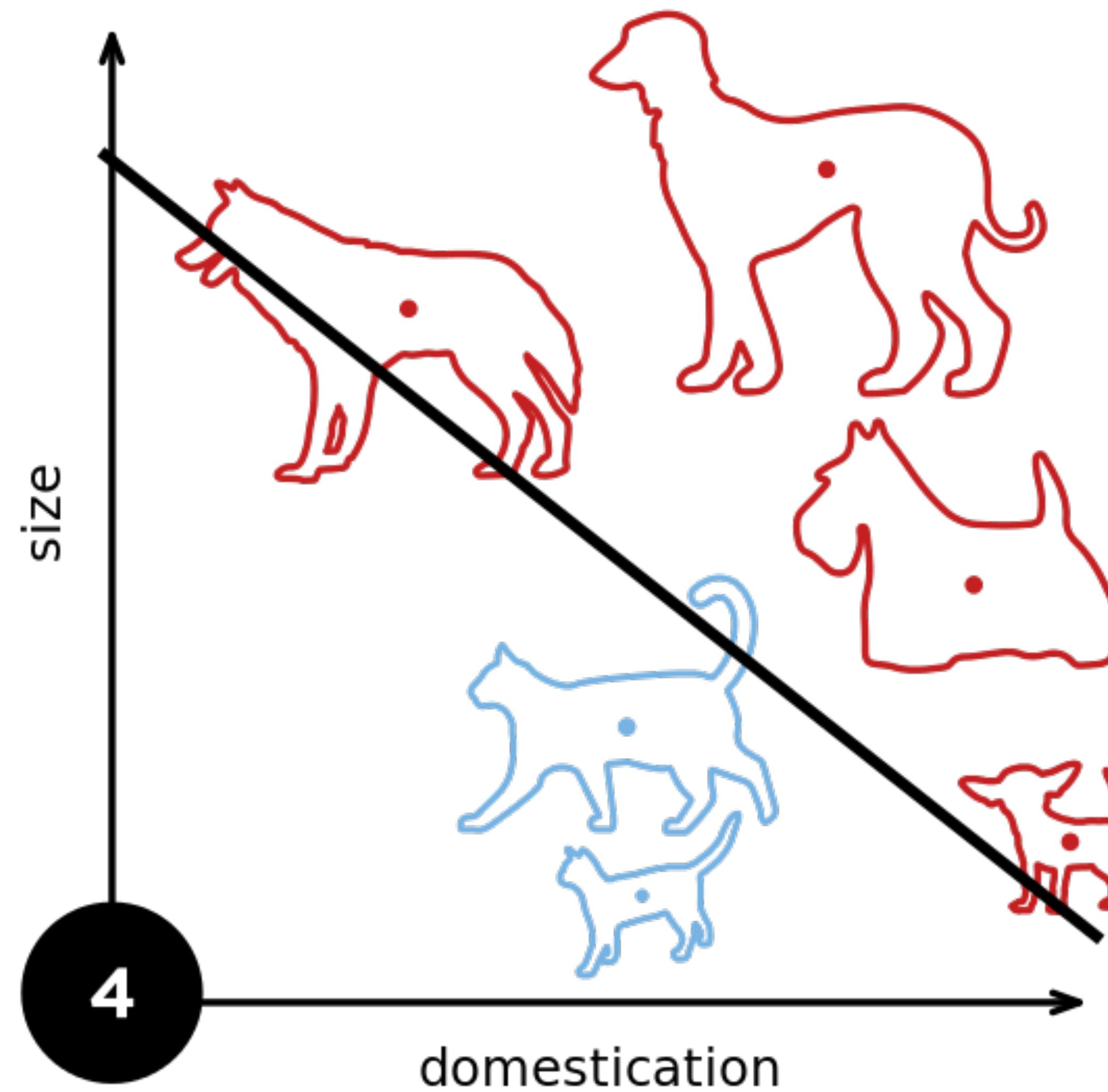
From wikipedia

Perceptron



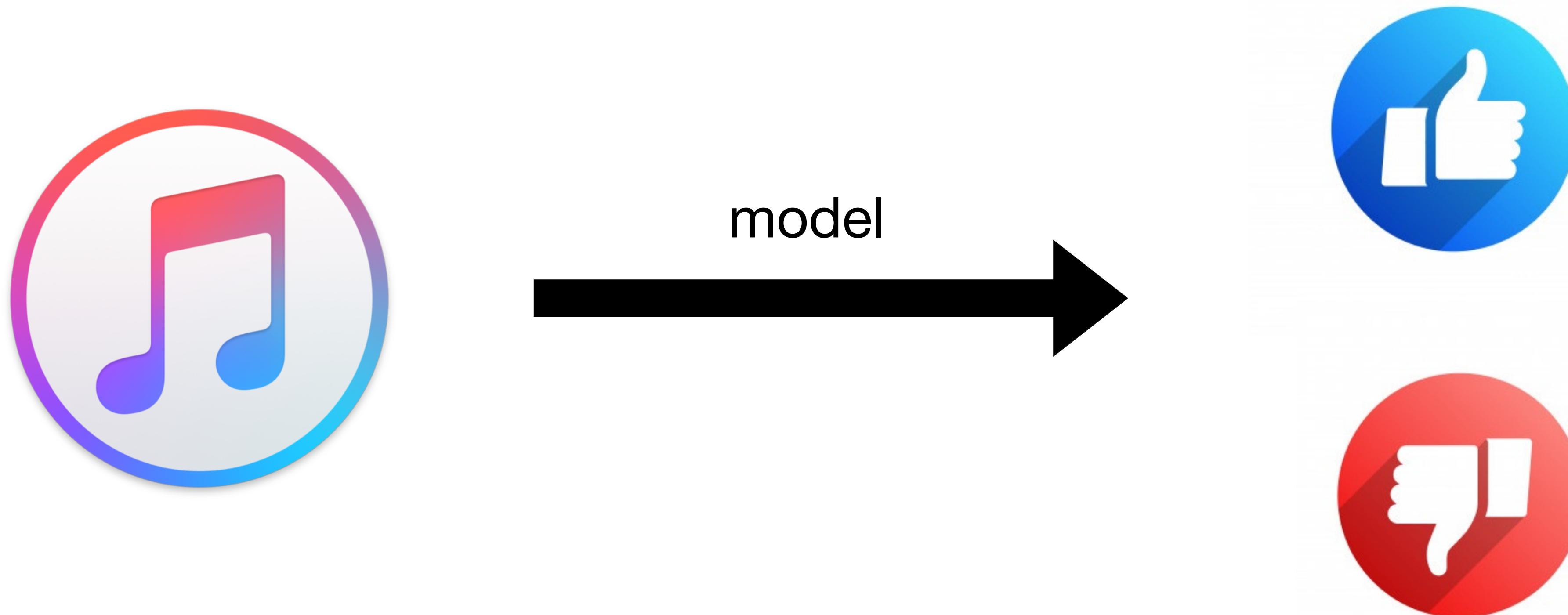
From wikipedia

Perceptron



From wikipedia

Example 2: Predict whether a user likes a song or not

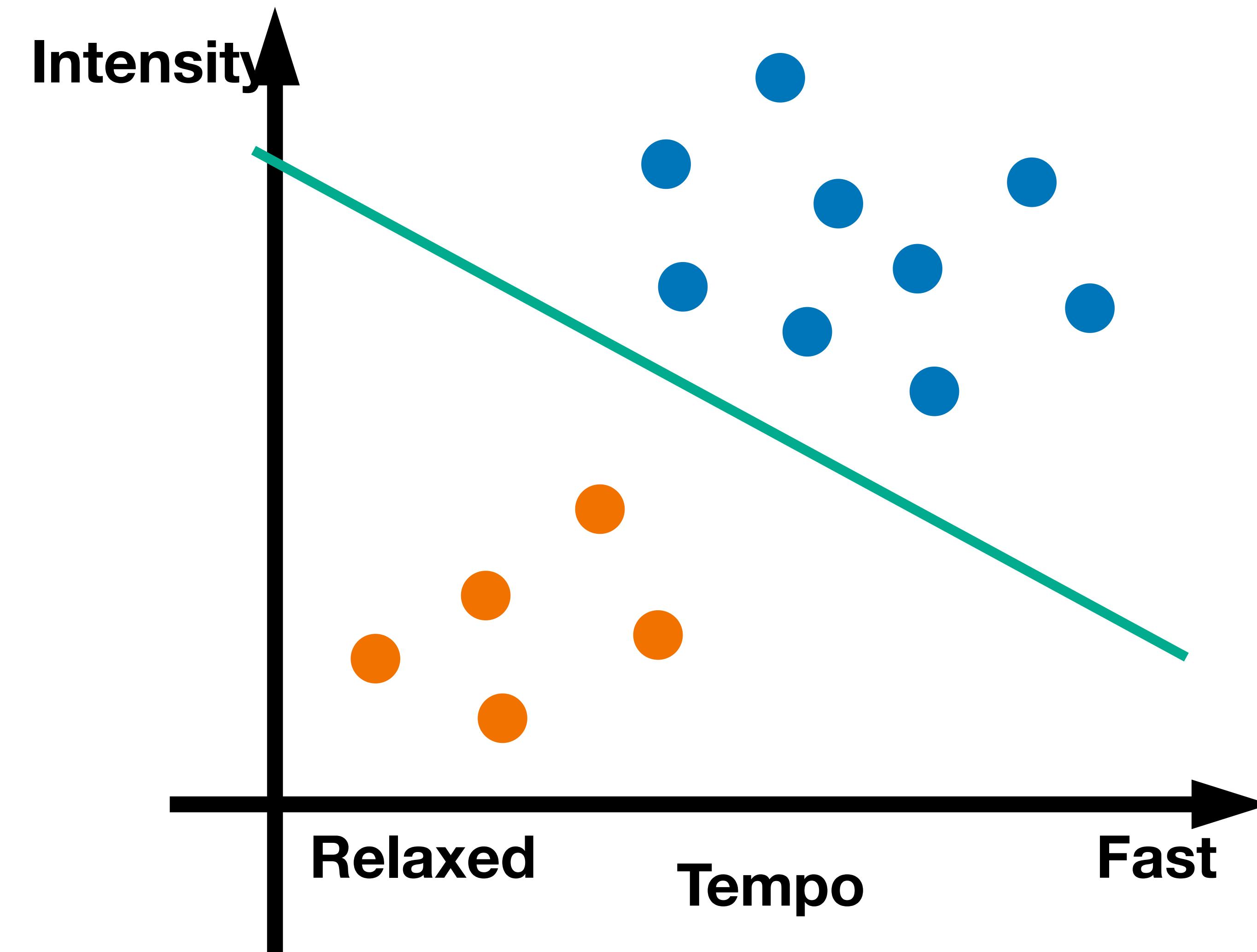


Example 2: Predict whether a user likes a song or not Using Perceptron



User Sharon

- The image consists of two large, solid-colored circles. The top circle is orange and positioned to the left of the word "DisLike". The bottom circle is blue and positioned to the left of the word "Like". Both words are written in a large, bold, black sans-serif font.



