

Supervised Learning: find f

- Given: Training set $\{(x_i, y_i) \mid i = 1 \dots N\}$
- Find: A good approximation to $f : X \rightarrow Y$

Examples: what are X and Y ?

- Spam Detection
 - Map email to {Spam, Not Spam}
- Digit recognition
 - Map pixels to {0,1,2,3,4,5,6,7,8,9}
- Stock Prediction
 - Map new, historic prices, etc. to \mathfrak{R} (the real numbers)

A Supervised Learning Problem

Dataset:

| Example | x_1 | x_2 | x_3 | x_4 | y |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

- Our goal is to find a function $f : X \rightarrow Y$
 - $X = \{0,1\}^4$
 - $Y = \{0,1\}$
- **Question 1:** How should we pick the *hypothesis space*, the set of possible functions f ?
- **Question 2:** How do we find the best f in the hypothesis space?

Most General Hypothesis Space

Consider all possible boolean functions over four input features!

- 2^{16} possible hypotheses
- 2^9 are consistent with our dataset
- How do we choose the best one?

| x_1 | x_2 | x_3 | x_4 | y |
|-------|-------|-------|-------|-----|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

Dataset:

| Example | x_1 | x_2 | x_3 | x_4 | y |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

A Restricted Hypothesis Space

Consider all conjunctive boolean functions.

- 16 possible hypotheses
- None are consistent with our dataset
- How do we choose the best one?

| Rule | Counterexample |
|------------------------------------------------------|----------------|
| $\Rightarrow y$ | 1 |
| $x_1 \Rightarrow y$ | 3 |
| $x_2 \Rightarrow y$ | 2 |
| $x_3 \Rightarrow y$ | 1 |
| $x_4 \Rightarrow y$ | 7 |
| $x_1 \wedge x_2 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_3 \wedge x_4 \Rightarrow y$ | 4 |
| $x_1 \wedge x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |

Dataset:

| Example | x_1 | x_2 | x_3 | x_4 | y |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

Occam's Razor Principle

- William of **Occam**: Monk living in the 14th century
- Principle of parsimony:

“One should not increase, beyond what is necessary, the number of entities required to explain anything”

- When **many** solutions are available for a given problem, we should select the **simplest** one
- But what do we mean by **simple**?
- We will use **prior knowledge** of the problem to solve to define what is a simple solution

Example of a prior: smoothness

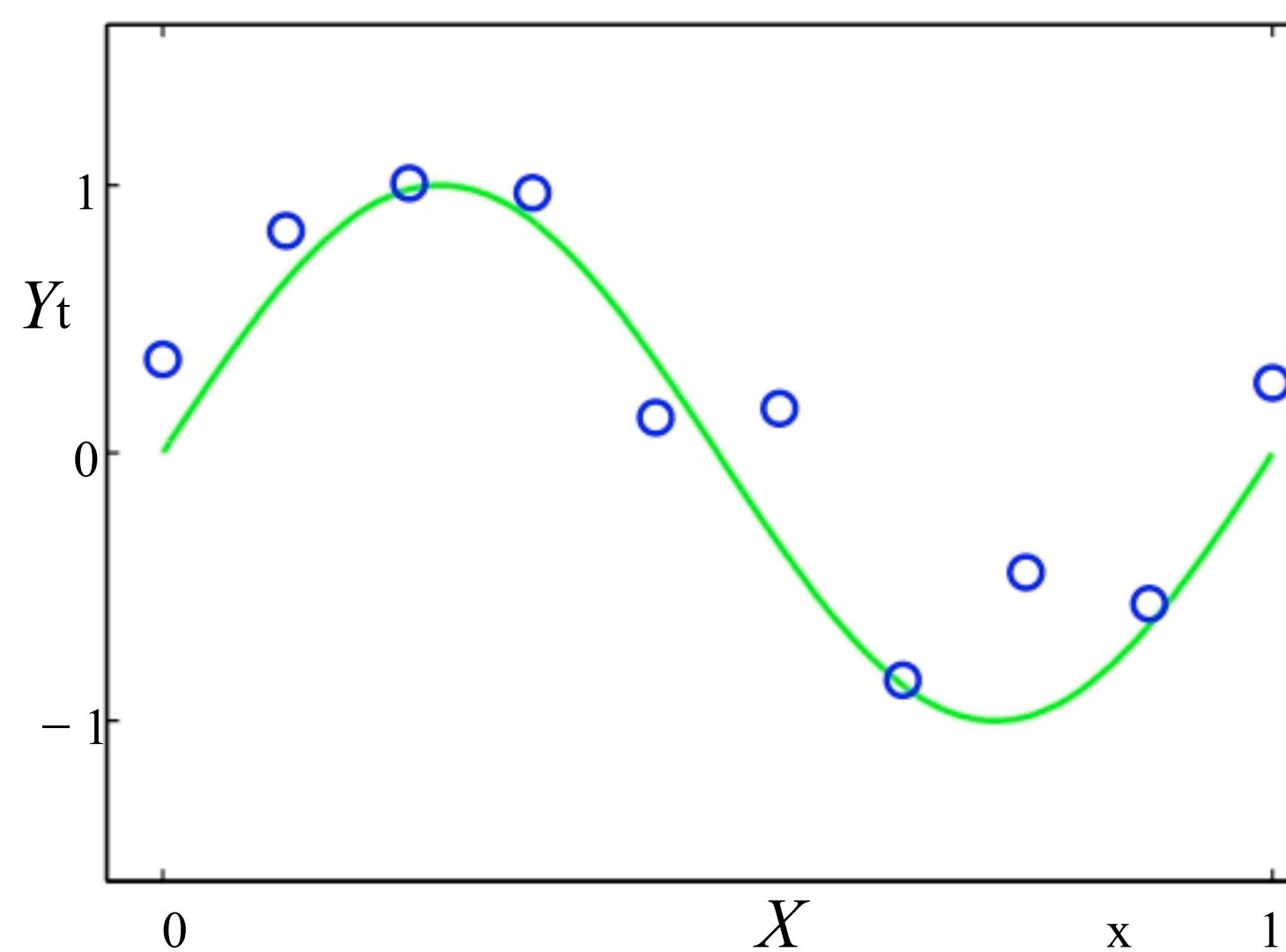
[Samy Bengio]

Key Issues in Machine Learning

- How do we choose a hypothesis space?
 - Often we use **prior knowledge** to guide this choice
- How can we gauge the accuracy of a hypothesis on unseen data?
 - **Occam's razor:** use the *simplest* hypothesis consistent with data!
This will help us avoid overfitting.
 - **Learning theory** will help us quantify our ability to **generalize** as a function of the amount of training data and the hypothesis space
- How do we find the best hypothesis?
 - This is an **algorithmic** question, the main topic of computer science
- How to model applications as machine learning problems?
(engineering challenge)

Second example: Regression

Dataset: 10 (X,Y) points generated
from a sin function, with noise



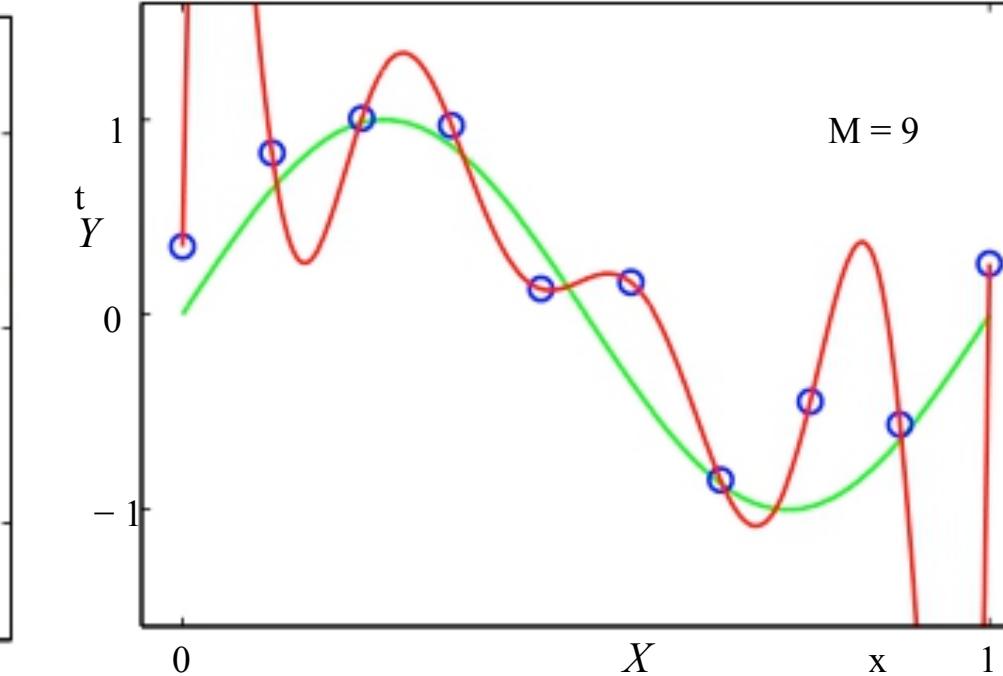
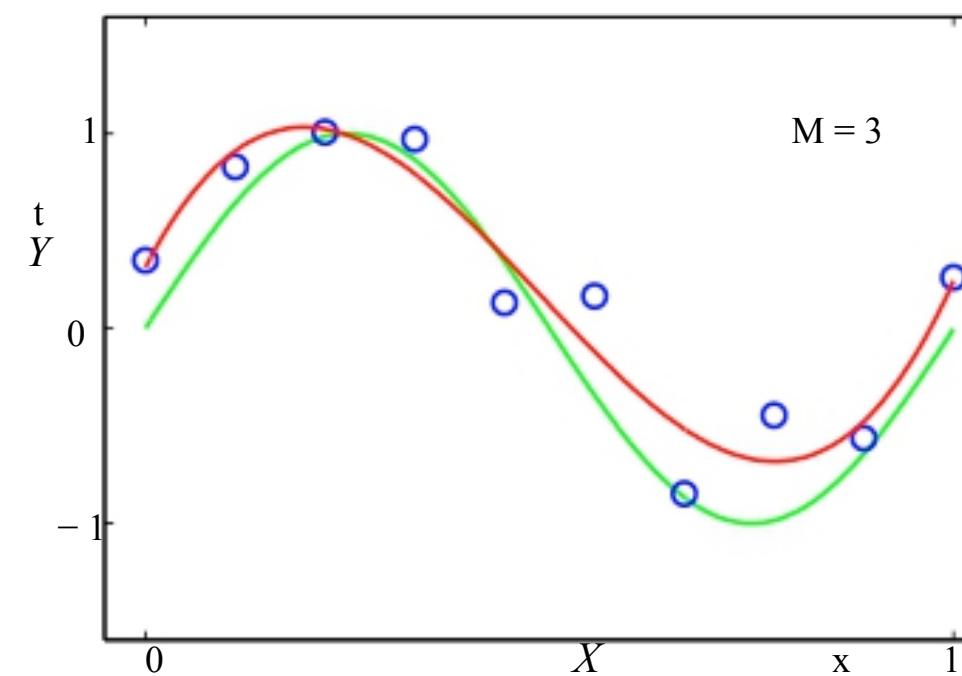
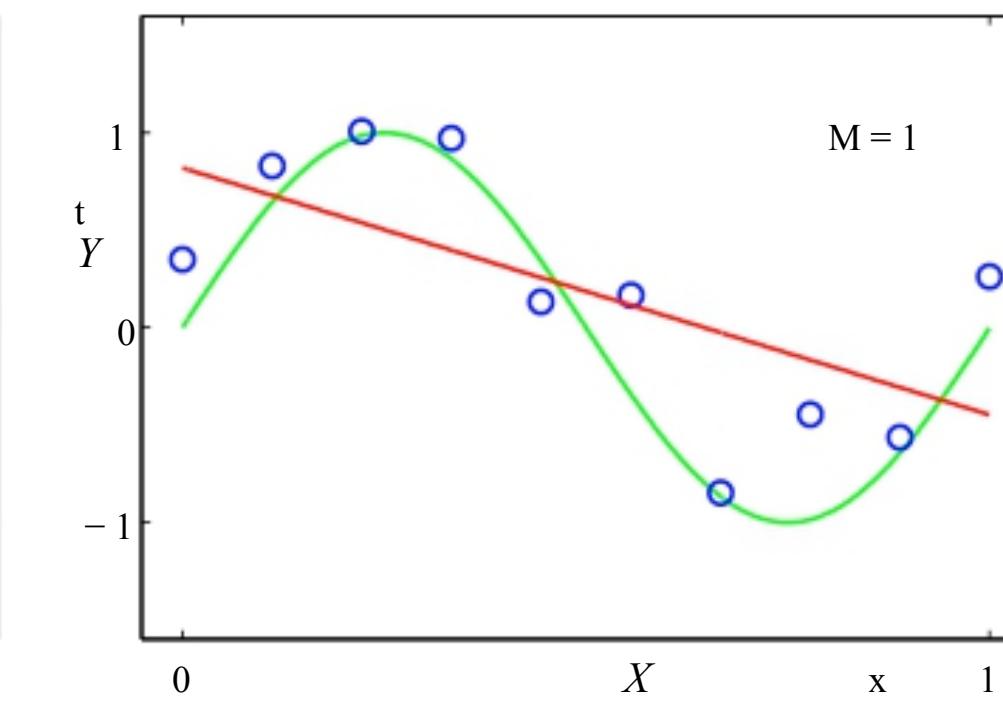
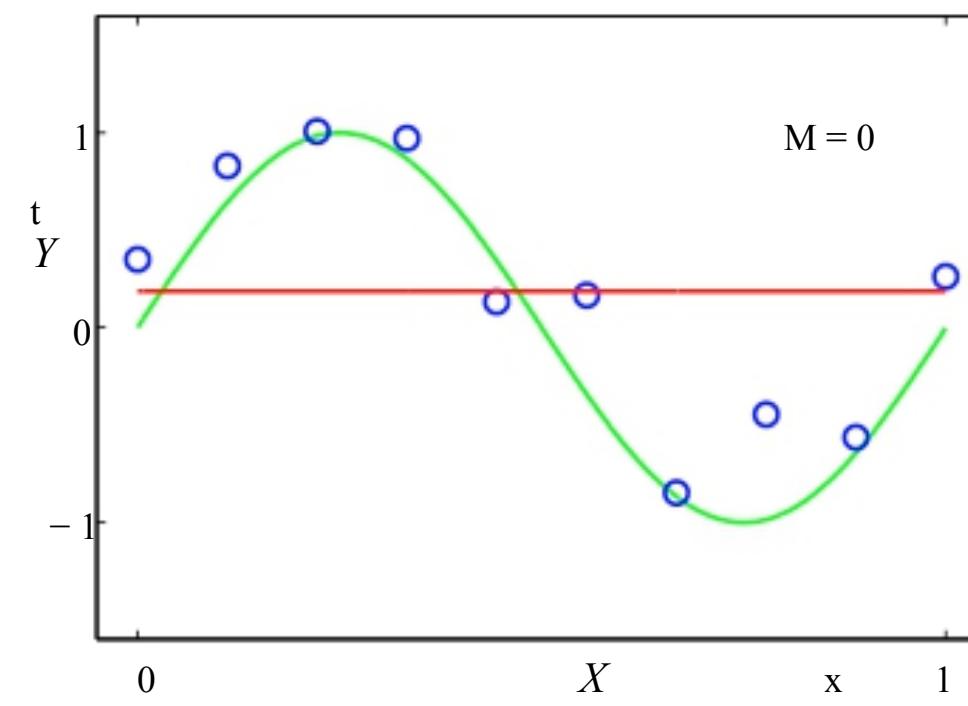
- Regression:
 - $f : X \rightarrow Y$
 - $X = \Re$
 - $Y = \Re$