Learning AND function using perceptron

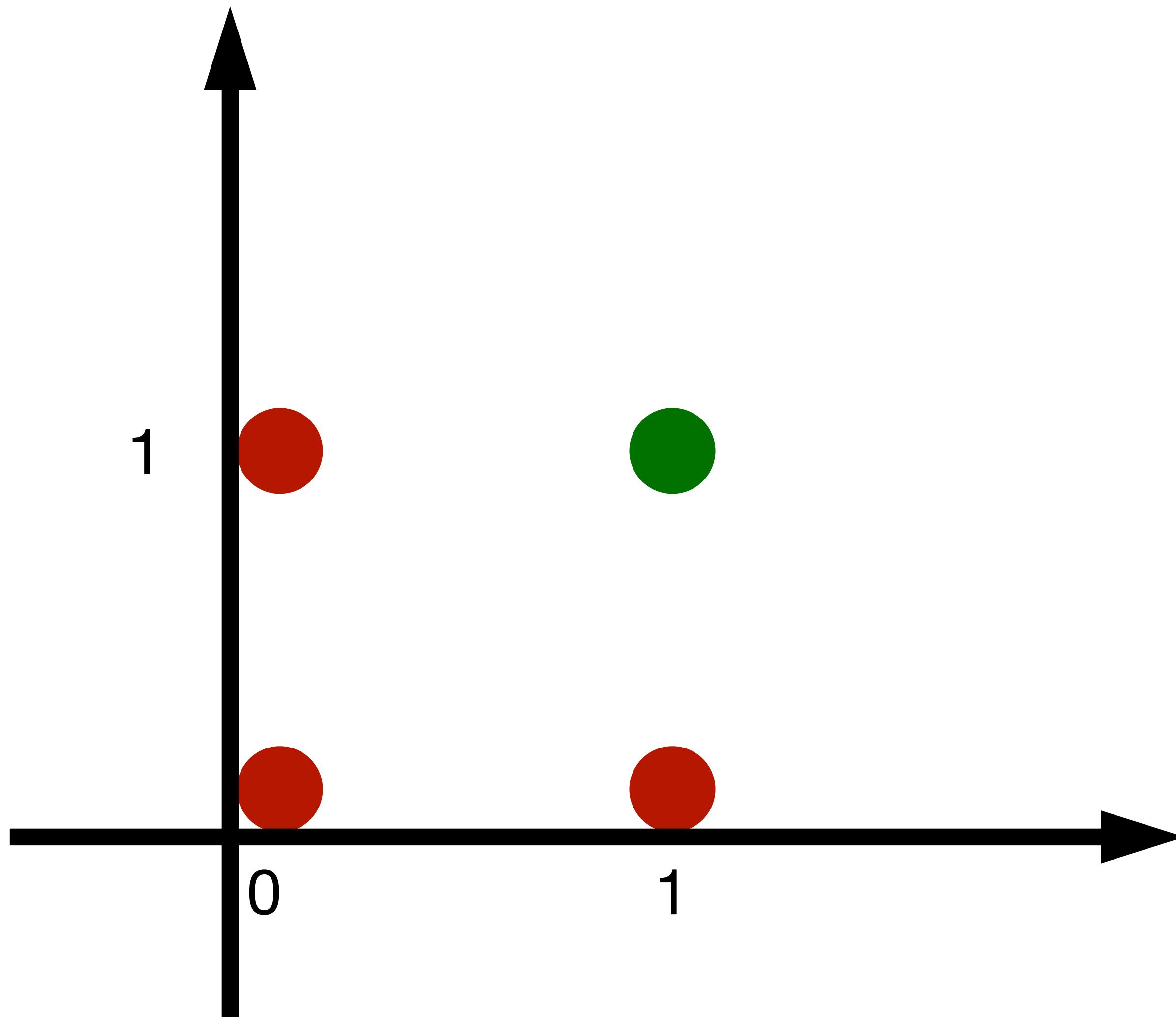
The perceptron can learn an AND function

$$x_1 = 1, x_2 = 1, y = 1$$

$$x_1 = 1, x_2 = 0, y = 0$$

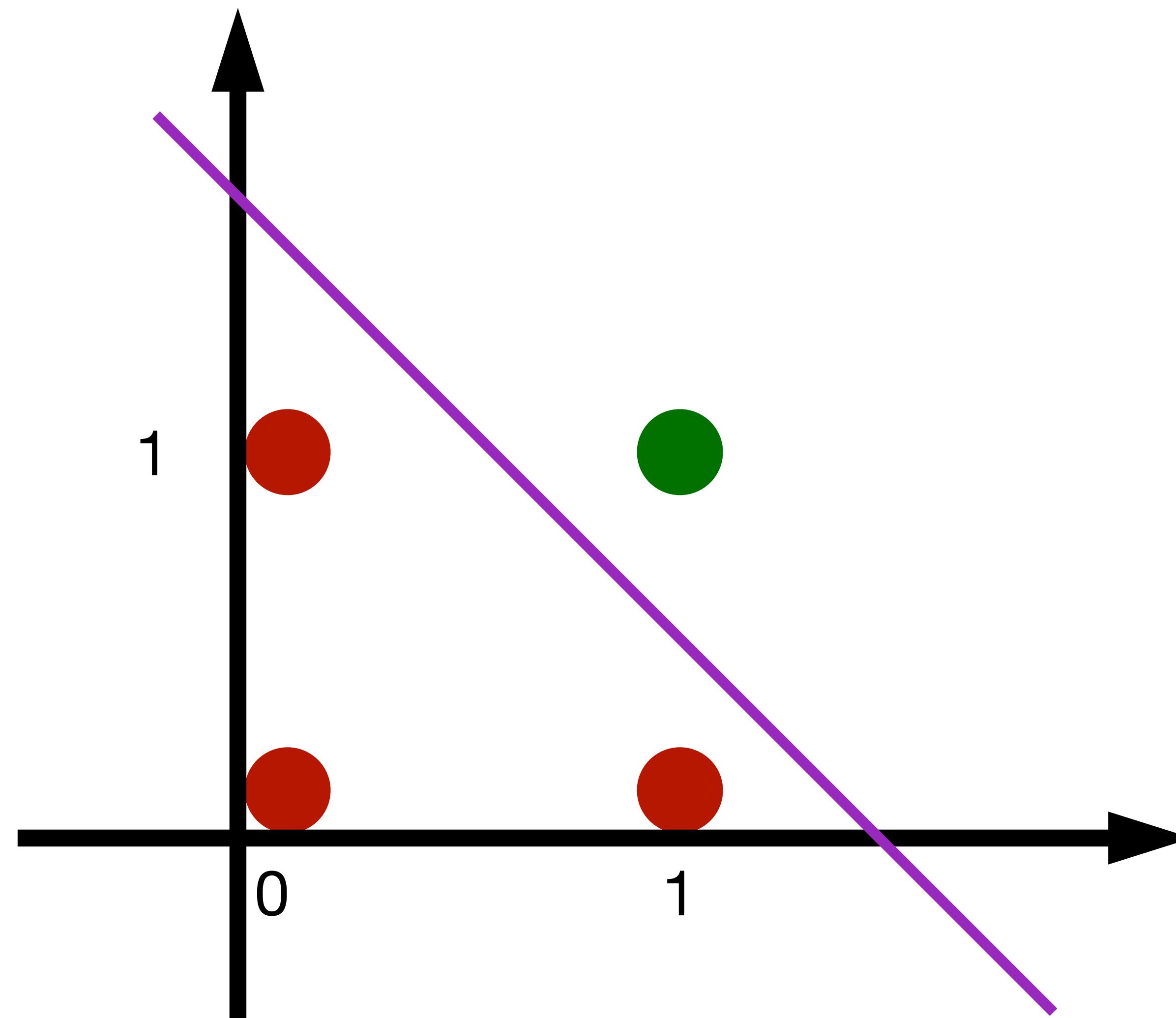
$$x_1 = 0, x_2 = 1, y = 0$$

$$x_1 = 0, x_2 = 0, y = 0$$



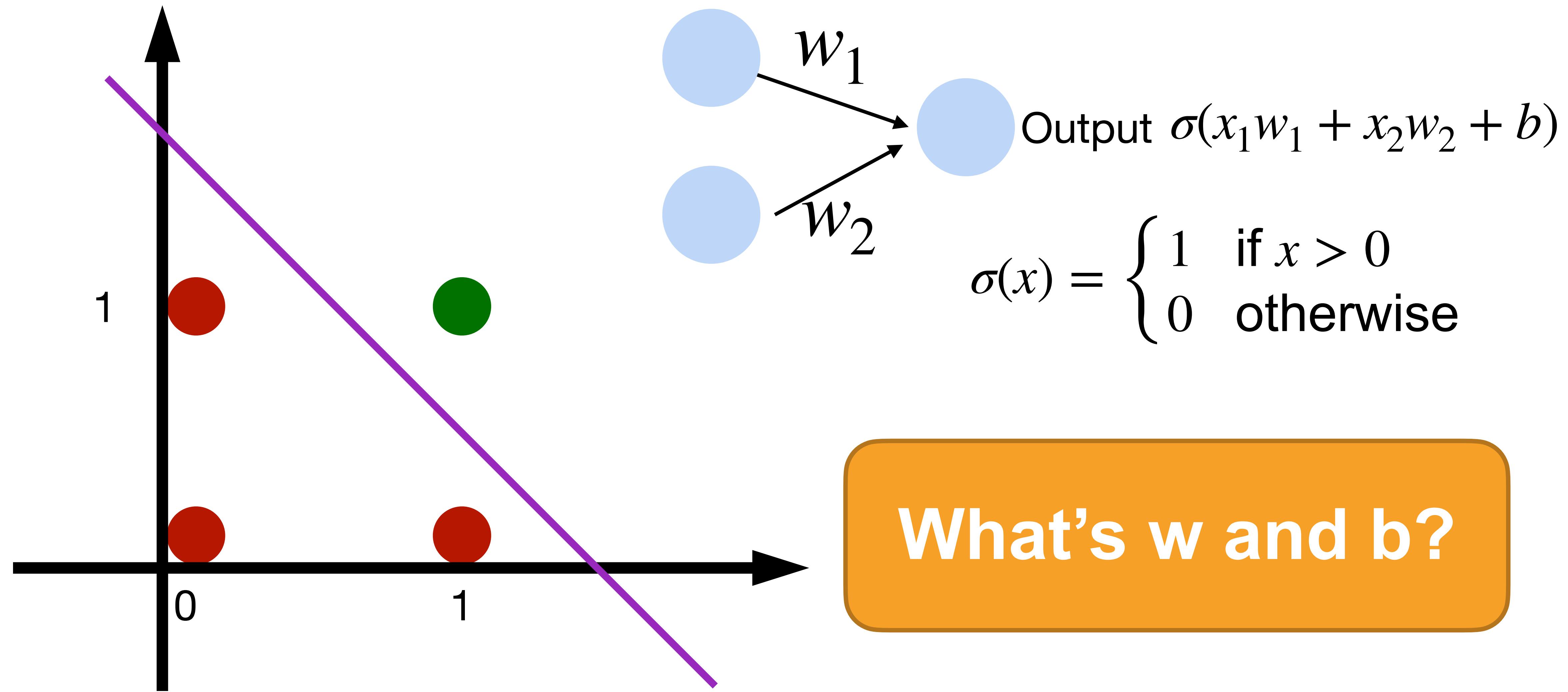
Learning AND function using perceptron

The perceptron can learn an AND function



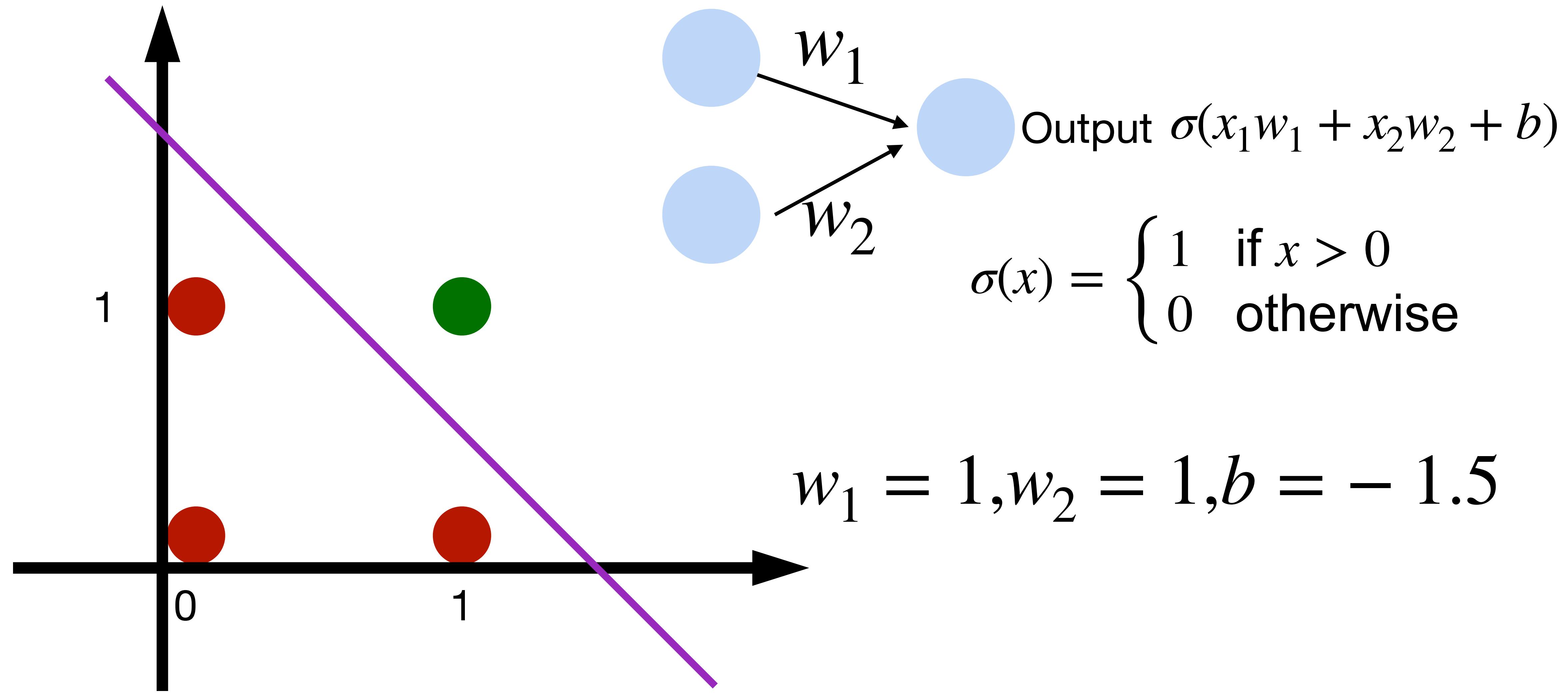
Learning AND function using perceptron

The perceptron can learn an AND function



Learning AND function using perceptron

The perceptron can learn an AND function



Learning OR function using perceptron

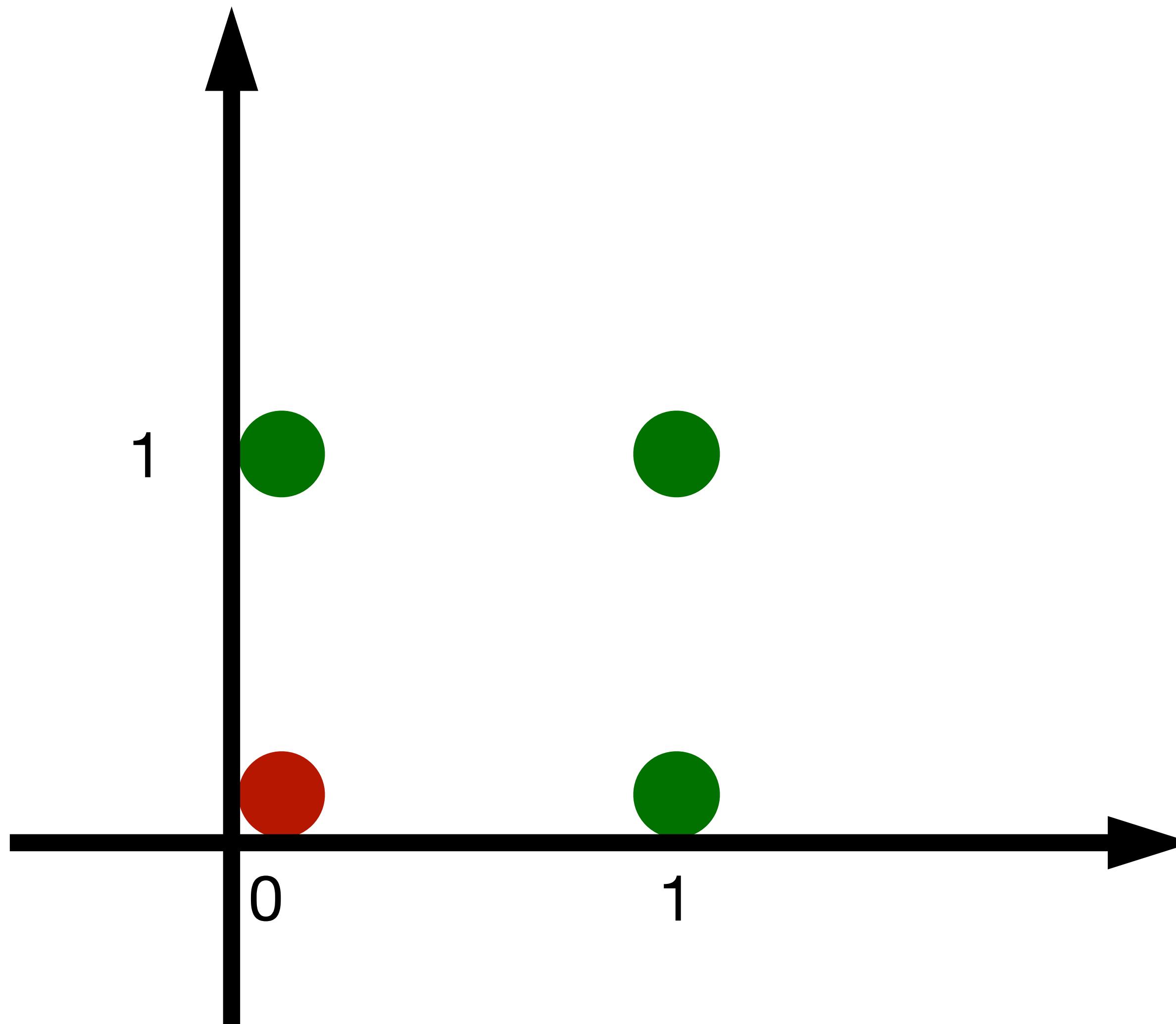
The perceptron can learn an AND function