[Bishop]

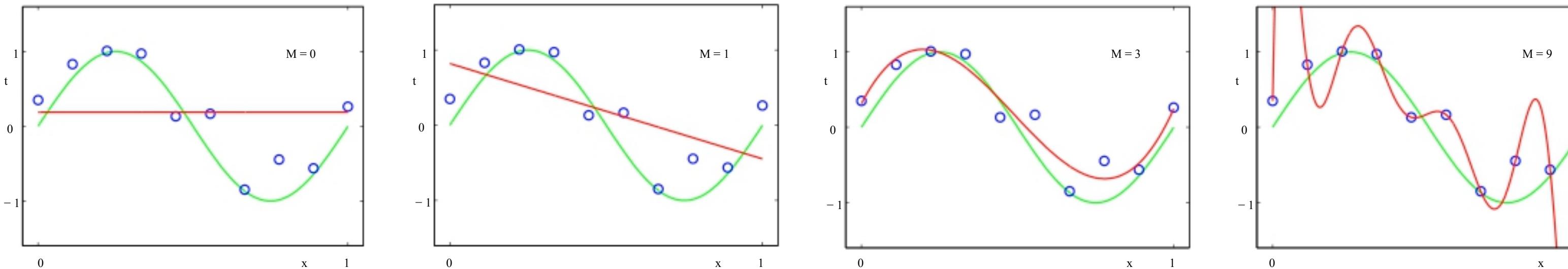
Degree-M Polynomials

How about letting f be a degree M polynomial?

- Which one is **best** ?



Hypo. Space: Degree-N Polynomials



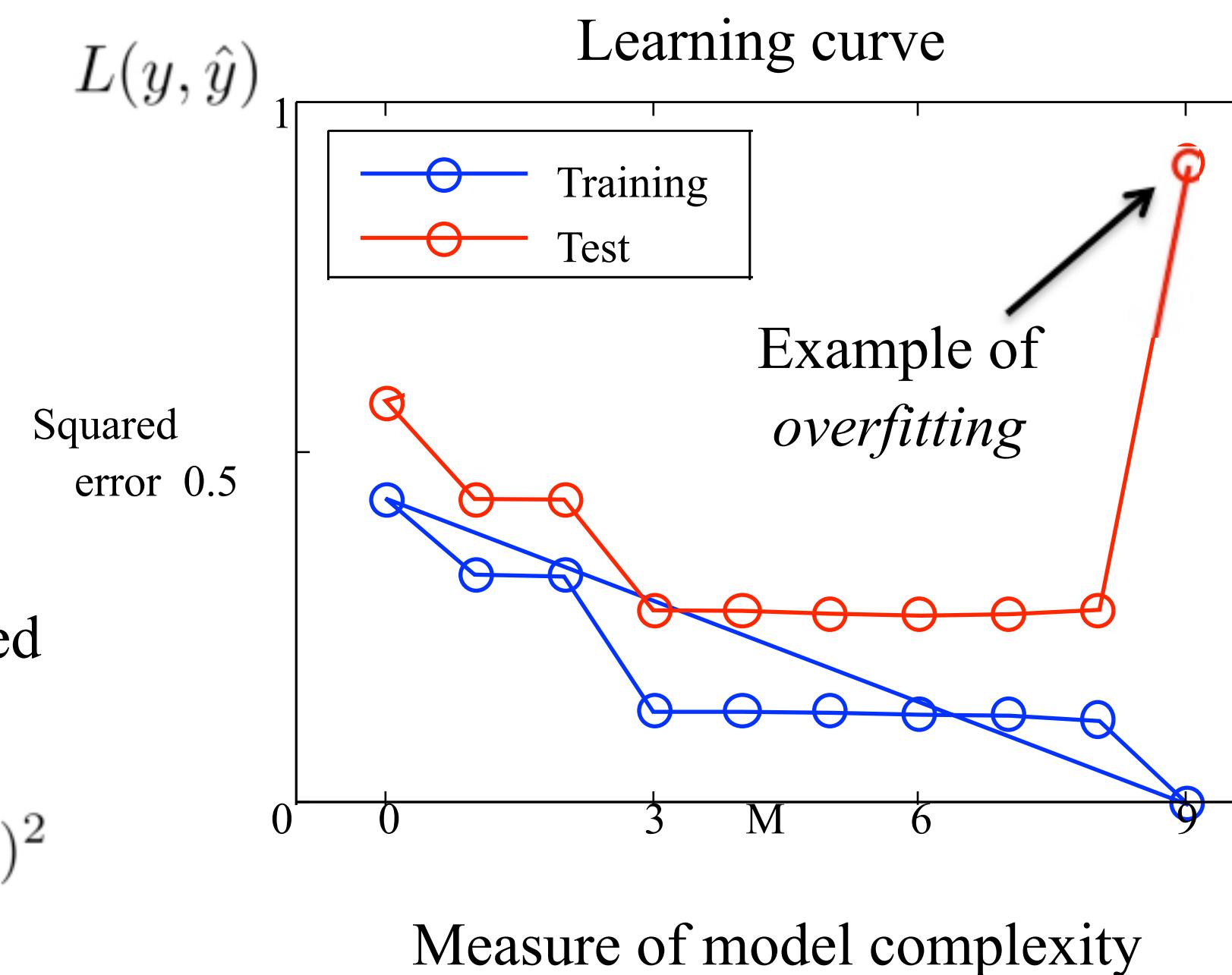
We measure error using a *loss function*

For regression, a common choice is squared loss:

$$L(y_i, f(x_i)) = (y_i - f(x_i))^2$$

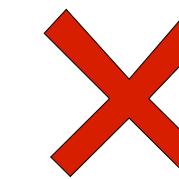
The *empirical loss* of the function f applied to the training data is then:

$$\frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$$



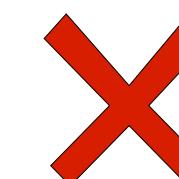
Binary classification

- Input: email
- Output: spam/ham
- Setup:
 - Get a large collection of example emails, each labeled “spam” or “ham”
 - Note: someone has to hand label all this data!
 - Want to learn to predict labels of new, future emails
- Features: The attributes used to make the ham / spam decision
 - Words: FREE!
 - Text Patterns: \$dd, CAPS
 - Non-text: SenderInContacts
 - ...



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE
MAILINGS, SIMPLY REPLY TO THIS
MESSAGE AND PUT "REMOVE" IN THE
SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



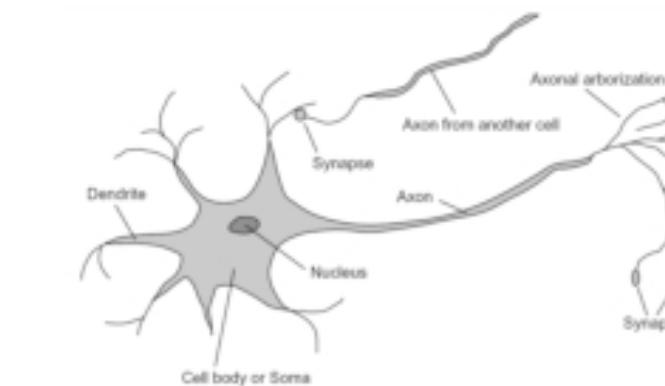
Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

The perceptron algorithm

- 1957: Perceptron algorithm invented by Rosenblatt

Wikipedia: “A handsome bachelor, he drove a classic MGA sports... for several years taught an interdisciplinary undergraduate honors course entitled "Theory of Brain Mechanisms" that drew students equally from Cornell's Engineering and Liberal Arts colleges...this course was a melange of ideas .. experimental brain surgery on epileptic patients while conscious, experiments on .. the visual cortex of cats, ... analog and digital electronic circuits that modeled various details of neuronal behavior (i.e. the perceptron itself, as a machine).”

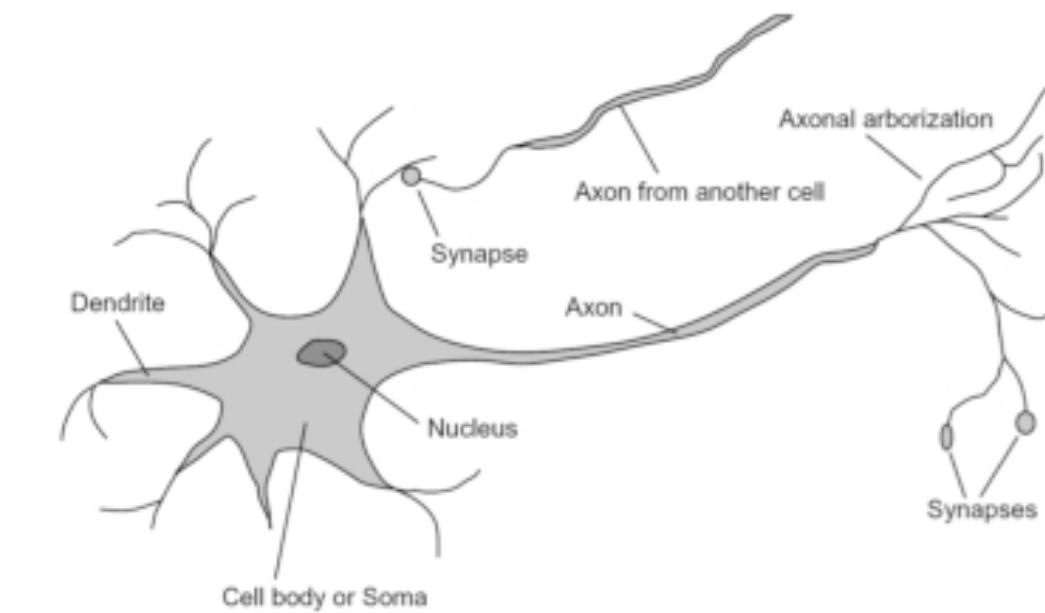
- Built on work of Hebb (1949); also developed by Widrow-Hoff (1960)
- 1960: Perceptron Mark 1 Computer – hardware implementation
- 1969: Minsky & Papert book shows perceptrons limited to *linearly separable* data
- 1970’s: Learning methods for two-layer neural networks



[William Cohen]

Linear Classifiers

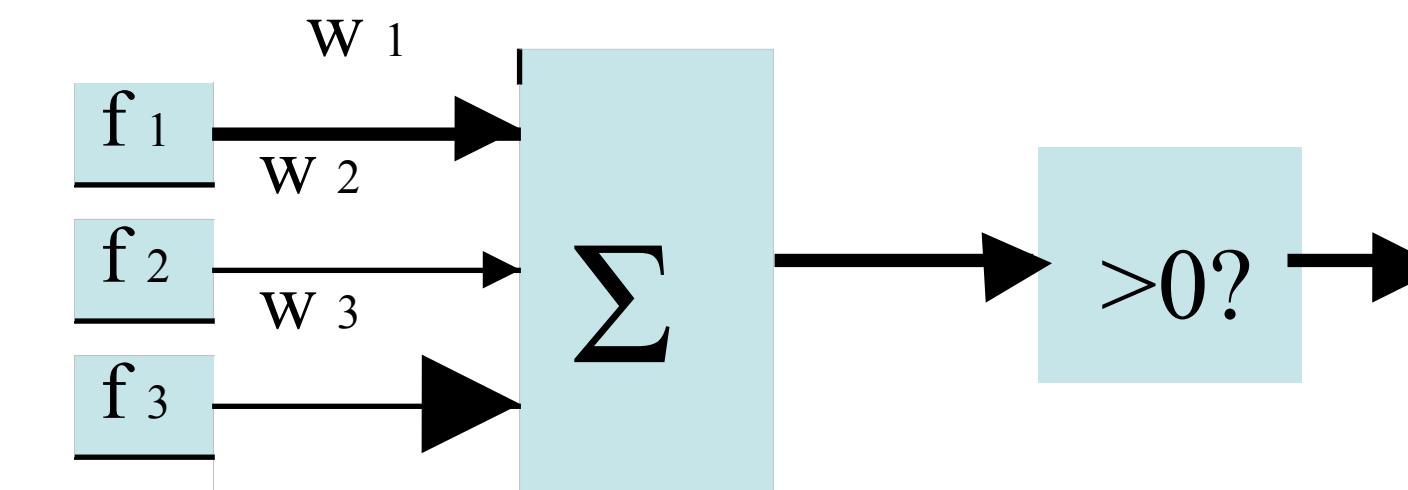
- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



Important note: changing notation!