$$x_1 = 1, x_2 = 1, y = 1$$

$$x_1 = 1, x_2 = 0, y = 1$$

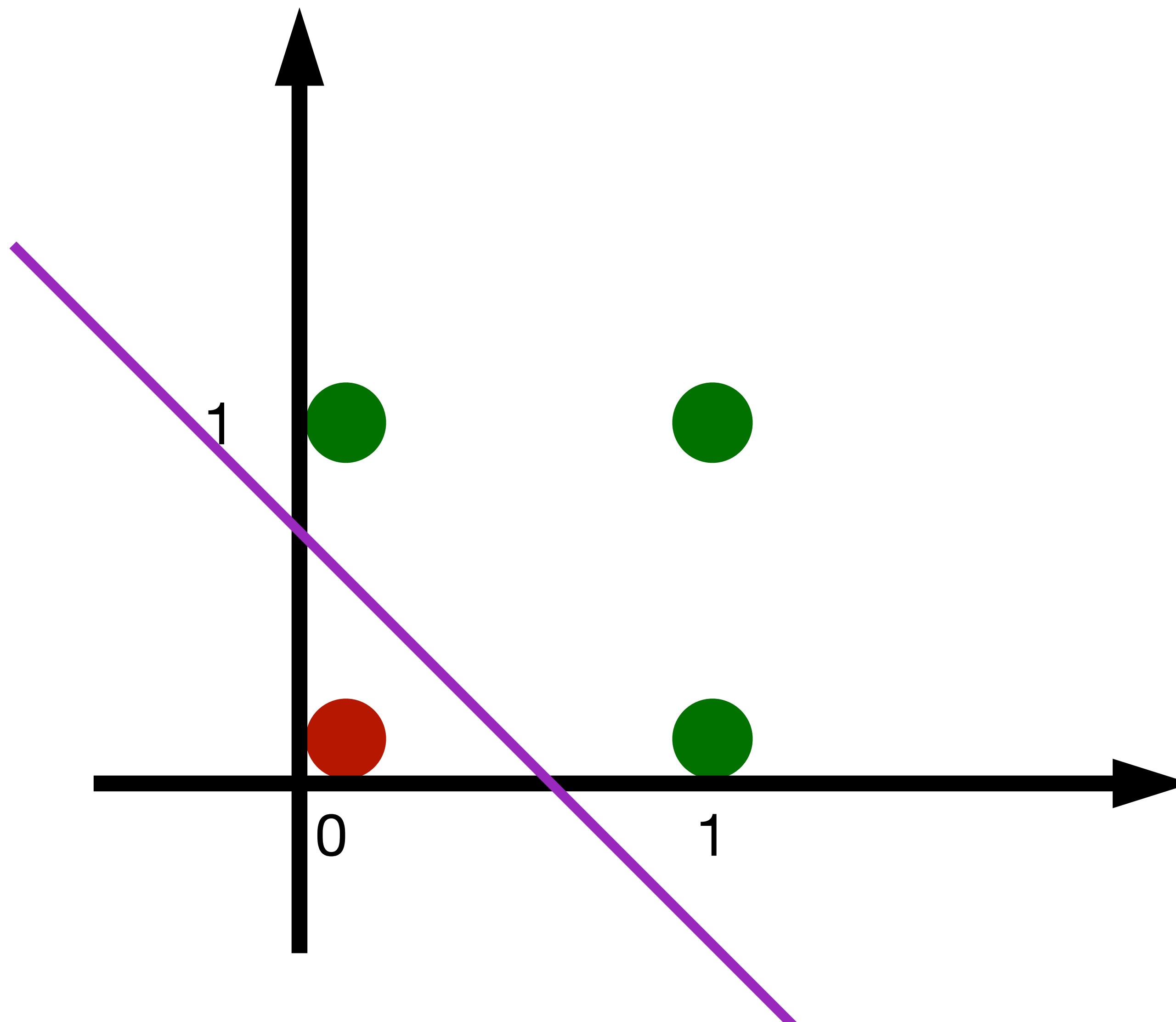
$$x_1 = 0, x_2 = 1, y = 1$$

$$x_1 = 0, x_2 = 0, y = 0$$



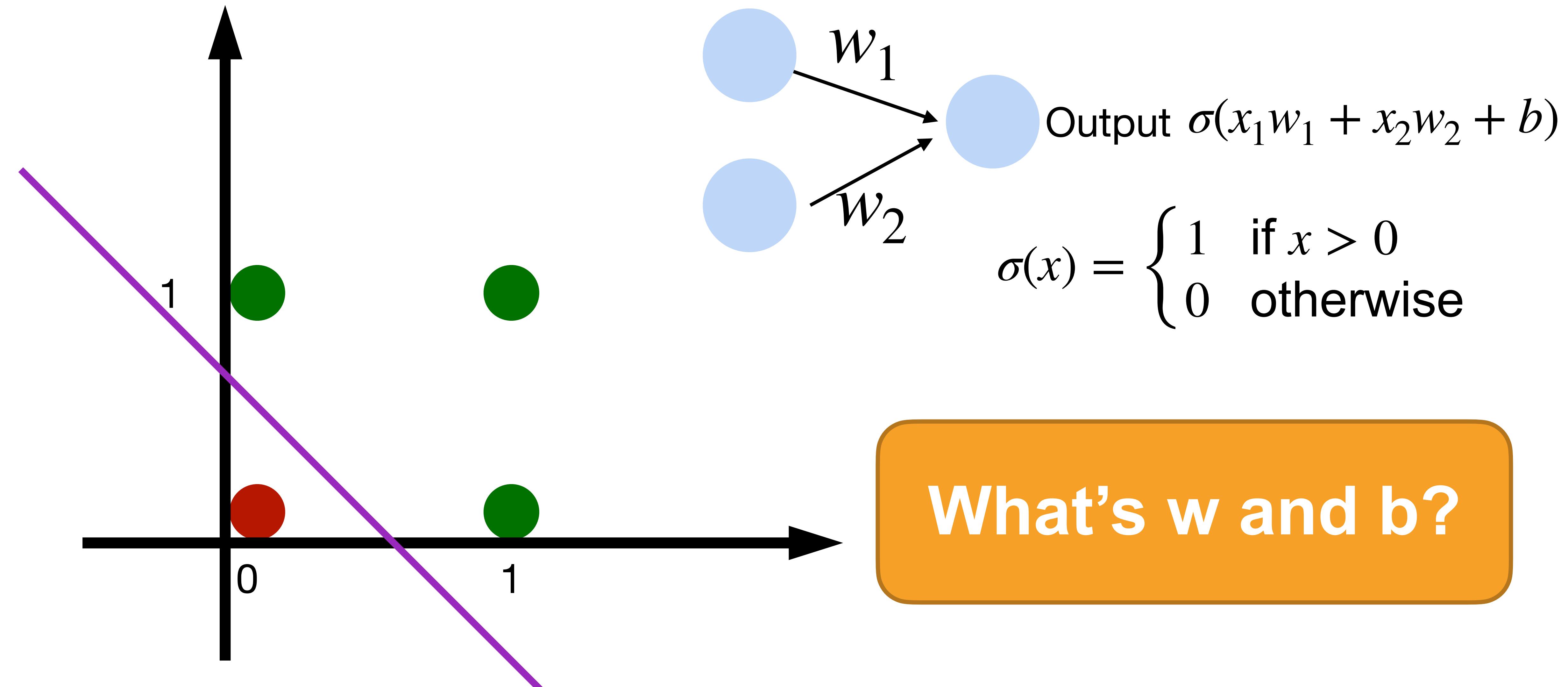
Learning OR function using perceptron

The perceptron can learn an AND function



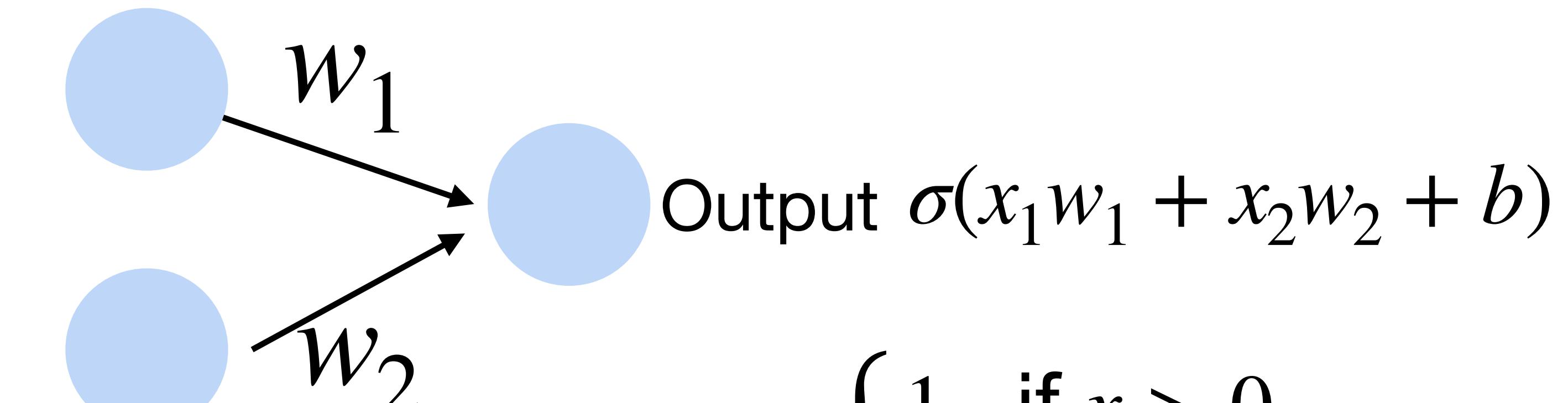
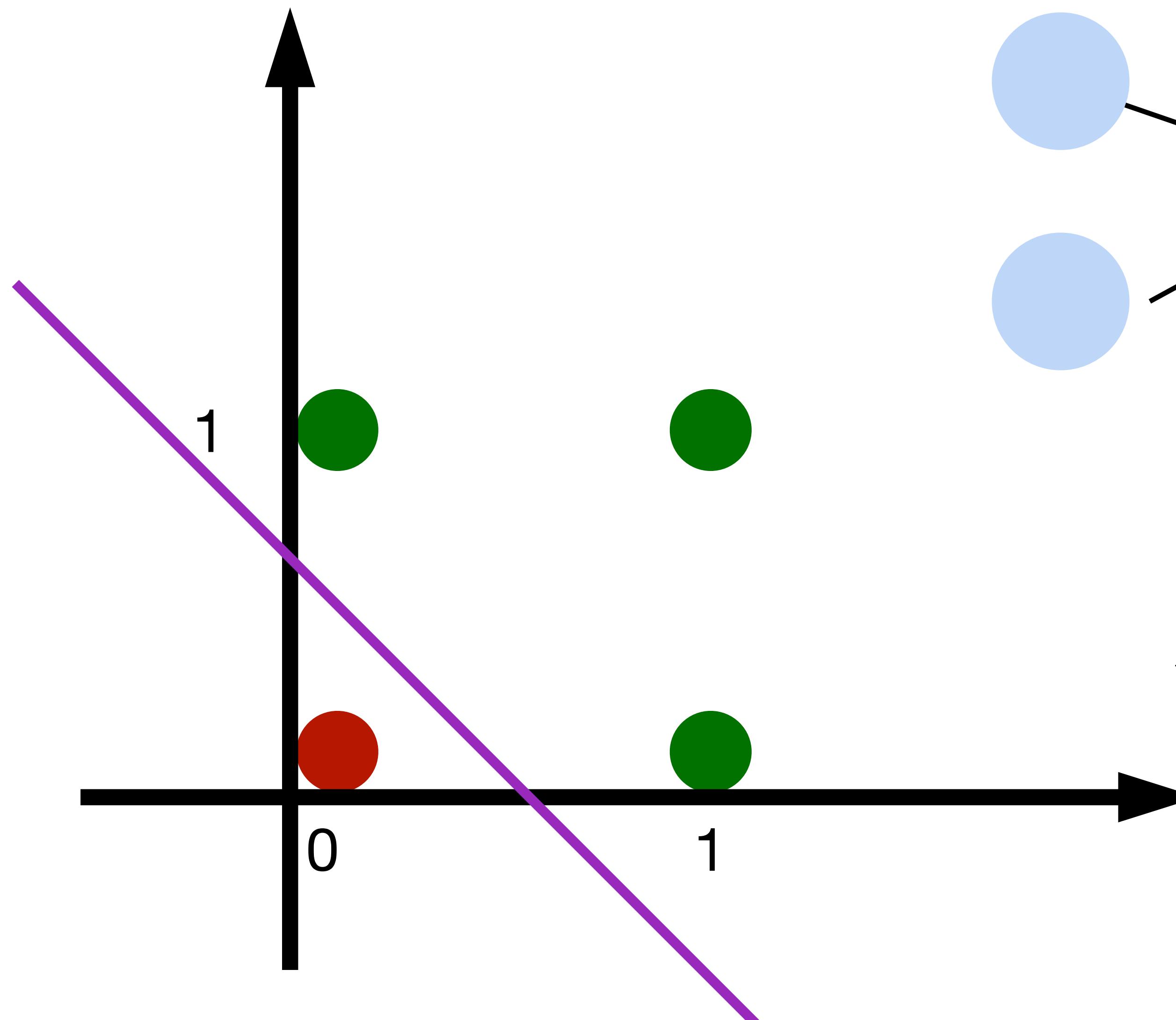
Learning OR function using perceptron

The perceptron can learn an AND function



Learning OR function using perceptron

The perceptron can learn an AND function

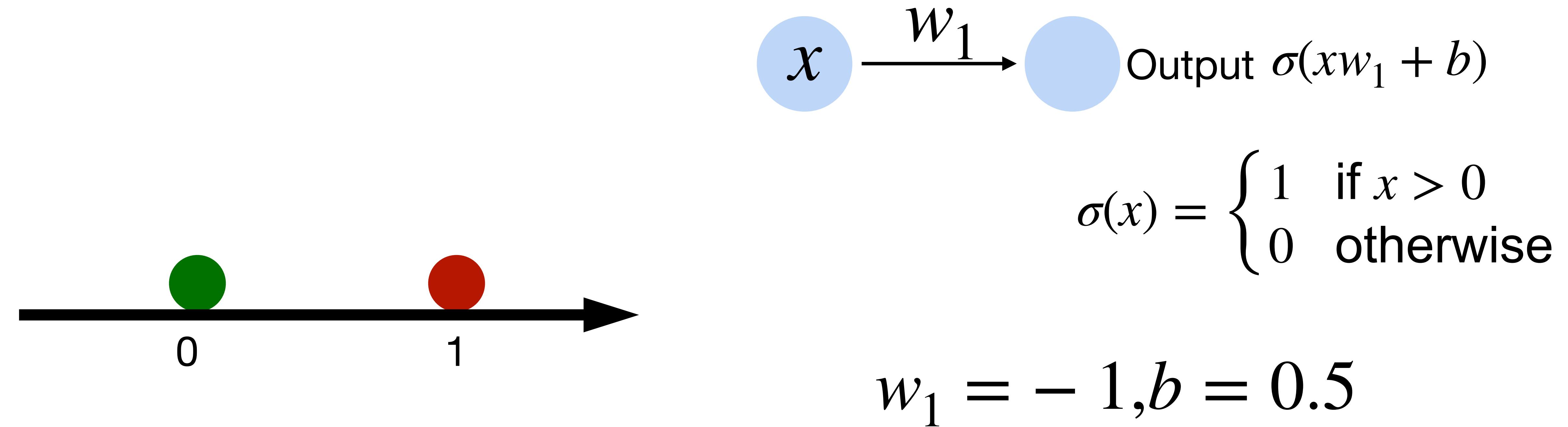


$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$w_1 = 1, w_2 = 1, b = -0.5$$