$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
 - Positive, output *class 1*
 - Negative, output *class 2*



Example: Spam

- Imagine 3 features (spam is “positive” class):

- free (number of occurrences of “free”)
- money (occurrences of “money”)
- BIAS (intercept, always has value 1)

x
“ free money”

| $f(x)$ | |
|--------|-----|
| BIAS | : 1 |
| free | : 1 |
| money | : 1 |
| ... | |

w

| | |
|-------|------|
| BIAS | : -3 |
| free | : 4 |
| money | : 2 |
| ... | |

$$\begin{aligned} & w \cdot f(x) \\ & \sum_i w_i \cdot f_i(x) \\ & (1)(-3) + \\ & (1)(4) + \\ & (1)(2) + \\ & \dots \\ & = 3 \end{aligned}$$

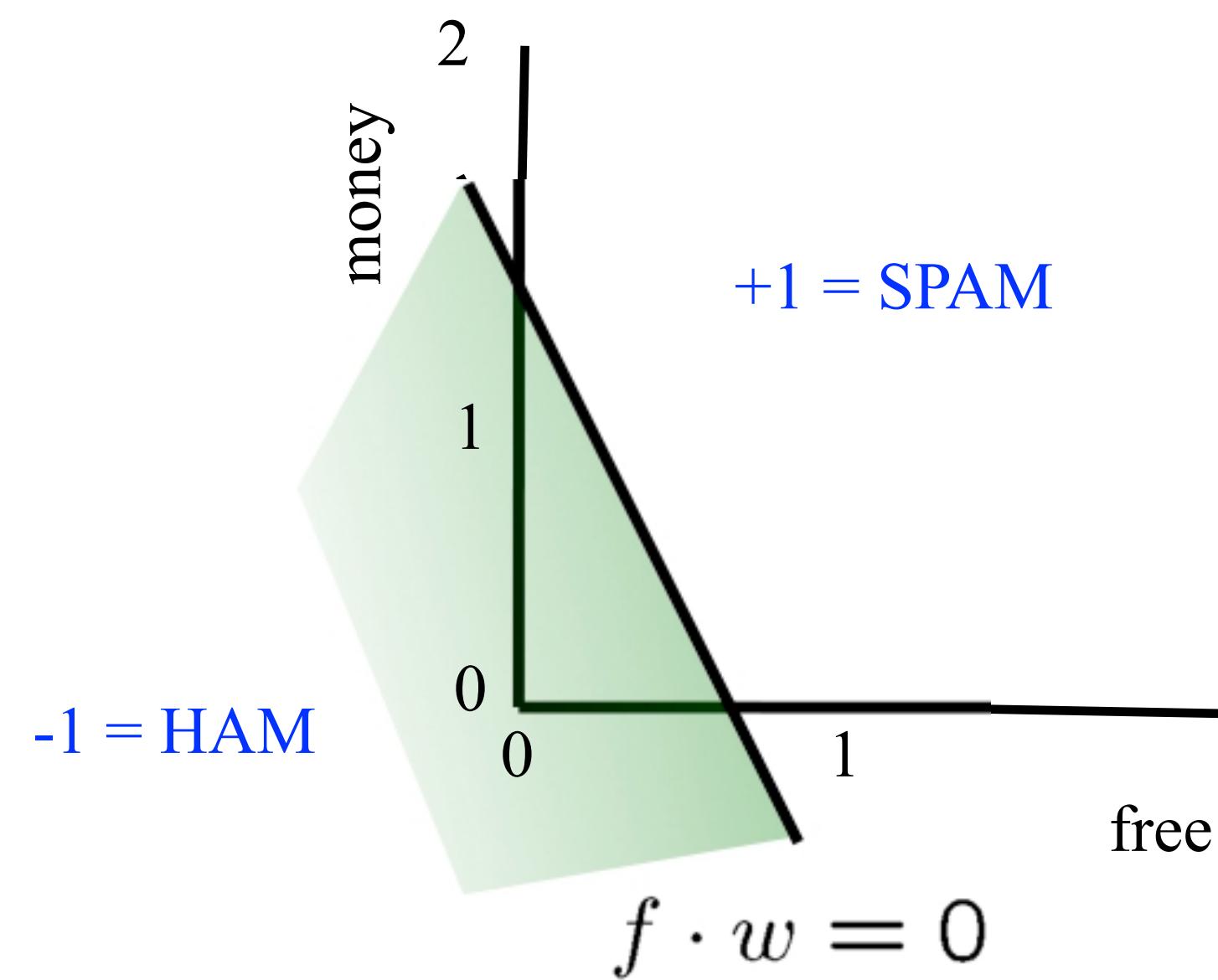
$w.f(x) > 0$ SPAM!!!

Binary Decision Rule

- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to $Y=+1$
 - Other corresponds to $Y=-1$

w

| |
|-----------|
| BIAS : -3 |
| free : 4 |
| money : 2 |
| ... |



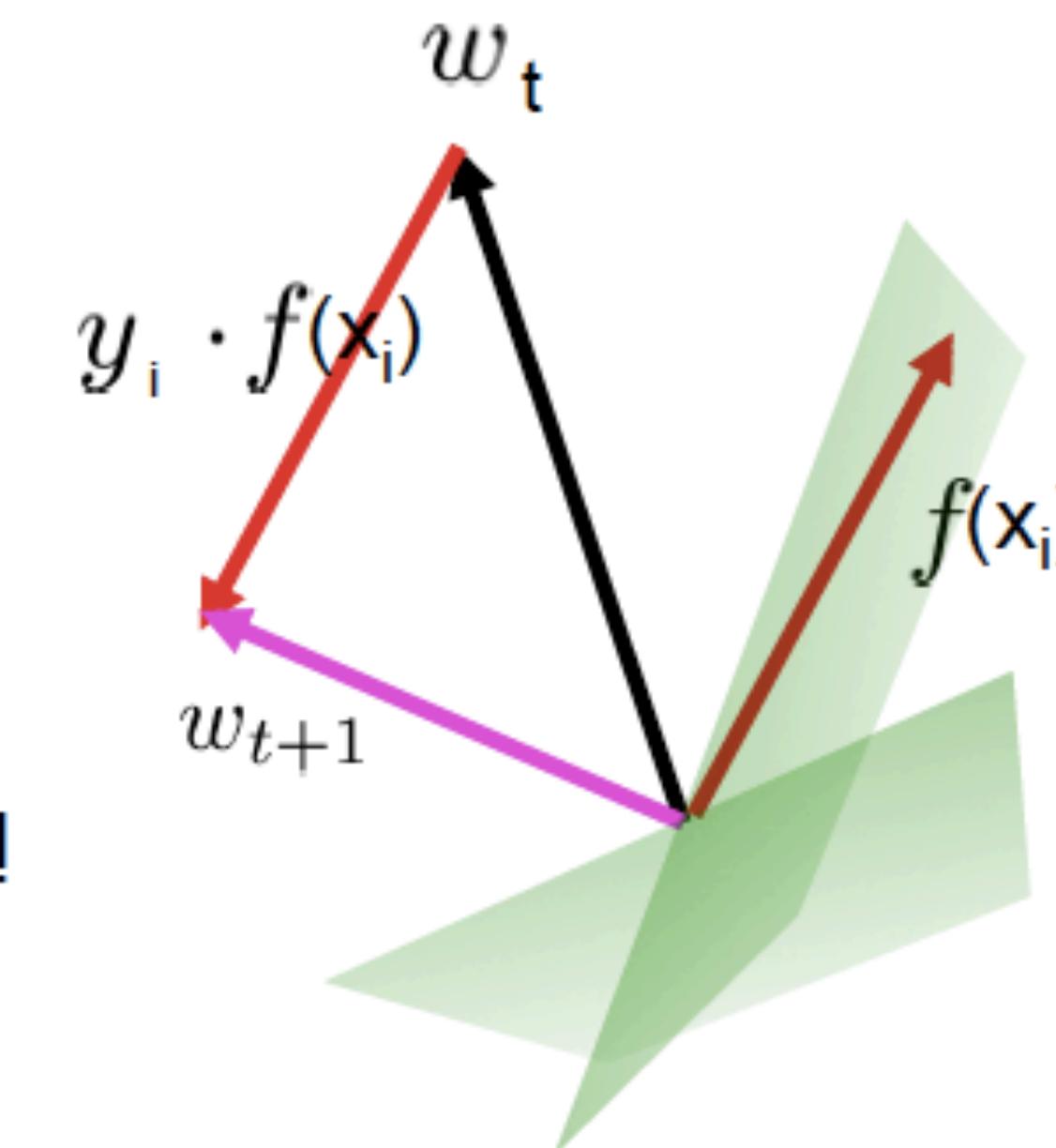
The Perception algorithm

- Start with weight vector = $\vec{0}$
- For each training instance (x_i, y_i) :
 - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x_i) \geq 0 \\ -1 & \text{if } w \cdot f(x_i) < 0 \end{cases}$$

- If correct (i.e., $y=y_i$), no change!
- If wrong: update

$$w = w + y_i f(x_i)$$



The Perceptron Algorithm:

1. Start with the all-zeroes weight vector $\mathbf{w}_1 = \mathbf{0}$, and initialize t to 1. Also let's automatically scale all examples \mathbf{x} to have (Euclidean) length 1, since this doesn't affect which side of the plane they are on.
2. Given example \mathbf{x} , predict positive iff $\mathbf{w}_t \cdot \mathbf{x} > 0$.
3. On a mistake, update as follows:
 - Mistake on positive: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \mathbf{x}$.
 - Mistake on negative: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \mathbf{x}$.

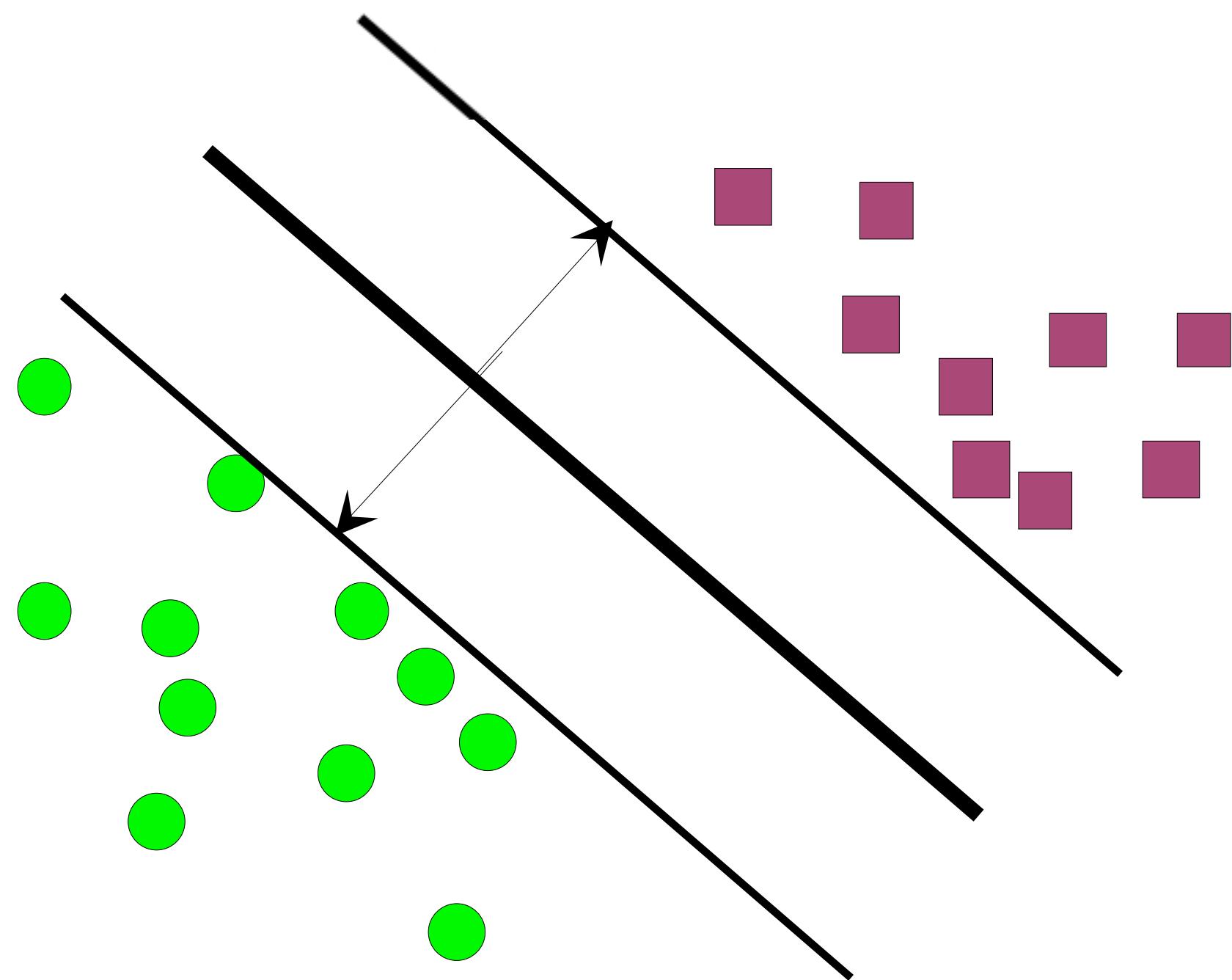
$$t \leftarrow t + 1.$$

Linearly Separable

$\exists \mathbf{w}$ such that $\forall t$

$$y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq \gamma > 0$$

Called the *functional margin*
with respect to the training set



Equivalently, for $y_t = +1$,

$$\mathbf{w} \cdot \mathbf{x}_t \geq \gamma$$

and for $y_t = -1$,

$$\mathbf{w} \cdot \mathbf{x}_t \leq -\gamma$$

Mistake Bound for Perceptron

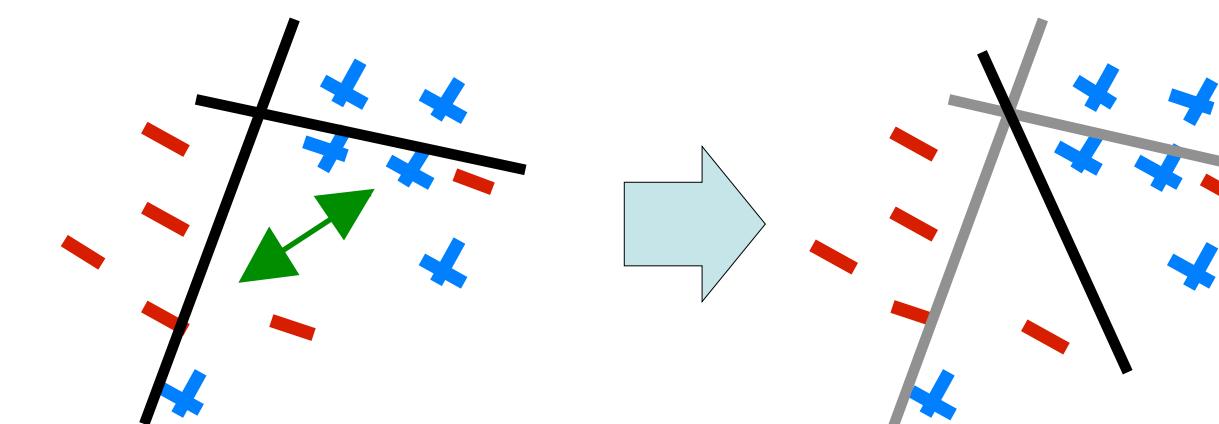
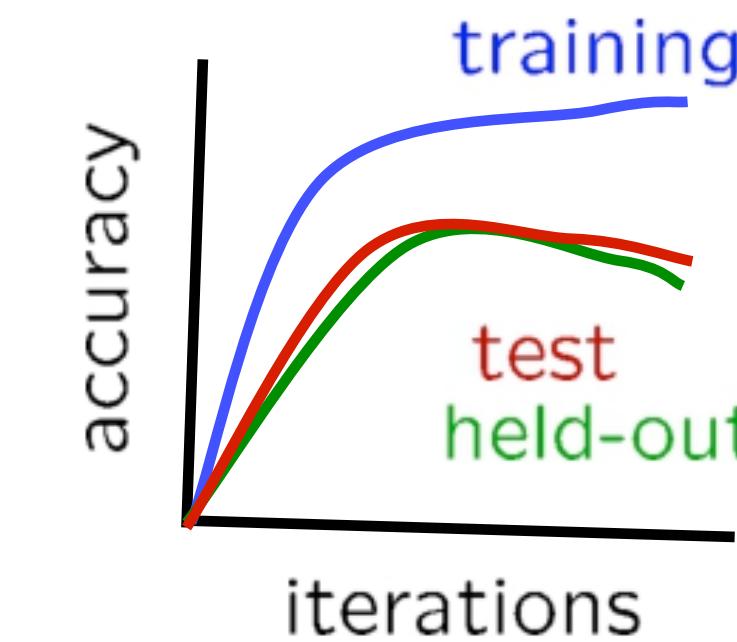
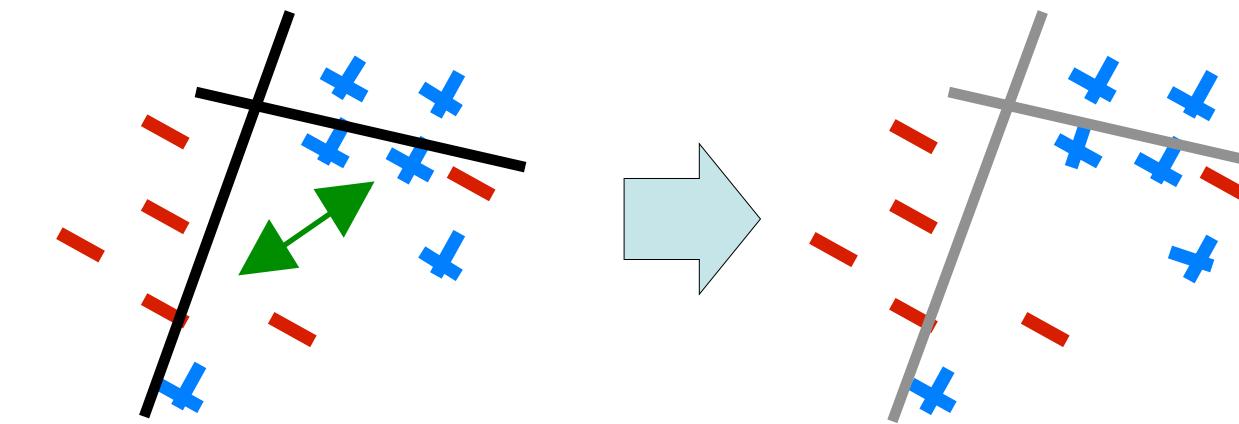
- Assume the data set D is linearly separable with *geometric* margin γ , i.e.,

$$\exists w^* \text{ s.t. } \|w^*\|_2 = 1 \text{ and } \forall t, y_t(w^* \cdot x_t) \geq \gamma$$

- Assume $\|x_t\|_2 \leq R, \forall t$
- Theorem: The maximum number of mistakes made by the perceptron algorithm is bounded by R^2/γ^2

Problems with the perceptron algorithm

- If the data isn't linearly separable, no guarantees of convergence or training accuracy
- Even if the training data is linearly separable, perceptron can overfit
- **Averaged** perceptron is an algorithmic modification that helps with both issues
 - Averages the weight vectors across all iterations

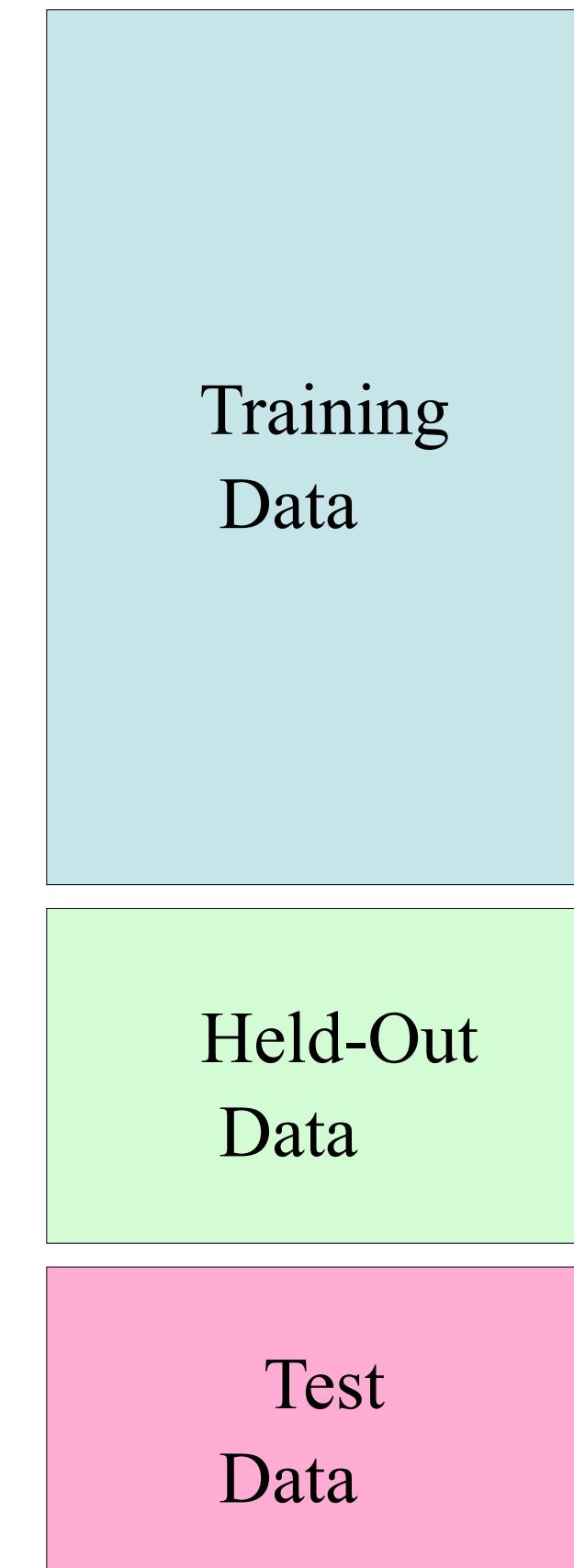


ML Methodology

- **Data:** labeled instances, e.g. emails marked spam/ham
 - Training set
 - Held out set (sometimes call Validation set)
 - Test set

Randomly allocate to these three, e.g. 60/20/20

- **Features:** attribute-value pairs which characterize each x
- **Experimentation cycle**
 - Select a hypothesis f
(Tune hyperparameters on held-out or *validation* set)
 - Compute accuracy of test set
 - Very important: never “peek” at the test set!
- **Evaluation**
 - Accuracy: fraction of instances predicted correctly



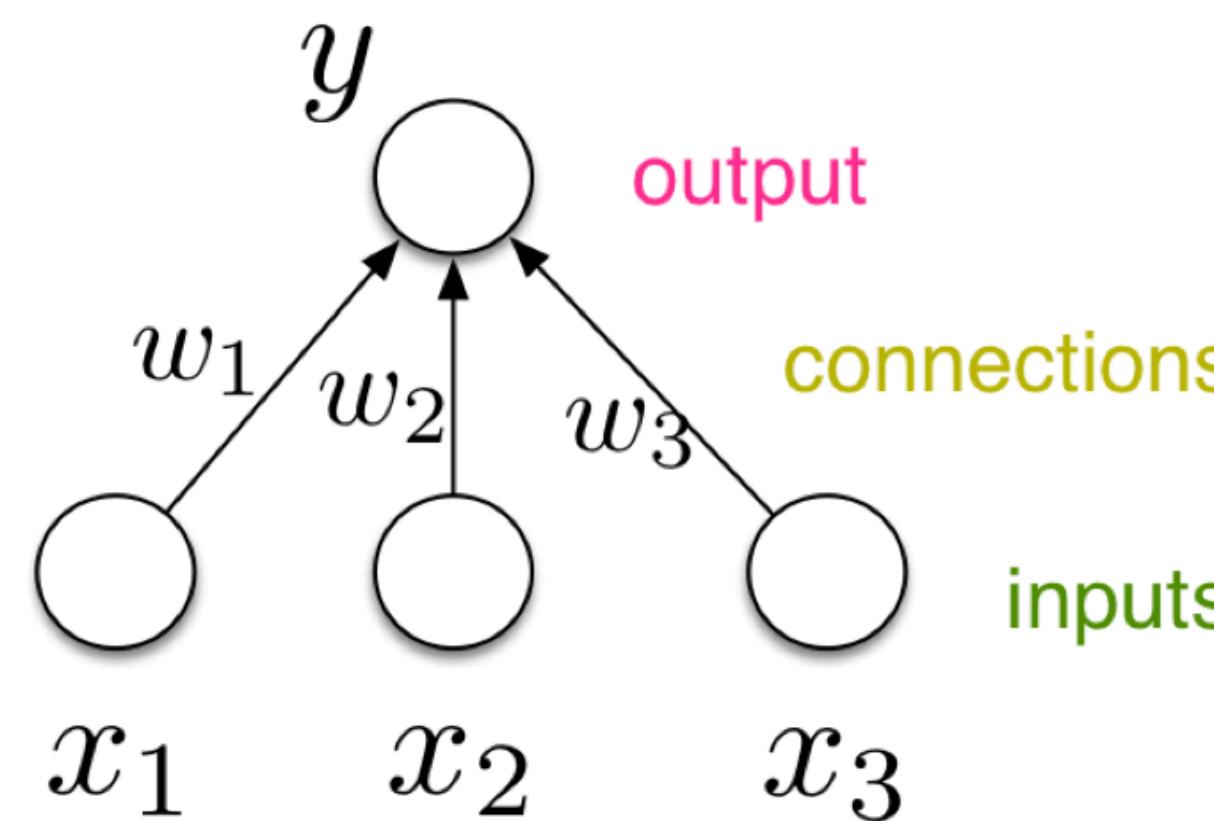
What are neural networks?