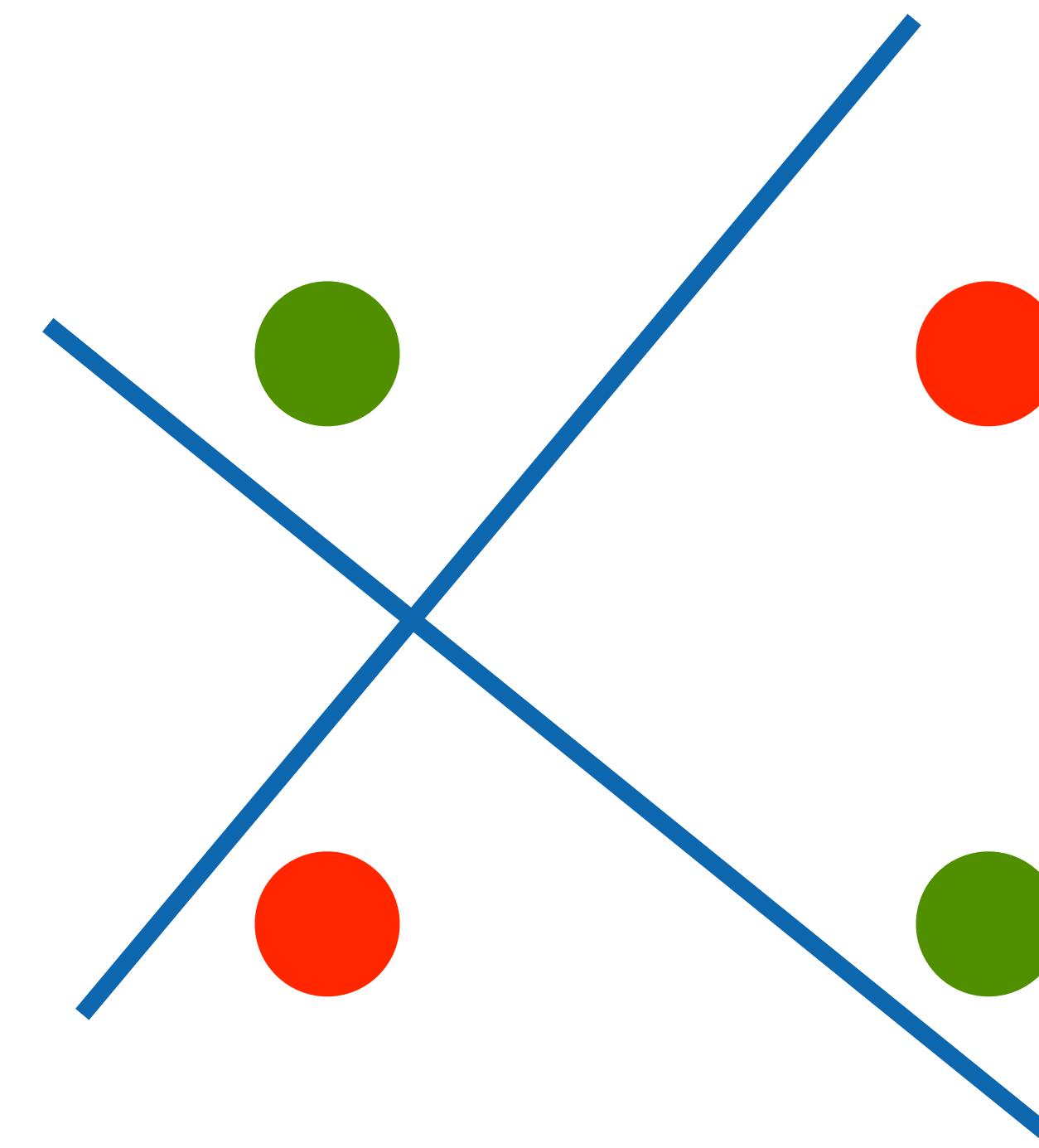
Learning NOT function using perceptron

The perceptron can learn NOT function (single input)



XOR Problem (Minsky & Papert, 1969)

The perceptron cannot learn an XOR function
(neurons can only generate linear separators)

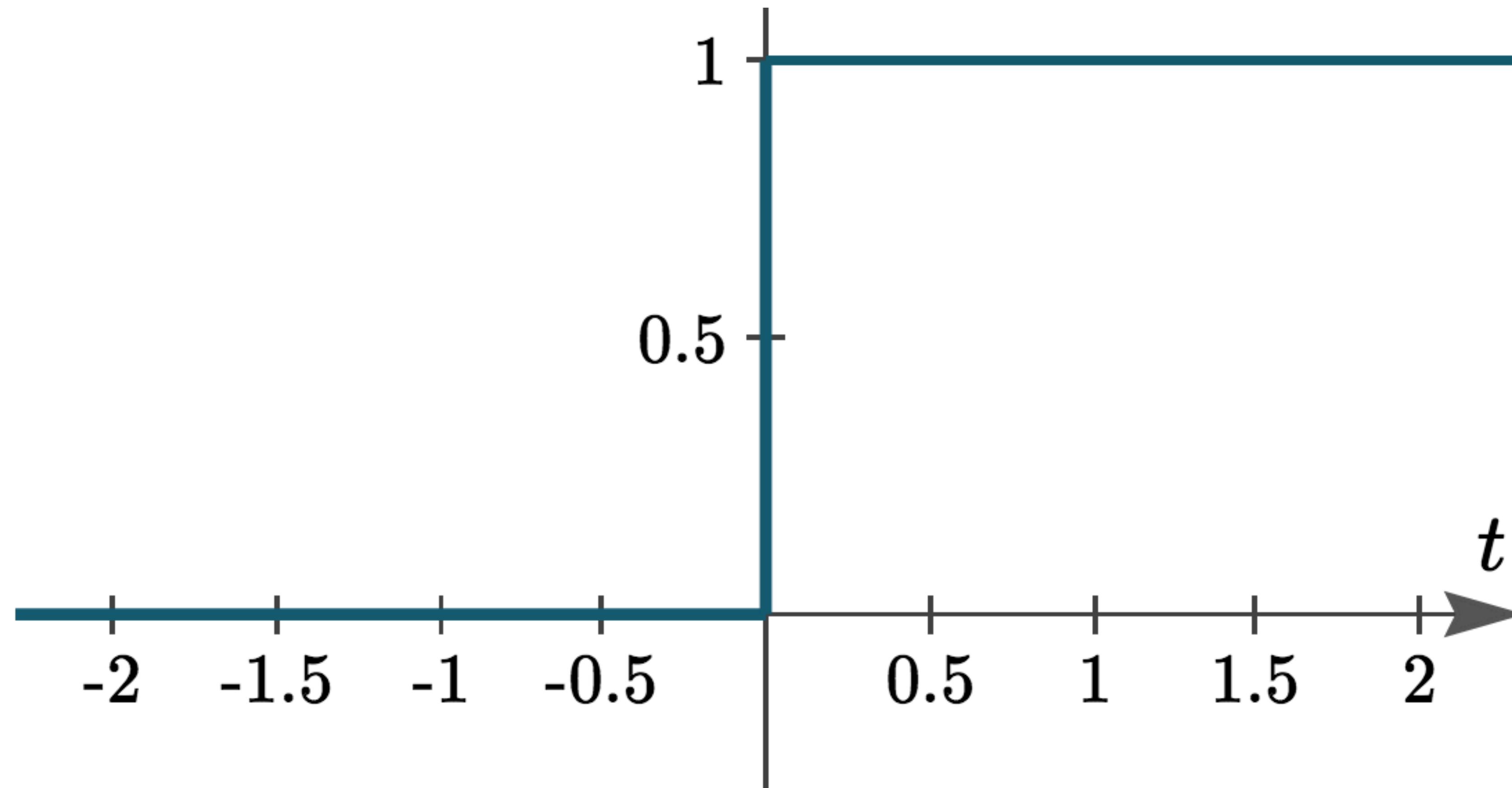


This contributed to the first AI winter

Step Function activation

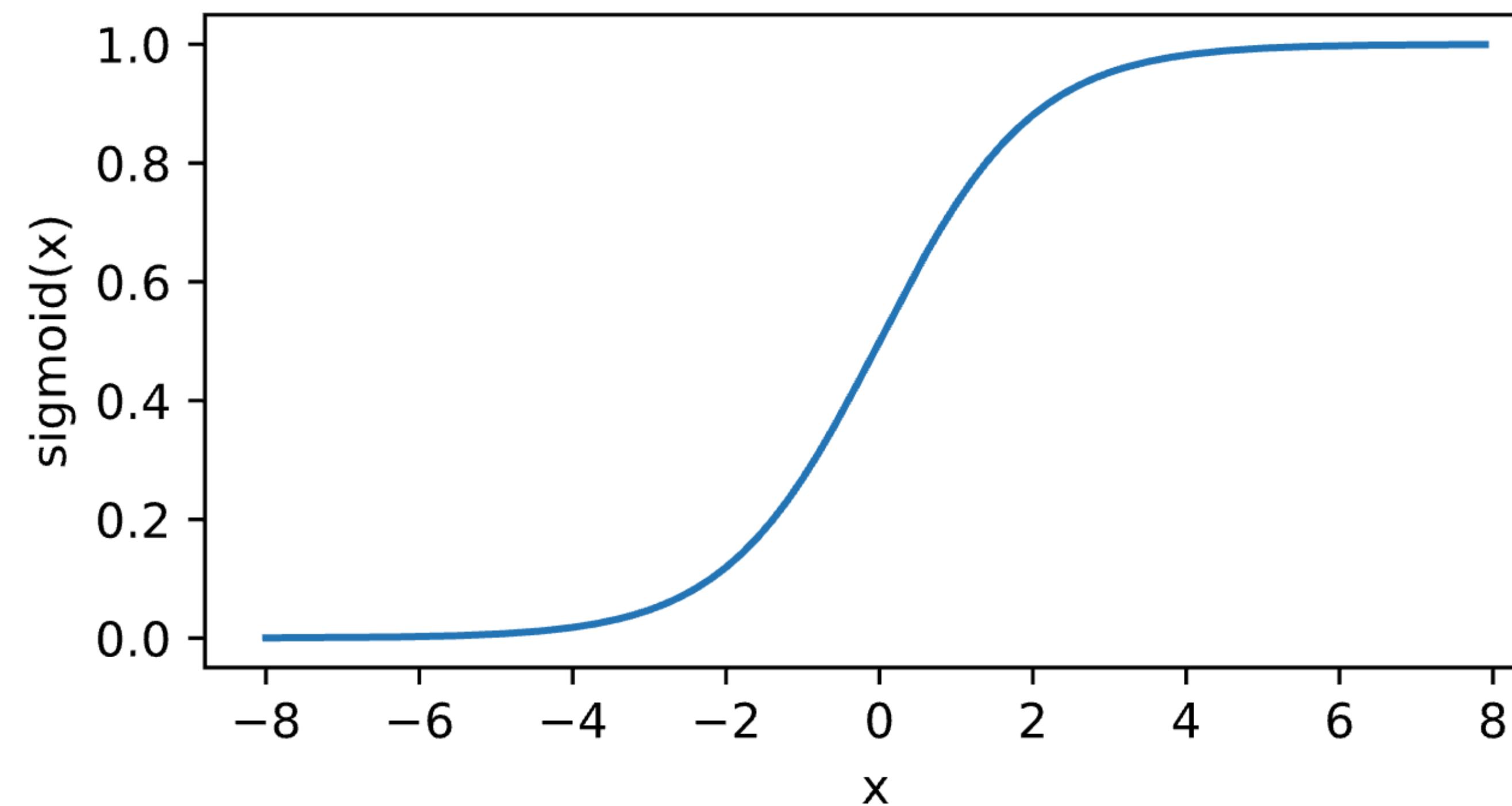
Step function is discontinuous, which cannot be used for gradient descent