Why neural nets?

- inspiration from the brain
 - proof of concept that a neural architecture can see and hear!
- very effective across a range of applications (vision, text, speech, medicine, robotics, etc.)
- widely used in both academia and the tech industry
- powerful software frameworks (PyTorch, TensorFlow, etc.) let us quickly implement sophisticated algorithms

What are neural networks?

- Most of the biological details aren't essential, so we use vastly simplified models of neurons.
- While neural nets originally drew inspiration from the brain, nowadays we mostly think about math, statistics, etc.



$$y = \phi(\mathbf{w}^T \mathbf{x} + b)$$

Annotations for the equation:

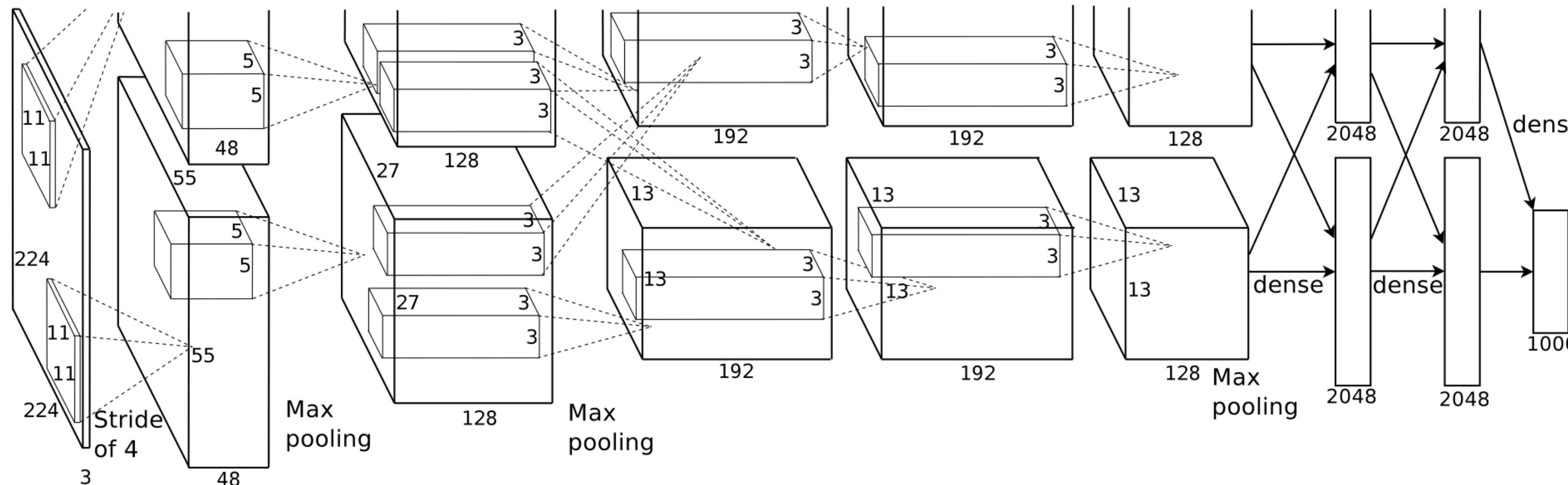
- "output" points to the variable y .
- "weights" points to the term \mathbf{w}^T .
- "bias" points to the term b .
- "activation function" points to the term ϕ .
- "inputs" points to the term \mathbf{x} .

- Neural networks are collections of thousands (or millions) of these simple processing units that together perform useful computations.

“Deep learning”

Deep learning: many layers (stages) of processing

E.g. this network which recognizes objects in images:

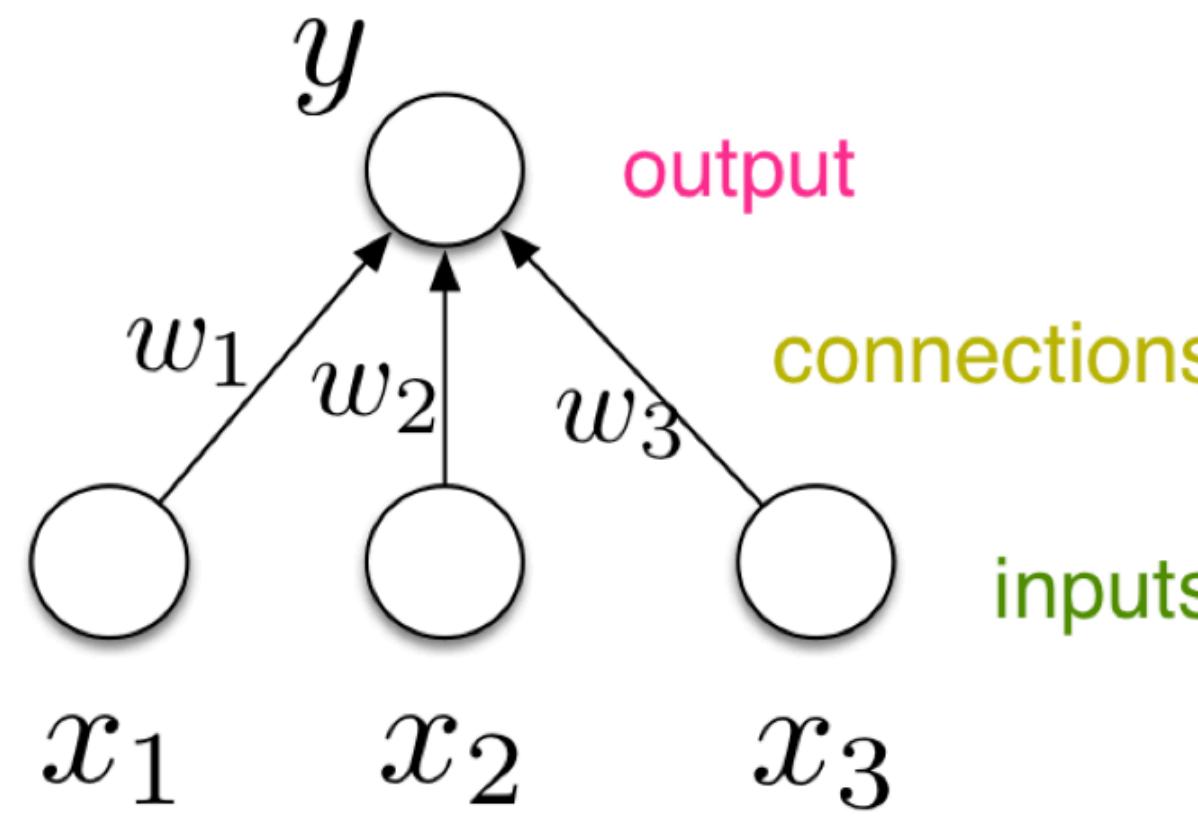


(Krizhevsky et al., 2012)

Each of the boxes consists of many neuron-like units similar to the one on the previous slide!

Overview

- Recall the simple neuron-like unit:



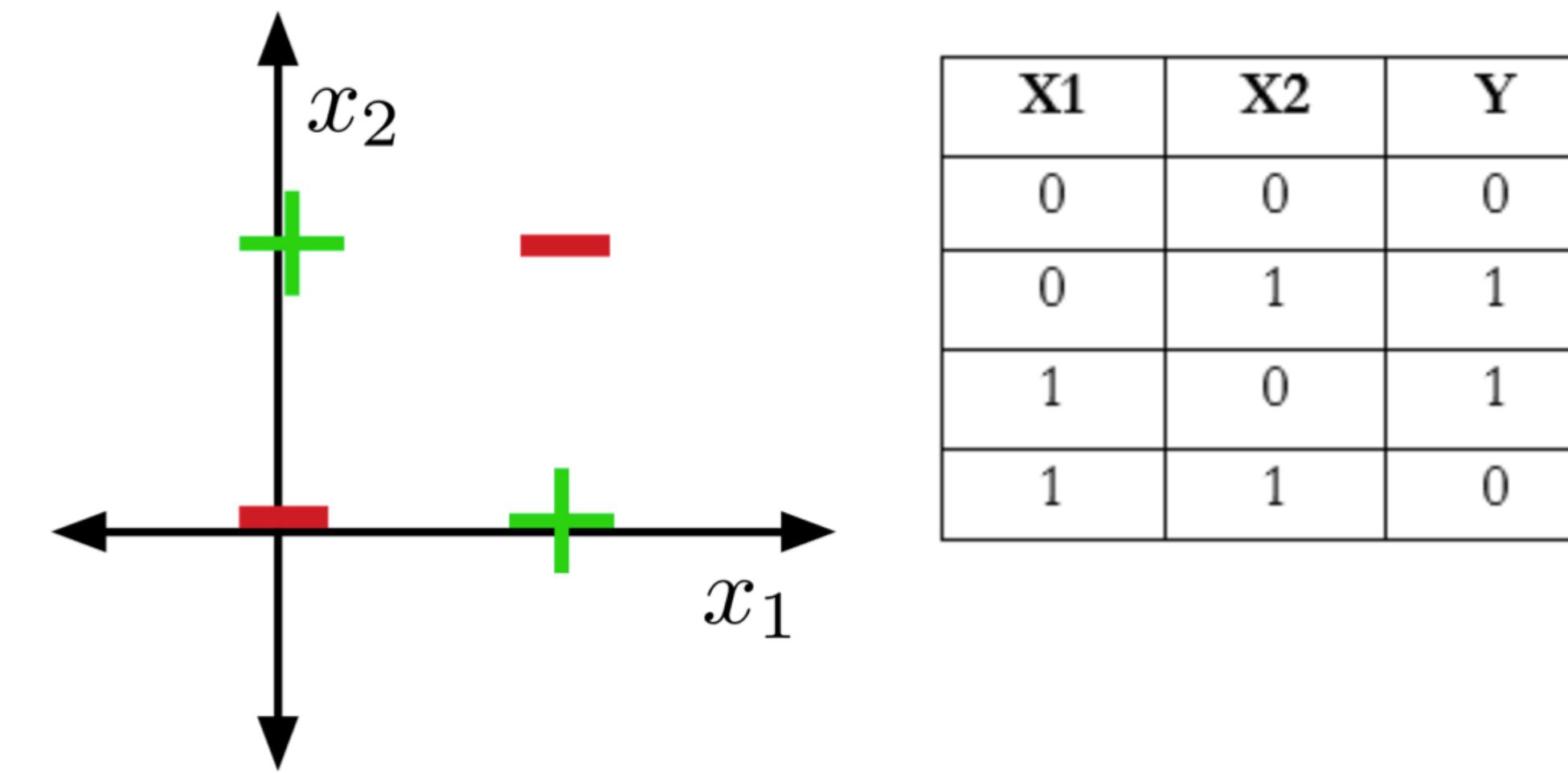
An equation for a simple neuron-like unit: $y = \phi(\mathbf{w}^\top \mathbf{x} + b)$. The equation is annotated with arrows pointing to its components:

- A pink arrow points down to the variable y , labeled "output".
- A blue arrow points down to the term $\mathbf{w}^\top \mathbf{x}$, labeled "weights".
- A blue arrow points down to the term b , labeled "bias".
- A red arrow points up to the activation function symbol ϕ , labeled "activation function".
- A green arrow points up to the input vector \mathbf{x} , labeled "inputs".

- Linear regression and logistic regression can each be viewed as a single unit.
- These units are much more powerful if we connect many of them into a neural network.

Limits of Linear Classification

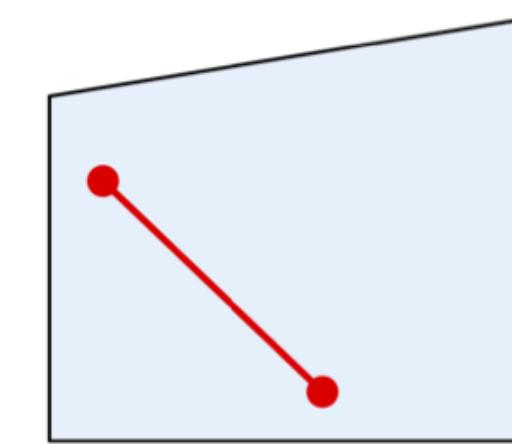
- Single neurons (linear classifiers) are very limited in expressive power.
- **XOR** is a classic example of a function that's not linearly separable.



- There's an elegant proof using convexity.

Limits of Linear Classification

Convex Sets



- A set \mathcal{S} is **convex** if any line segment connecting points in \mathcal{S} lies entirely within \mathcal{S} . Mathematically,

$$\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{S} \implies \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{S} \quad \text{for } 0 \leq \lambda \leq 1.$$

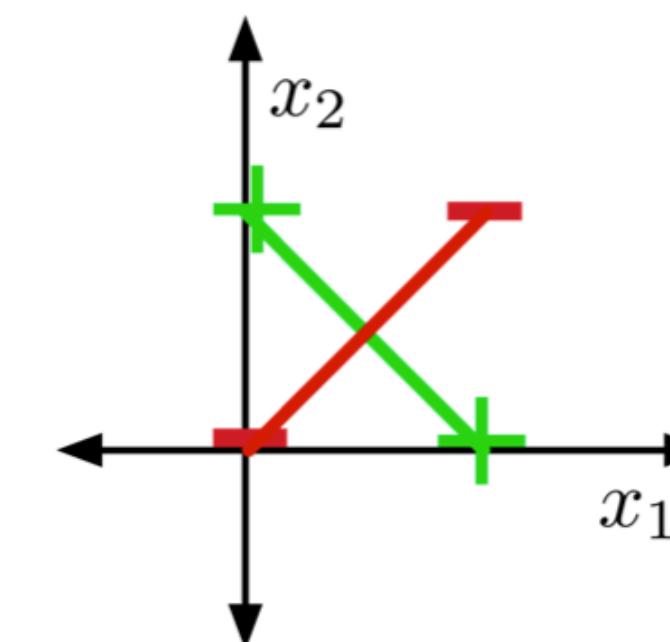
- A simple inductive argument shows that for $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{S}$, **weighted averages**, or **convex combinations**, lie within the set:

$$\lambda_1 \mathbf{x}_1 + \dots + \lambda_N \mathbf{x}_N \in \mathcal{S} \quad \text{for } \lambda_i > 0, \lambda_1 + \dots + \lambda_N = 1.$$

Limits of Linear Classification

Showing that XOR is not linearly separable

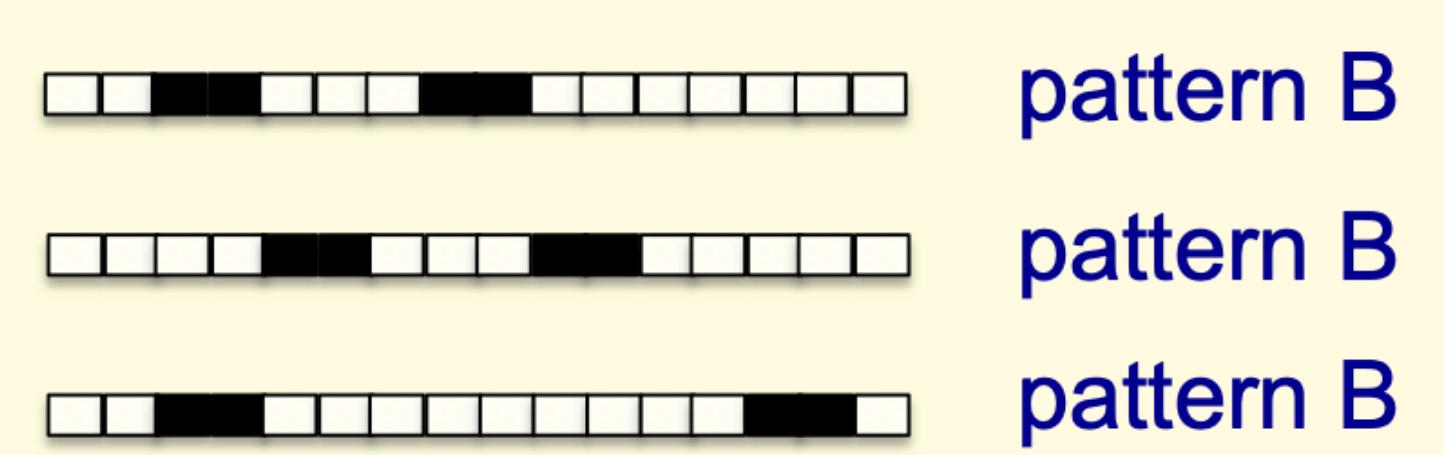
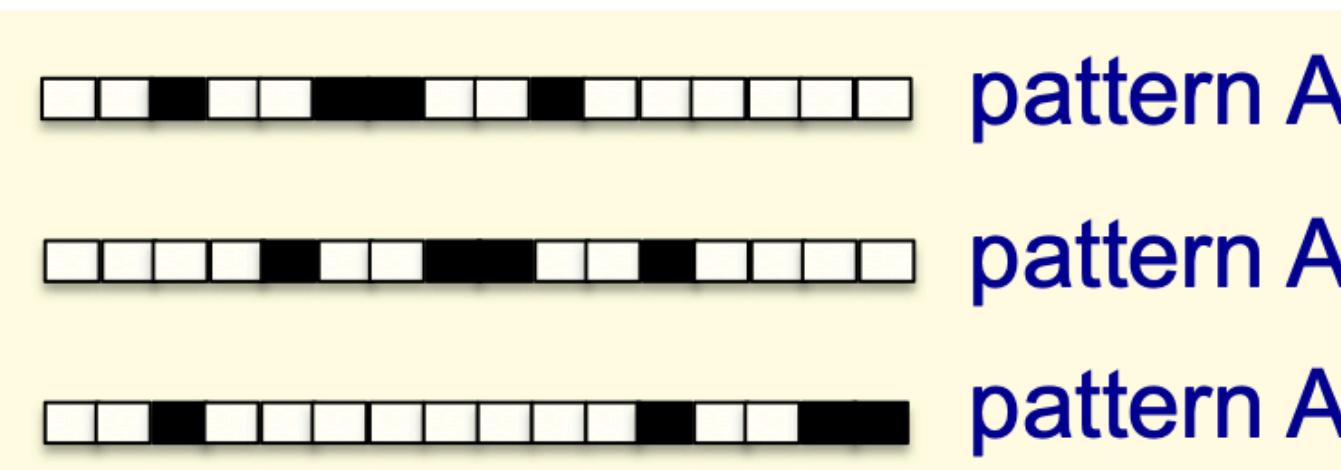
- Half-spaces are obviously convex.
- Suppose there were some feasible hypothesis. If the positive examples are in the positive half-space, then the green line segment must be as well.
- Similarly, the red line segment must lie within the negative half-space.



- But the intersection can't lie in both half-spaces. Contradiction!

Limits of Linear Classification

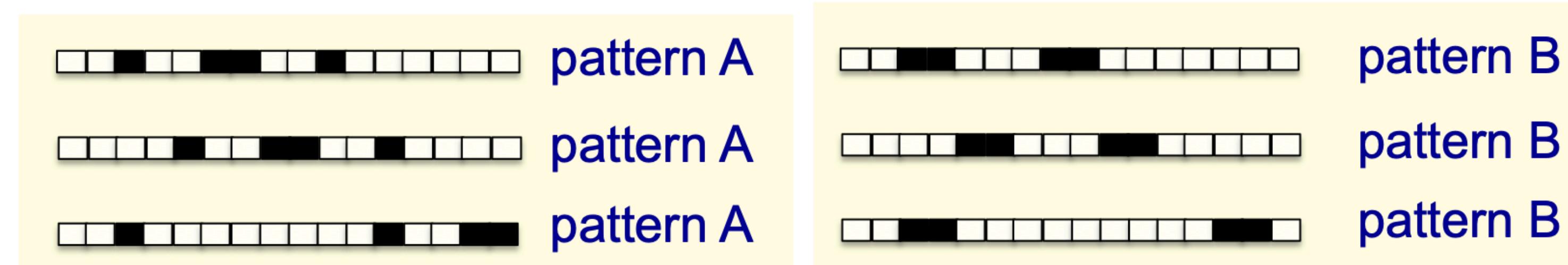
A more troubling example



- These images represent 16-dimensional vectors. White = 0, black = 1.
- Want to distinguish patterns A and B in all possible translations (with wrap-around)
- Translation invariance is commonly desired in vision!

Limits of Linear Classification

A more troubling example



- These images represent 16-dimensional vectors. White = 0, black = 1.
- Want to distinguish patterns A and B in all possible translations (with wrap-around)
- Translation invariance is commonly desired in vision!
- Suppose there's a feasible solution. The average of all translations of A is the vector $(0.25, 0.25, \dots, 0.25)$. Therefore, this point must be classified as A.
- Similarly, the average of all translations of B is also $(0.25, 0.25, \dots, 0.25)$. Therefore, it must be classified as B. Contradiction!

Limits of Linear Classification

- Sometimes we can overcome this limitation using feature maps, just like for linear regression. E.g., for **XOR**:

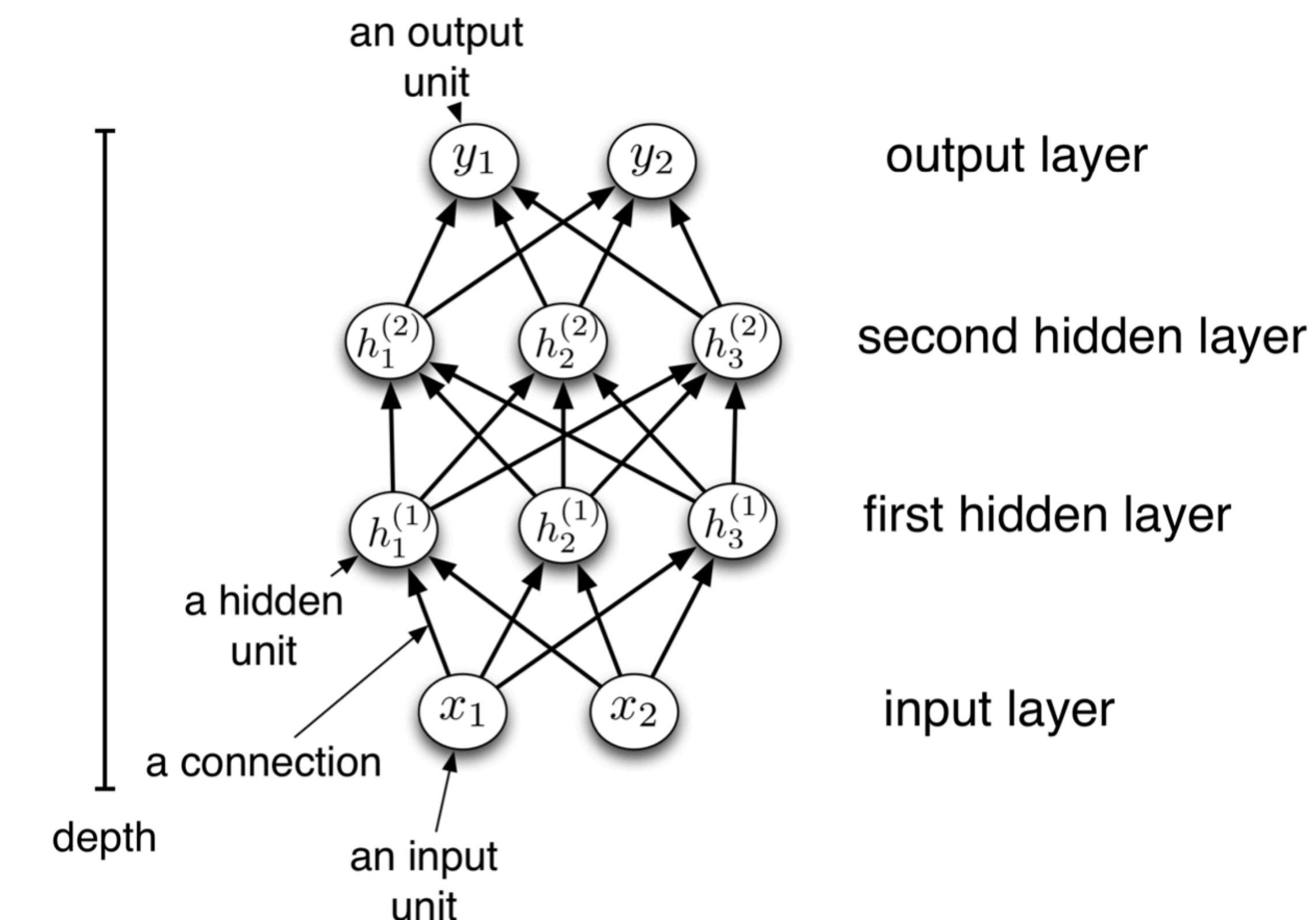
$$\psi(\mathbf{x}) = \begin{pmatrix} x_1 \\ x_2 \\ x_1x_2 \end{pmatrix}$$

| x_1 | x_2 | $\phi_1(\mathbf{x})$ | $\phi_2(\mathbf{x})$ | $\phi_3(\mathbf{x})$ | t |
|-------|-------|----------------------|----------------------|----------------------|-----|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 1 | 1 | 0 |

- This is linearly separable. (Try it!)
- Not a general solution: it can be hard to pick good basis functions.
Instead, we'll use neural nets to learn nonlinear hypotheses directly.