$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$



Sigmoid/Logistic Activation

Map input into $[0, 1]$, a **soft** version of $\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

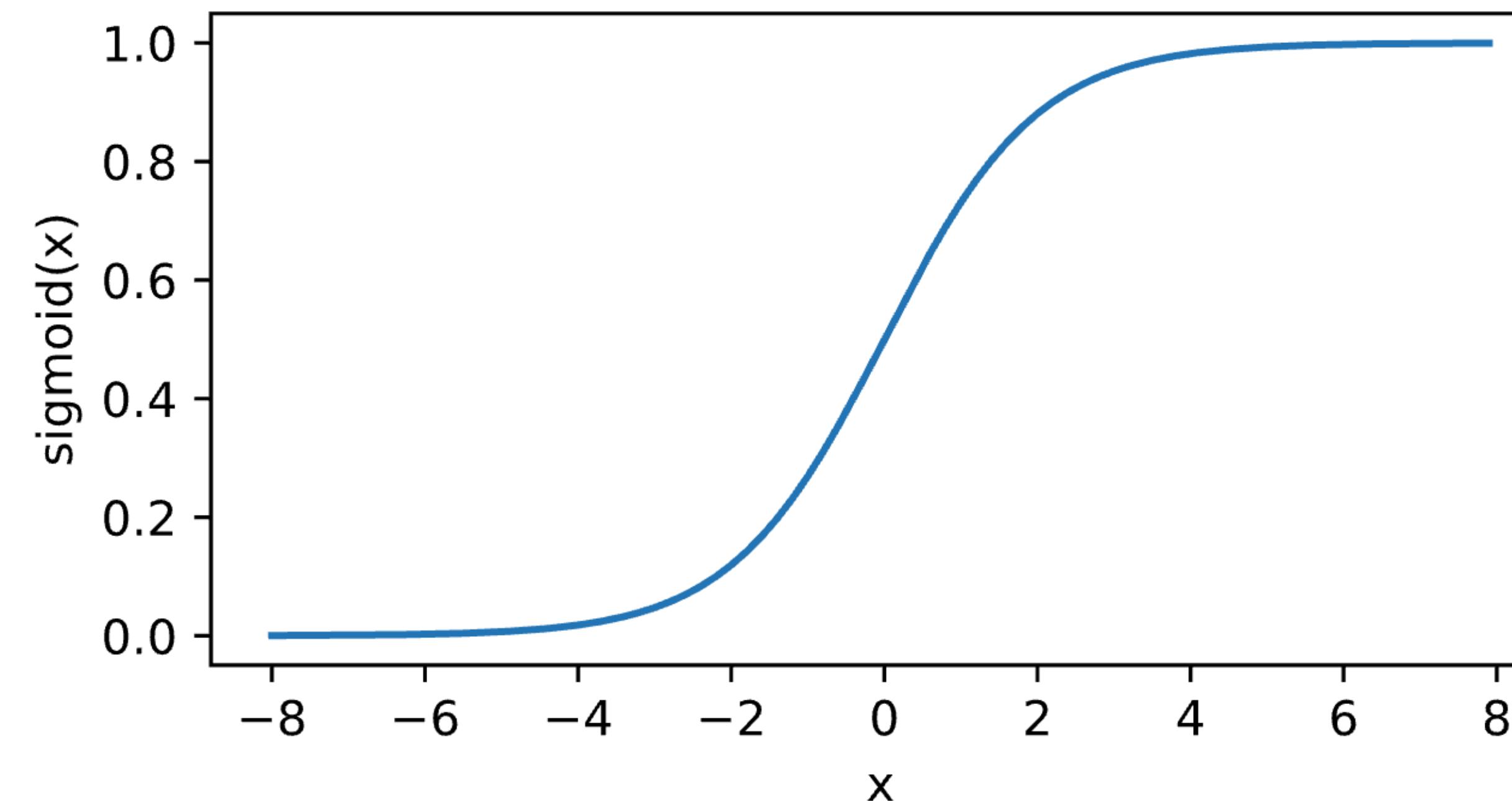
$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$


Logistic regression

$\mathbf{x} \in \mathbb{R}^d, y = \{-1, +1\}$

$$p(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

$$p(y = -1 | \mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x})}$$



Logistic regression

Given: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

Training: maximize likelihood estimate (on the conditional probability)

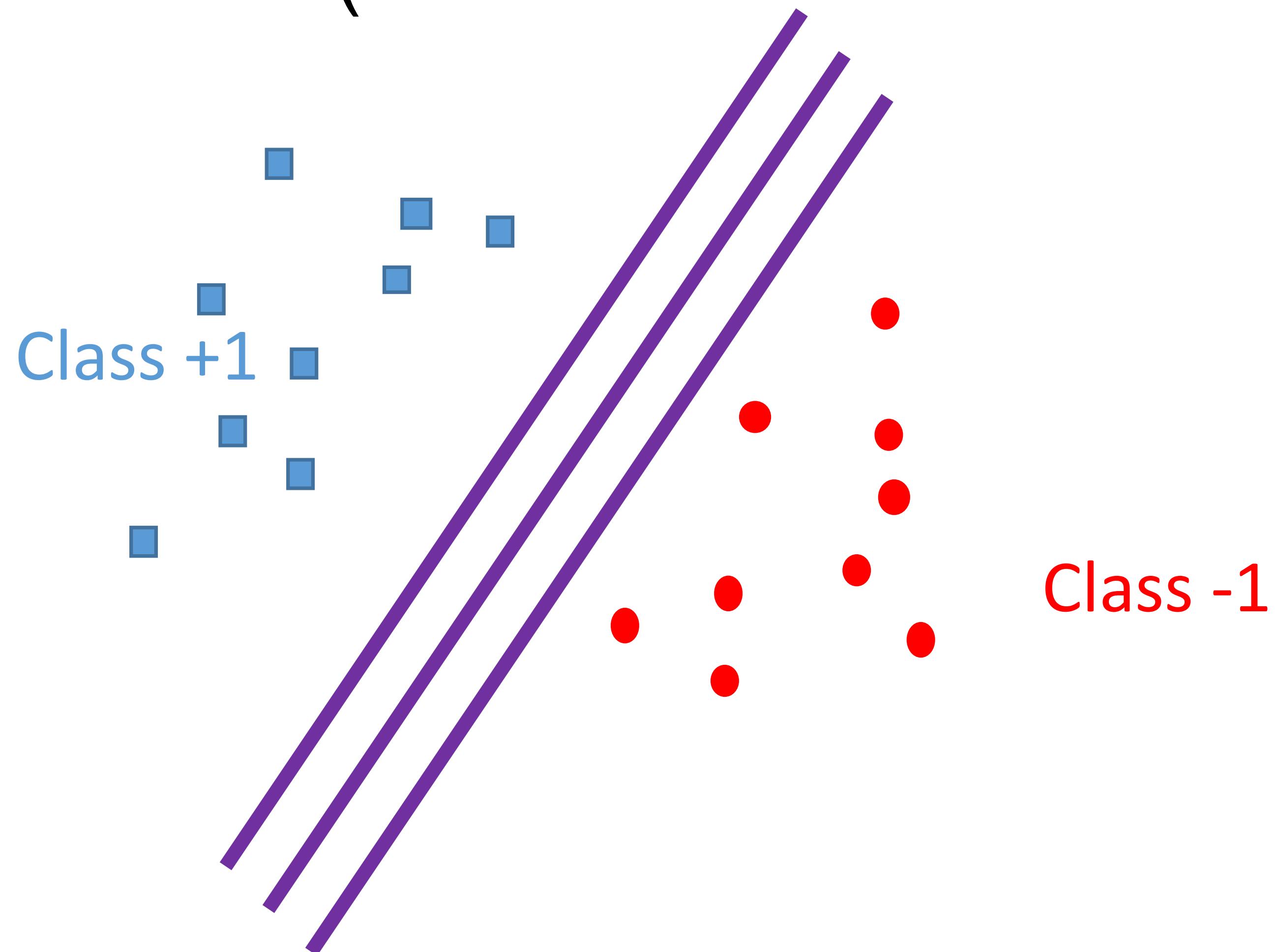
$$\max_{\mathbf{w}} \sum_i \log \frac{1}{1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)}$$

Logistic regression

Given: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

Training: maximize likelihood estimate (on the conditional probability)

When training data is linearly separable, many solutions



Logistic regression

Given: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

Training: maximum A posteriori (MAP)

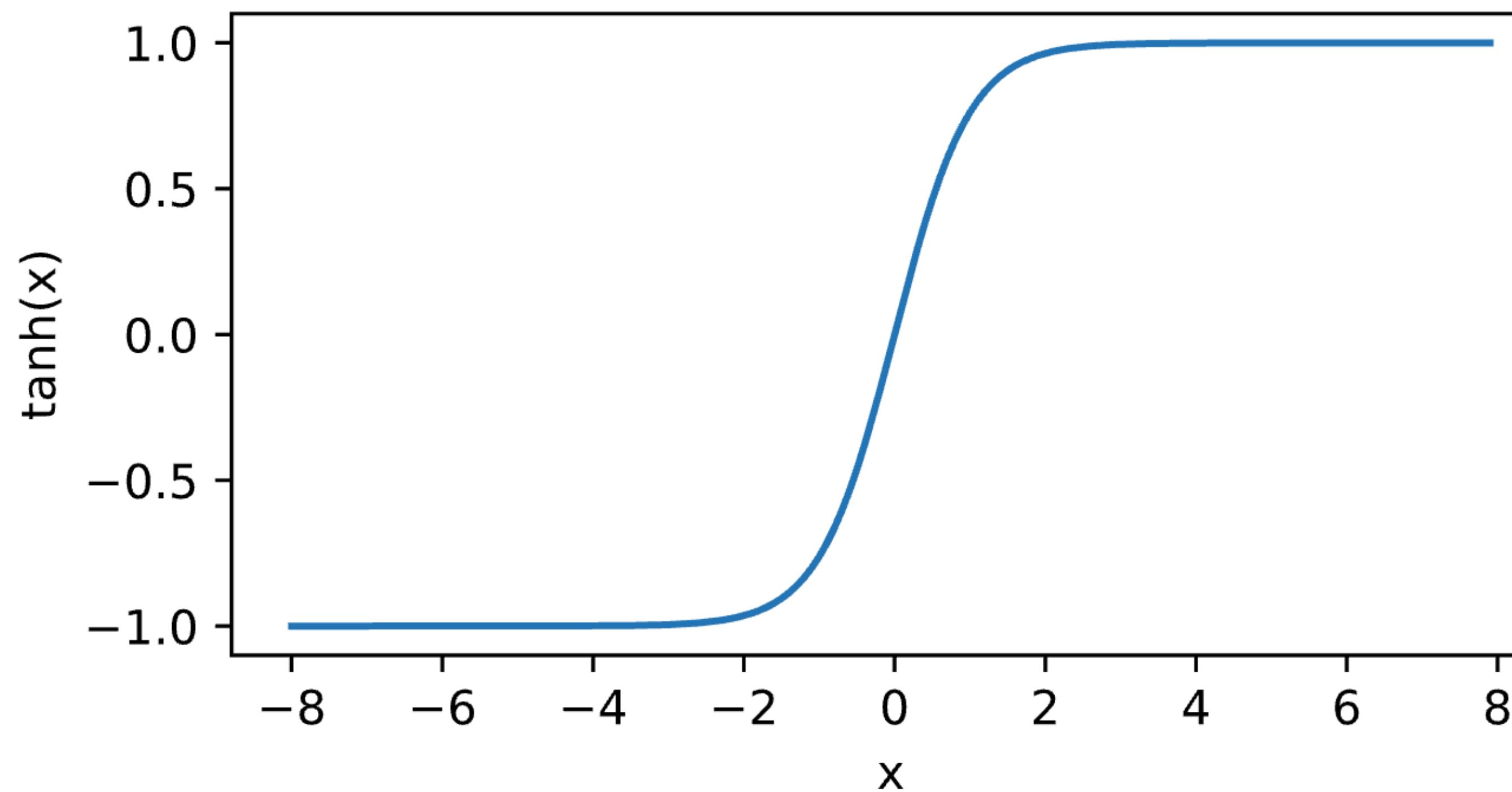
$$\min_{\mathbf{w}} \sum_i -\log \frac{1}{1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Convex optimization
- Solve via (stochastic) gradient descent

Tanh Activation

Map inputs into (-1, 1)

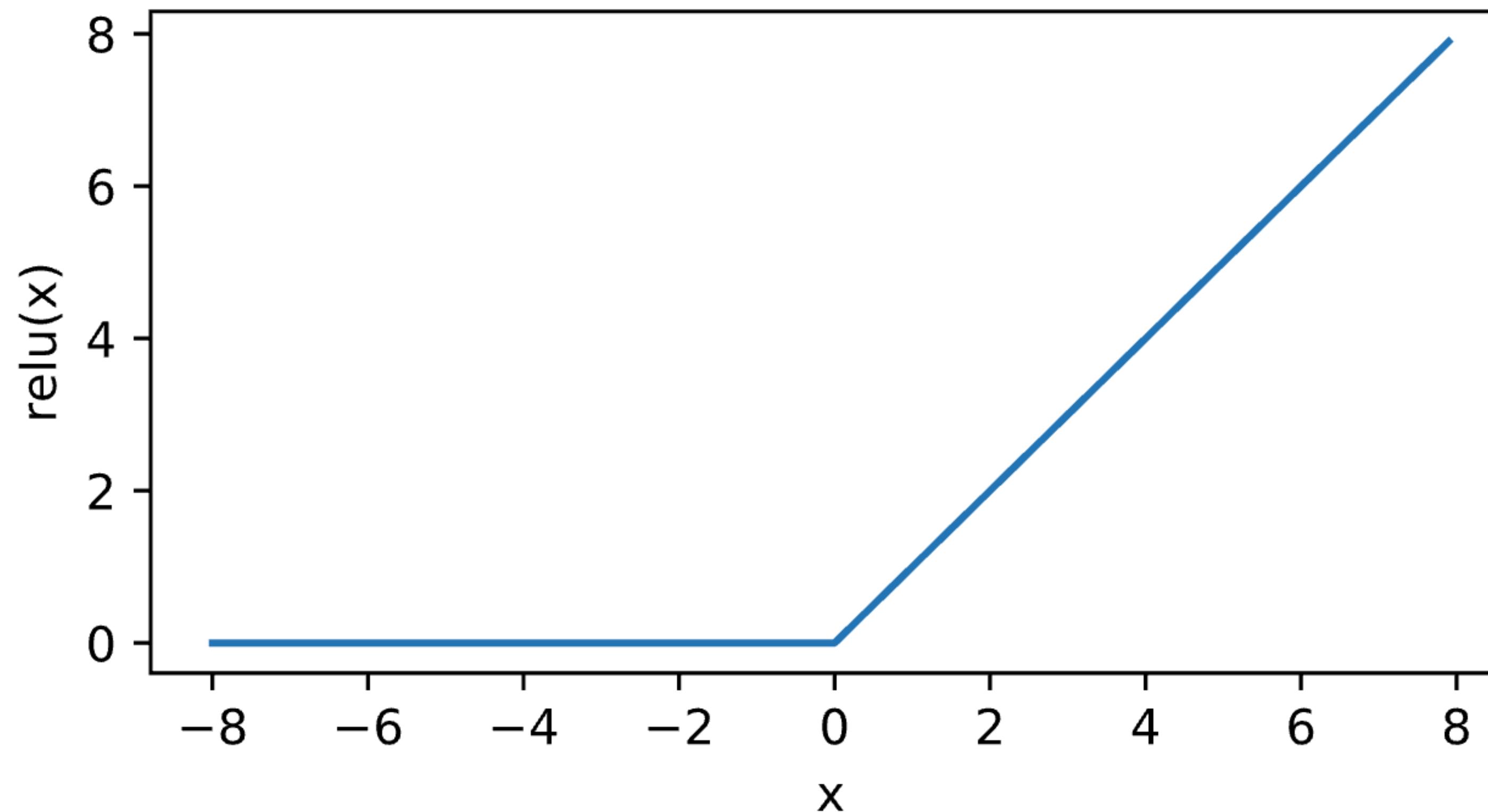
$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}$$



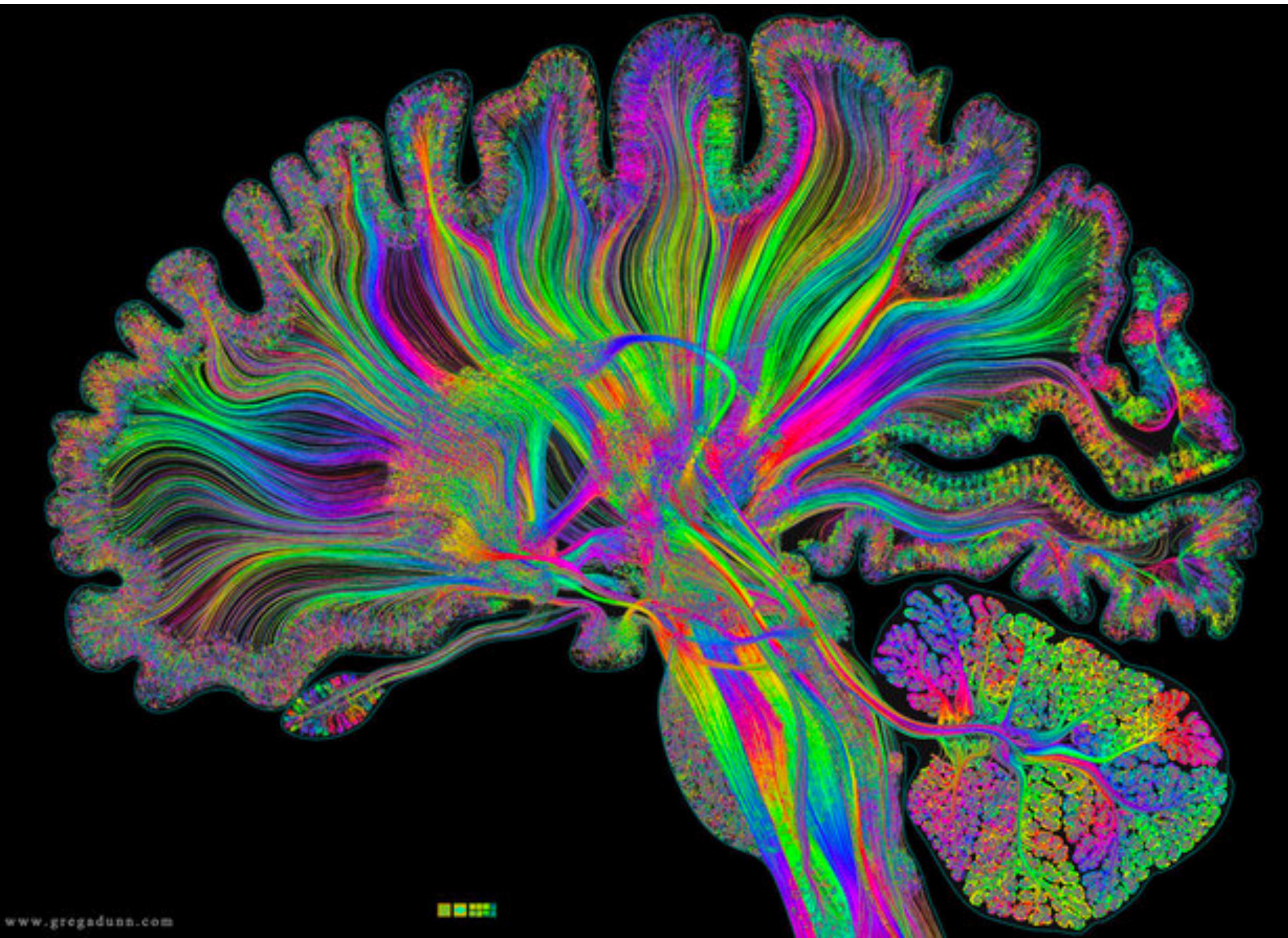
ReLU Activation

ReLU: rectified linear unit (commonly used in modern neural networks)