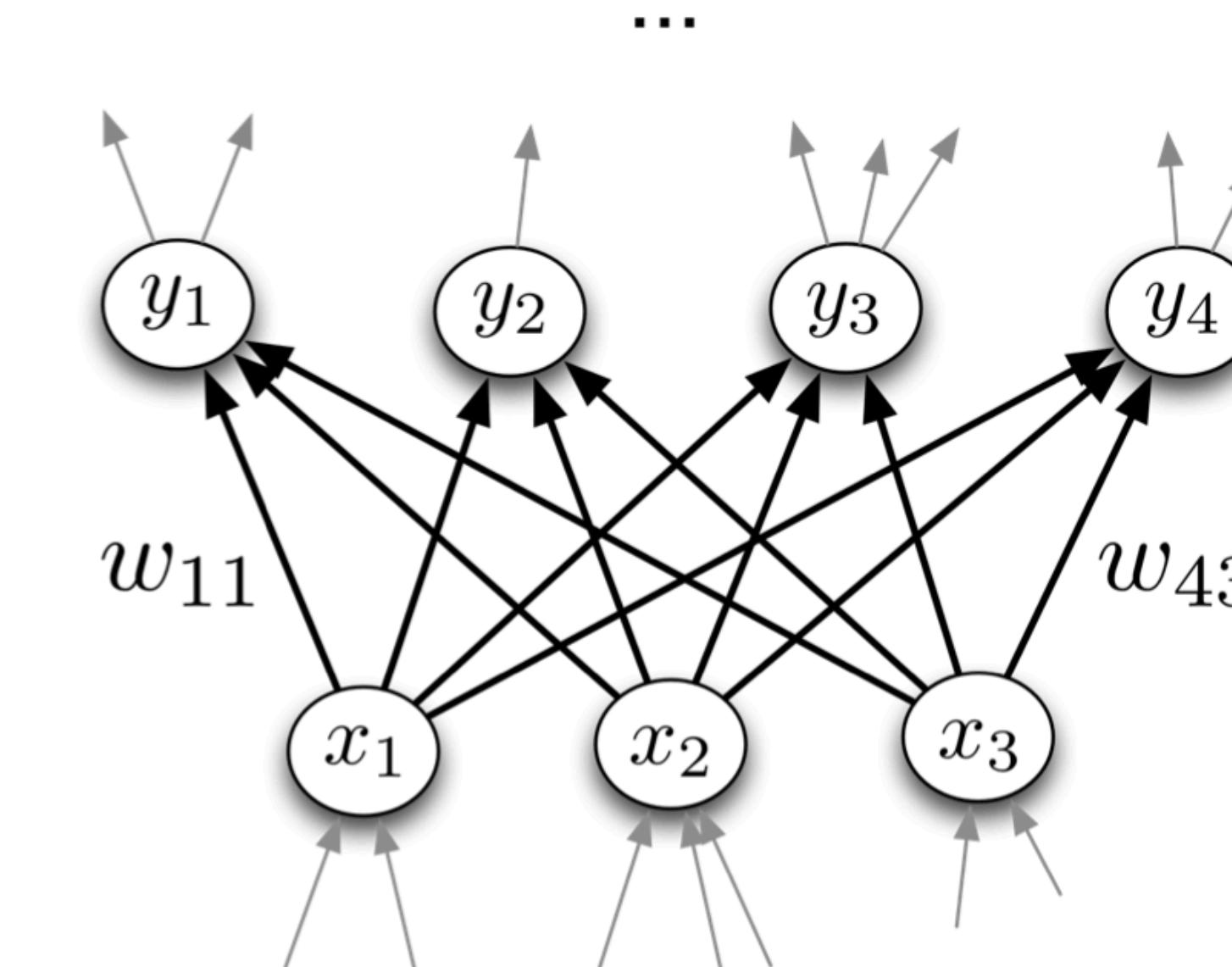
Multilayer Perceptrons

- We can connect lots of units together into a directed acyclic graph.
- This gives a **feed-forward neural network**. That's in contrast to **recurrent neural networks**, which can have cycles. (We'll talk about those later.)
- Typically, units are grouped together into layers.



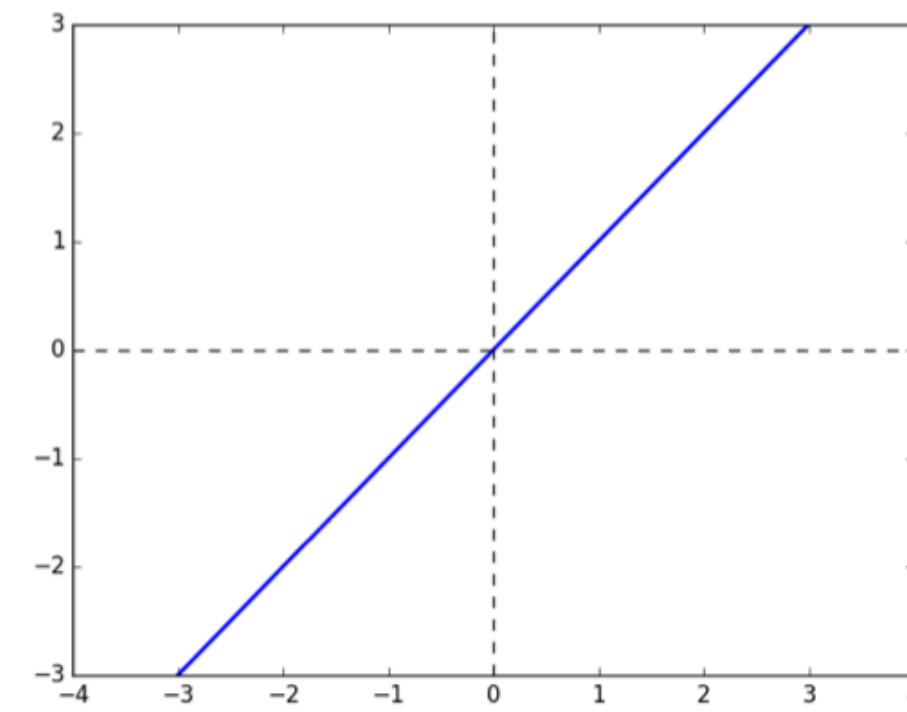
Multilayer Perceptrons

- Each layer connects N input units to M output units.
 - In the simplest case, all input units are connected to all output units. We call this a **fully connected layer**. We'll consider other layer types later.
 - Note: the inputs and outputs for a layer are distinct from the inputs and outputs to the network.
-
- Recall from softmax regression: this means we need an $M \times N$ weight matrix.
 - The output units are a function of the input units:
$$\mathbf{y} = f(\mathbf{x}) = \phi(\mathbf{Wx} + \mathbf{b})$$
 - A multilayer network consisting of fully connected layers is called a **multilayer perceptron**. Despite the name, it has nothing to do with perceptrons!



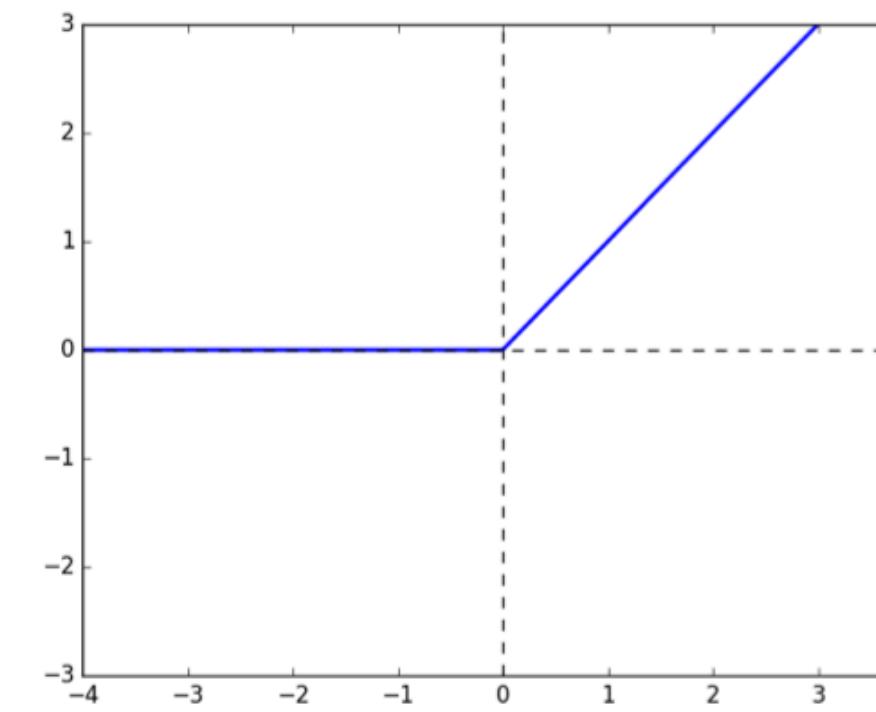
Multilayer Perceptrons

Some activation functions:



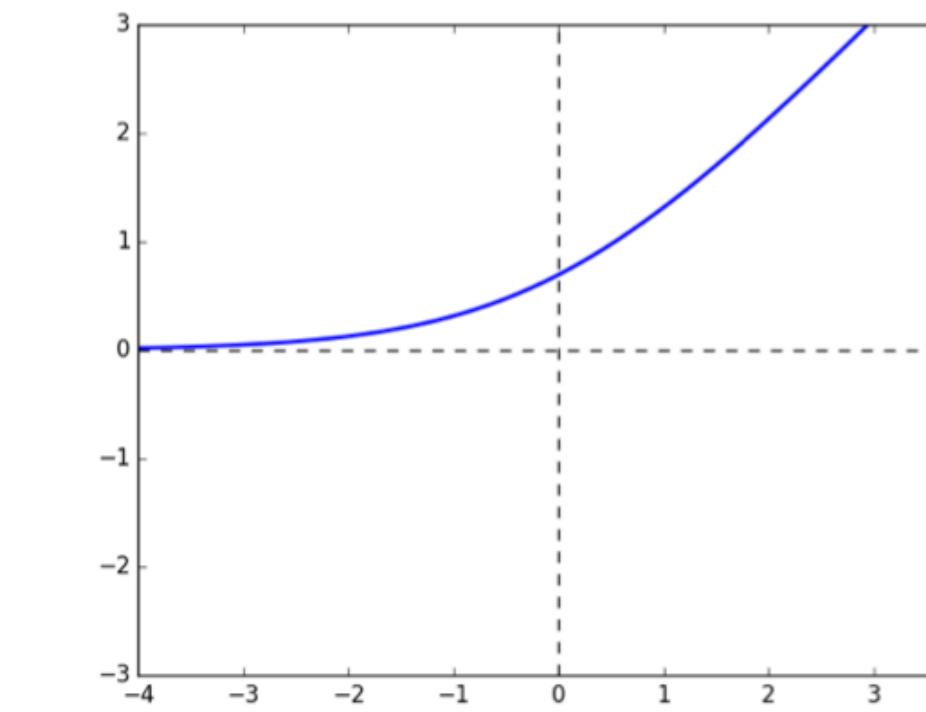
Linear

$$y = z$$



**Rectified Linear Unit
(ReLU)**

$$y = \max(0, z)$$

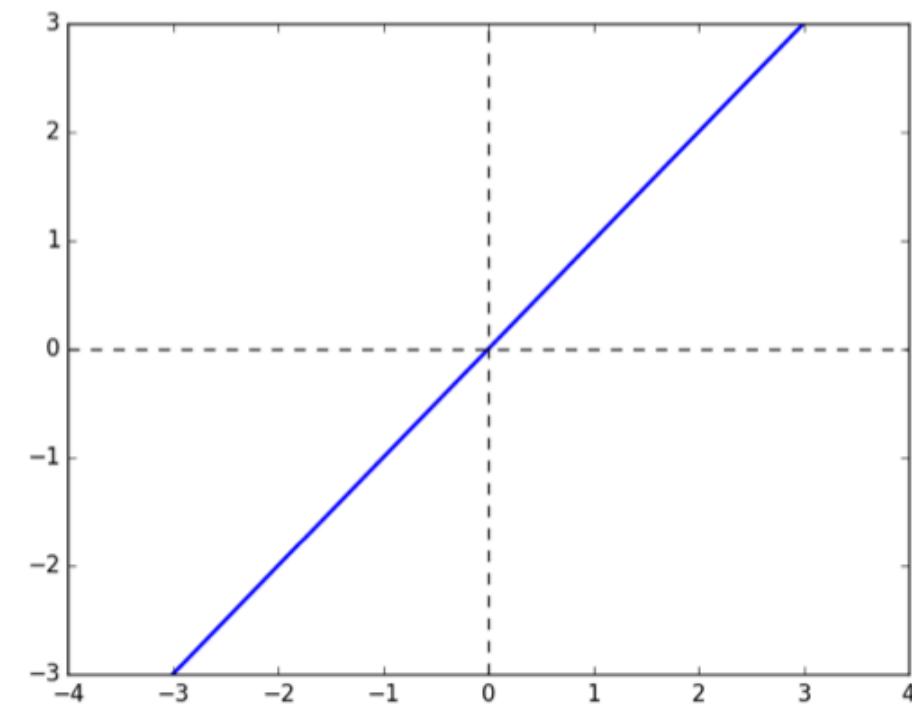


Soft ReLU

$$y = \log(1 + e^z)$$

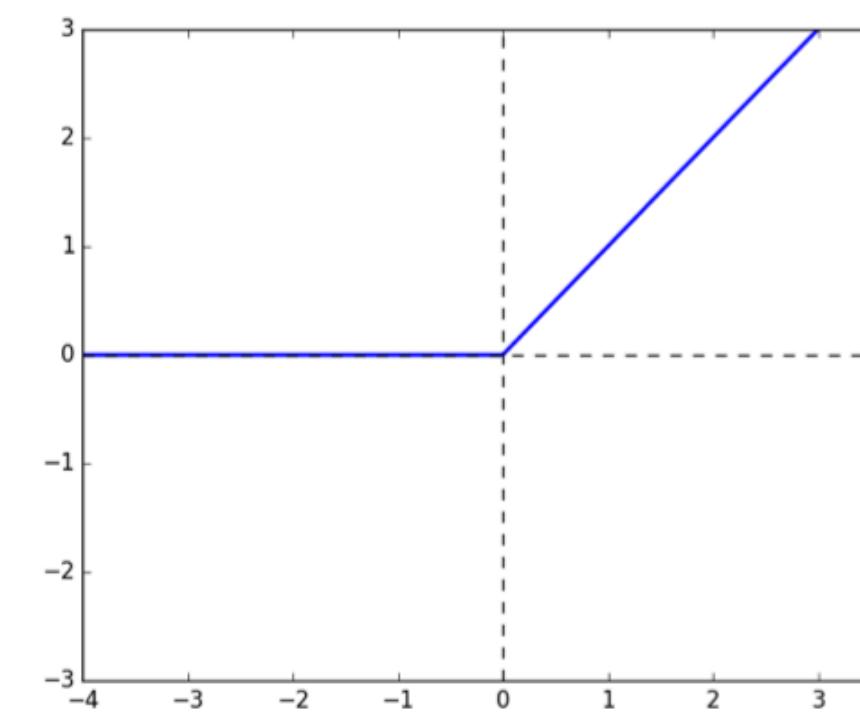
Multilayer Perceptrons

Some activation functions:



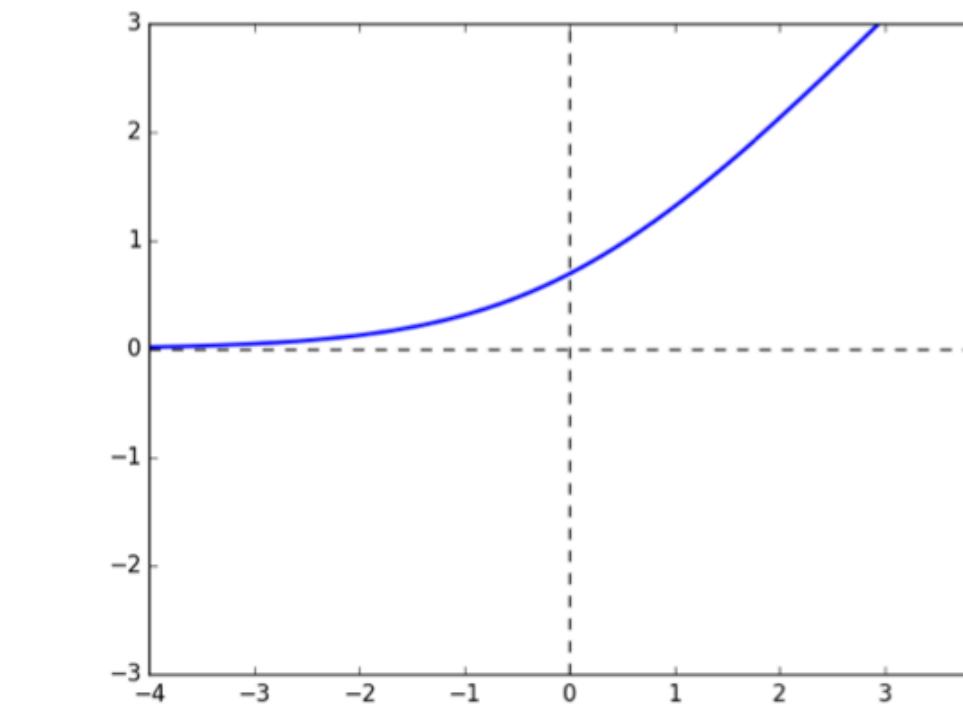
Linear

$$y = z$$



**Rectified Linear Unit
(ReLU)**

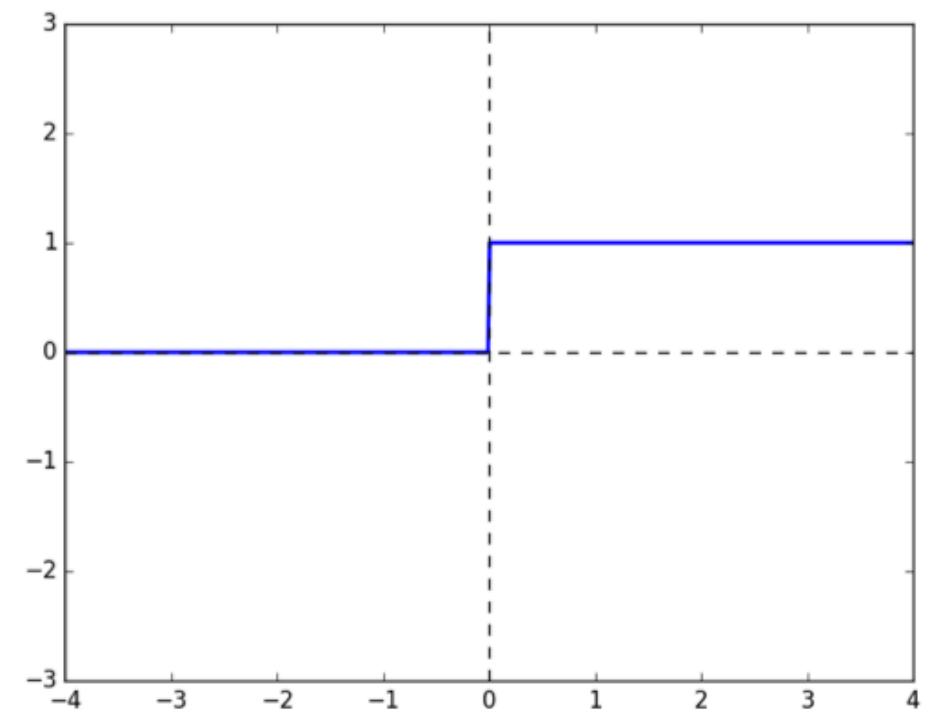
$$y = \max(0, z)$$



Soft ReLU

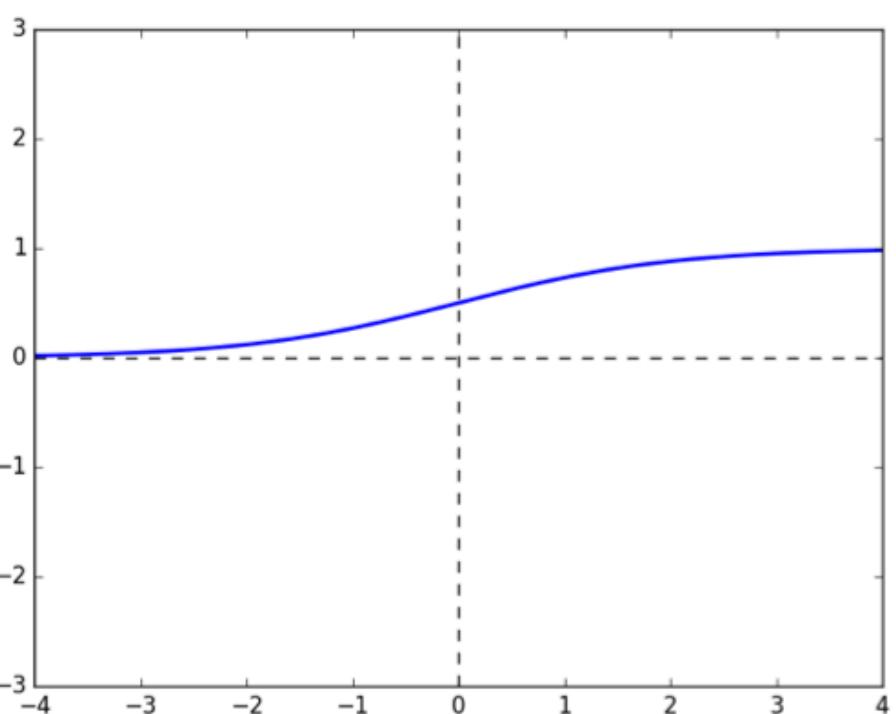
$$y = \log(1 + e^z)$$

Some activation functions:



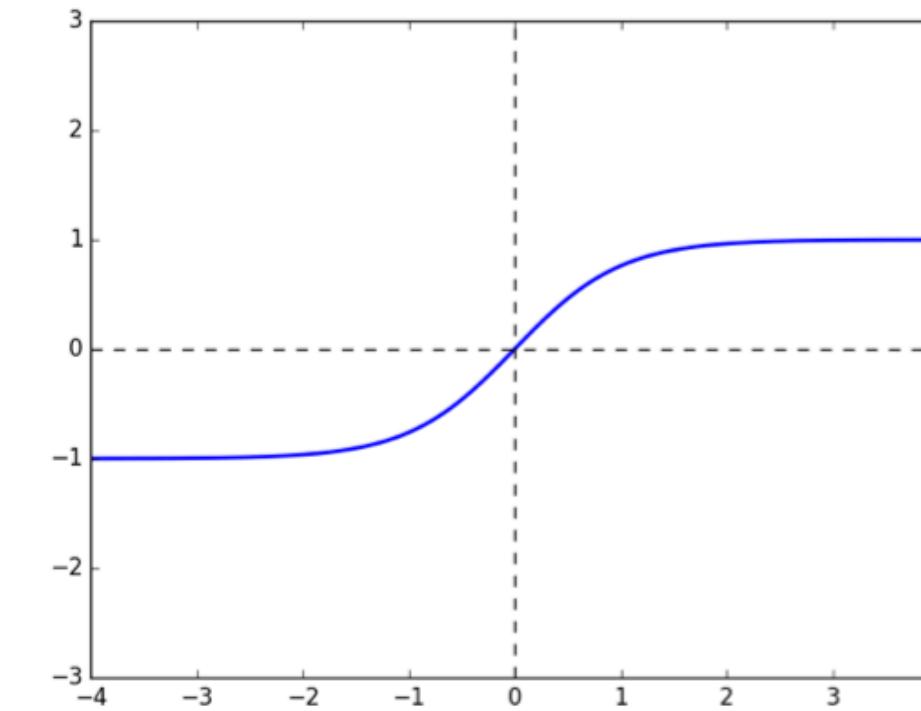
Hard Threshold

$$y = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$



Logistic

$$y = \frac{1}{1 + e^{-z}}$$



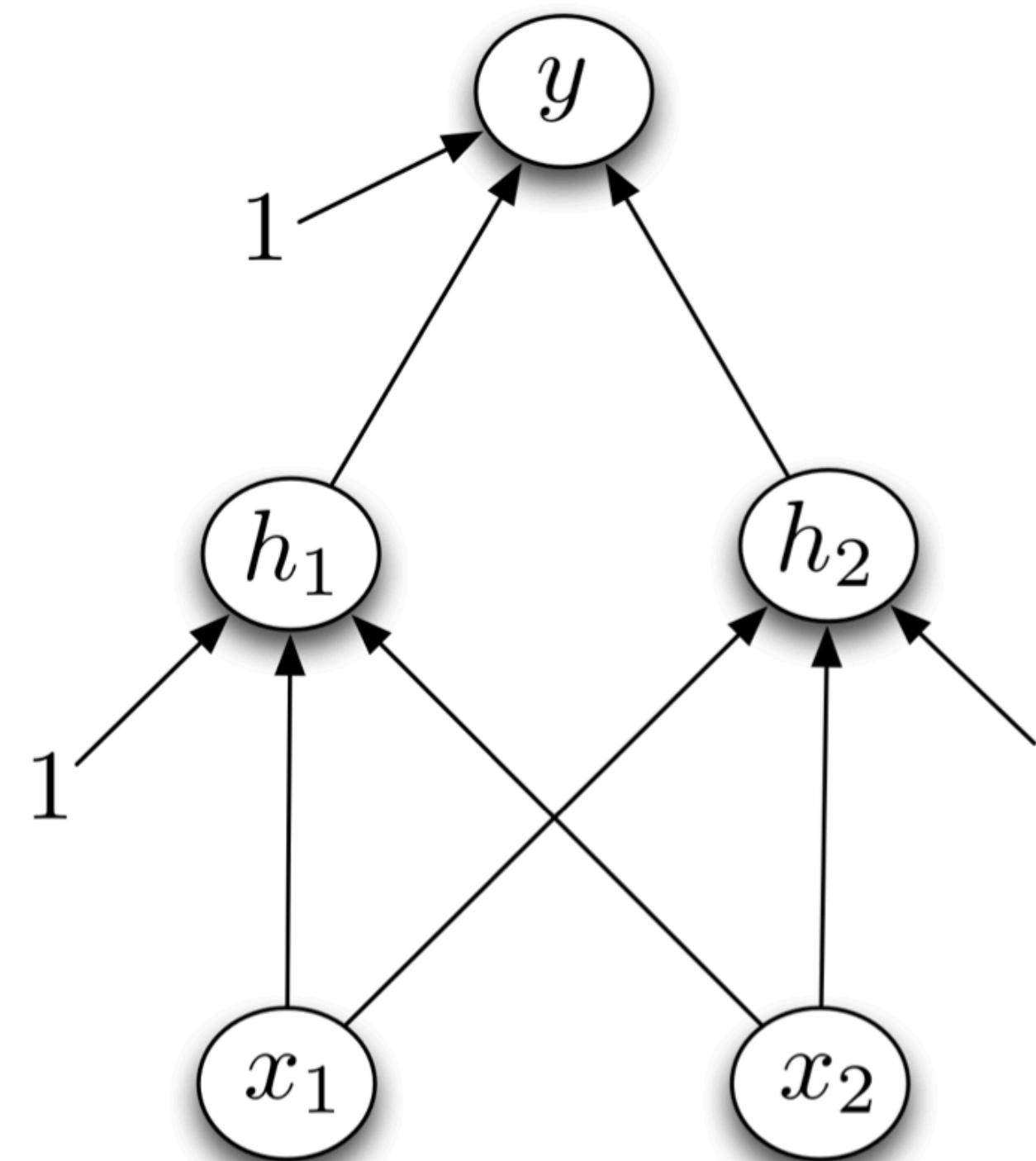
**Hyperbolic Tangent
(tanh)**

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

Multilayer Perceptrons

Designing a network to compute XOR:

Assume hard threshold activation function



Multilayer Perceptrons