$$\text{ReLU}(x) = \max(x, 0)$$

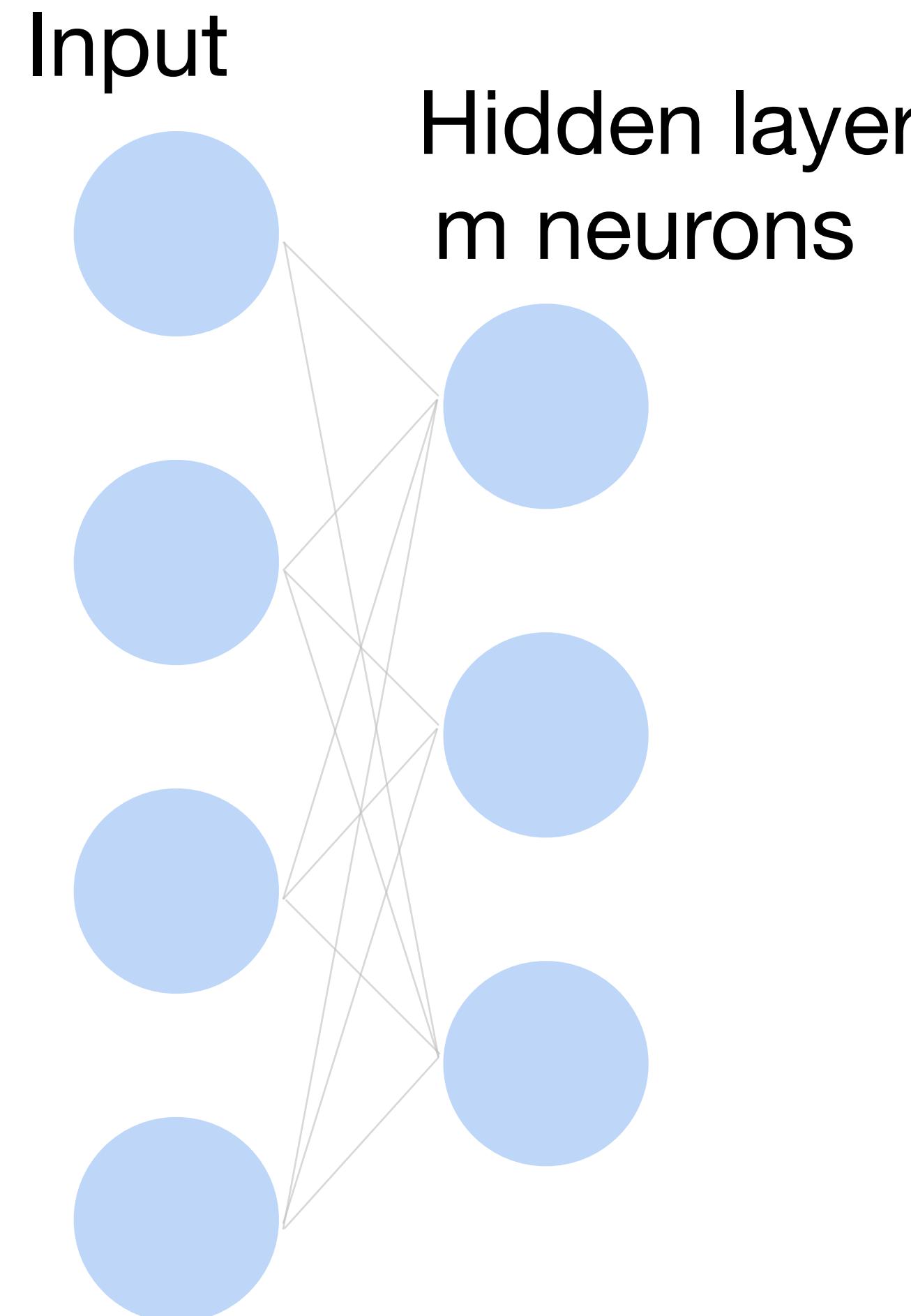


Multilayer Perceptron



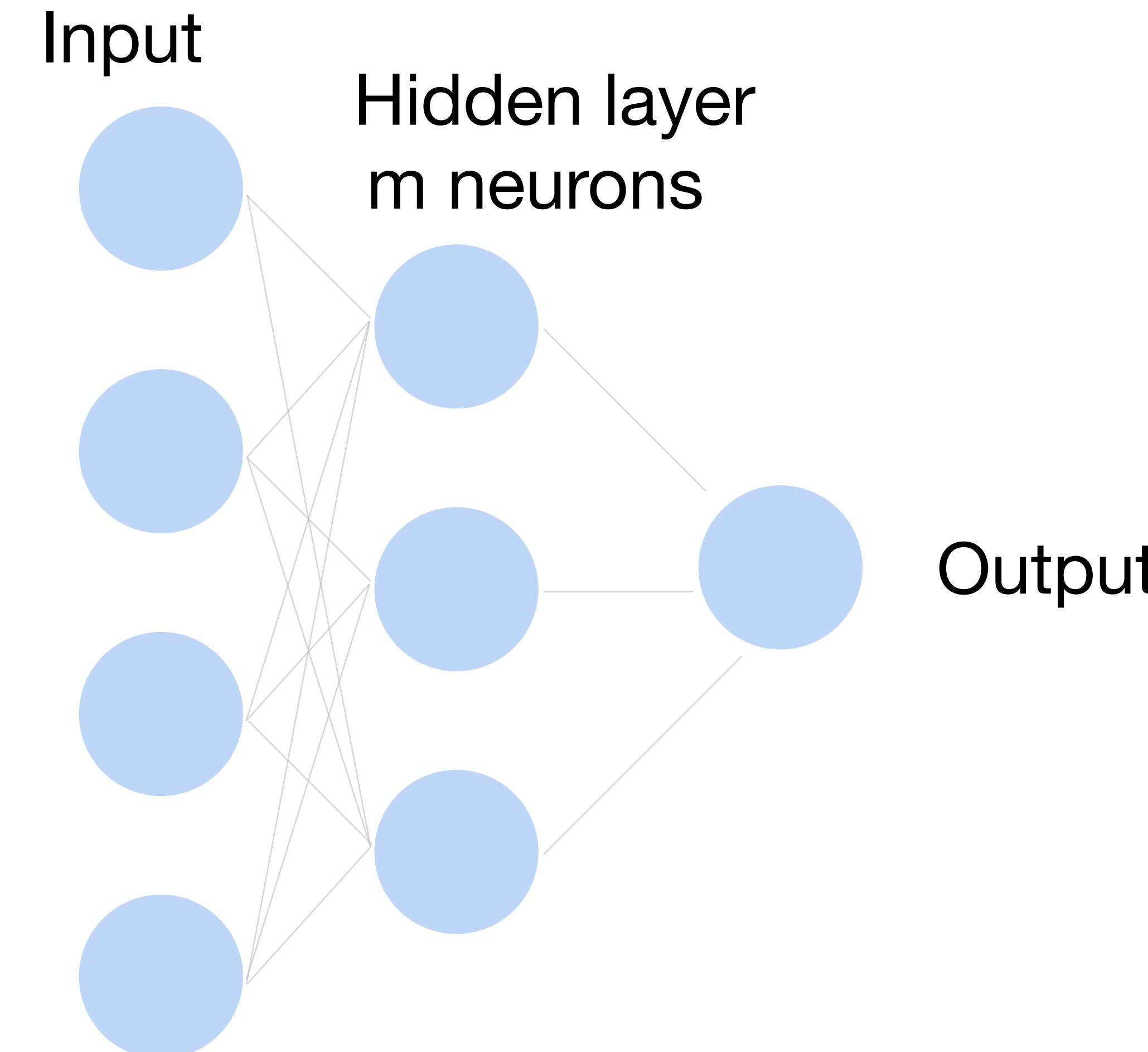
Single Hidden Layer

**How to classify
Cats vs. dogs?**

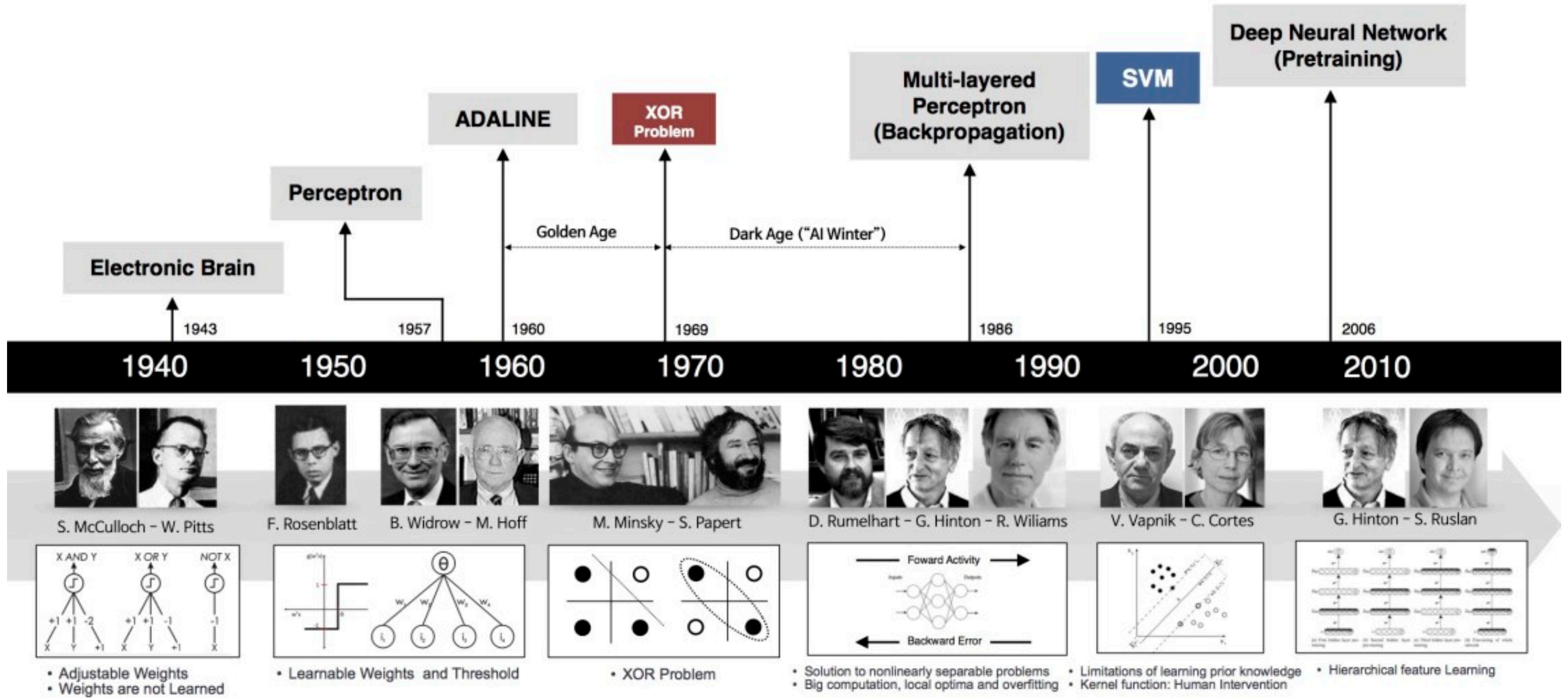


Single Hidden Layer

**How to classify
Cats vs. dogs?**



Brief history of neural networks



What we've learned today...

- Single-layer Perceptron
 - Motivation
 - Activation function
 - Representing AND, OR
- Brief history of neural networks



Thanks!

Based on slides from Xiaojin (Jerry) Zhu and Yingyu Liang (<http://pages.cs.wisc.edu/~jerryzhu/cs540.html>),
and Alex Smola: <https://courses.d2l.ai/berkeley-stat-157/units/mlp.html>