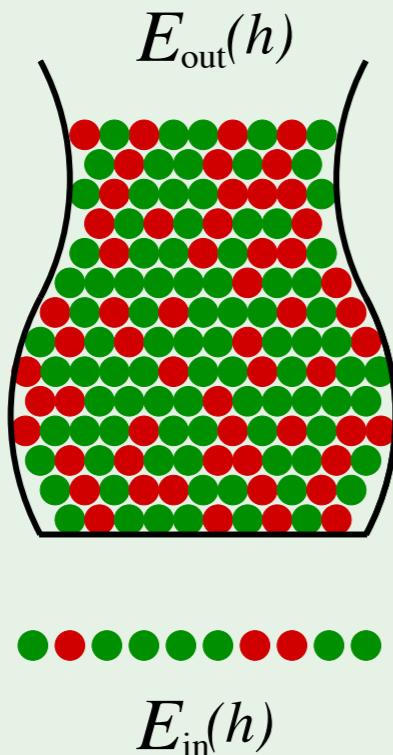


## Review of Lecture 2

Is Learning feasible?

Yes, in a **probabilistic** sense.



$$\mathbb{P} [ |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon ] \leq 2e^{-2\epsilon^2 N}$$

Since  $g$  has to be one of  $h_1, h_2, \dots, h_M$ , we conclude that

If:

$$|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon$$

Then:

$$|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon \quad \text{or}$$

$$|E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon \quad \text{or}$$

...

$$|E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| > \epsilon$$

This gives us an added **M** factor.

# Learning From Data

Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 3: Linear Models I

Sponsored by Caltech's Provost Office, E&AS Division, and IST

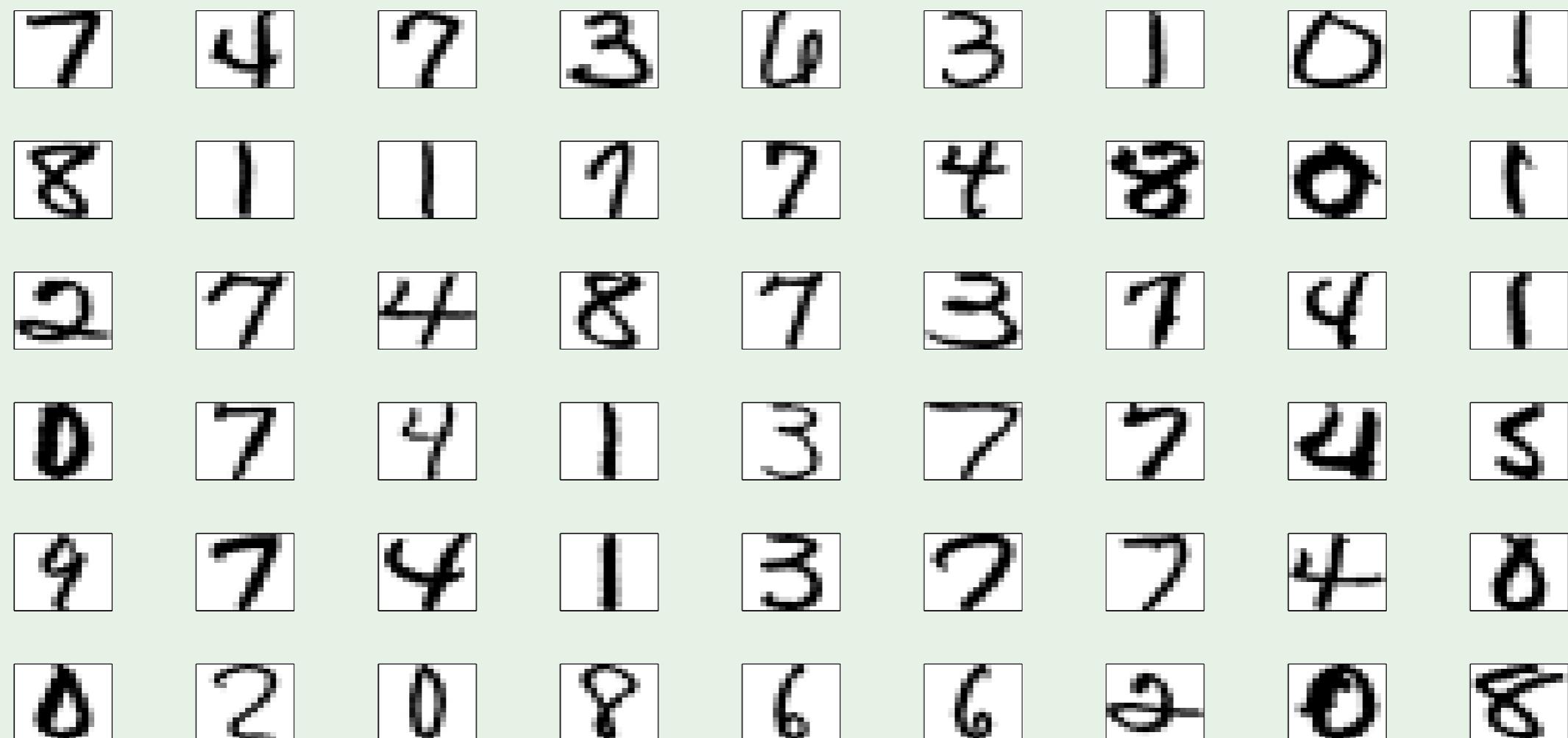


Tuesday, April 10, 2012

# Outline

- Input representation
- Linear Classification
- Linear Regression
- Nonlinear Transformation

## A real data set



## Input representation

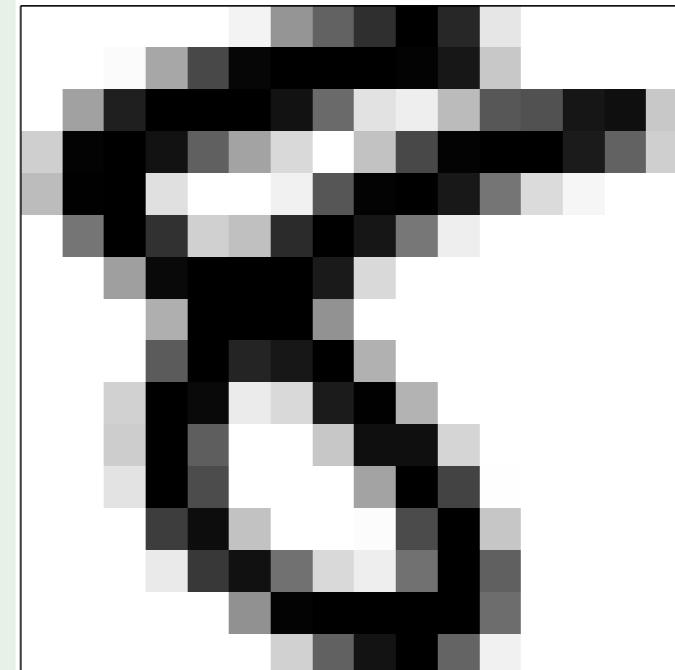
'raw' input  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_{256})$

linear model:  $(w_0, w_1, w_2, \dots, w_{256})$

**Features:** Extract useful information, e.g.,

intensity and symmetry  $\mathbf{x} = (x_0, x_1, x_2)$

linear model:  $(w_0, w_1, w_2)$

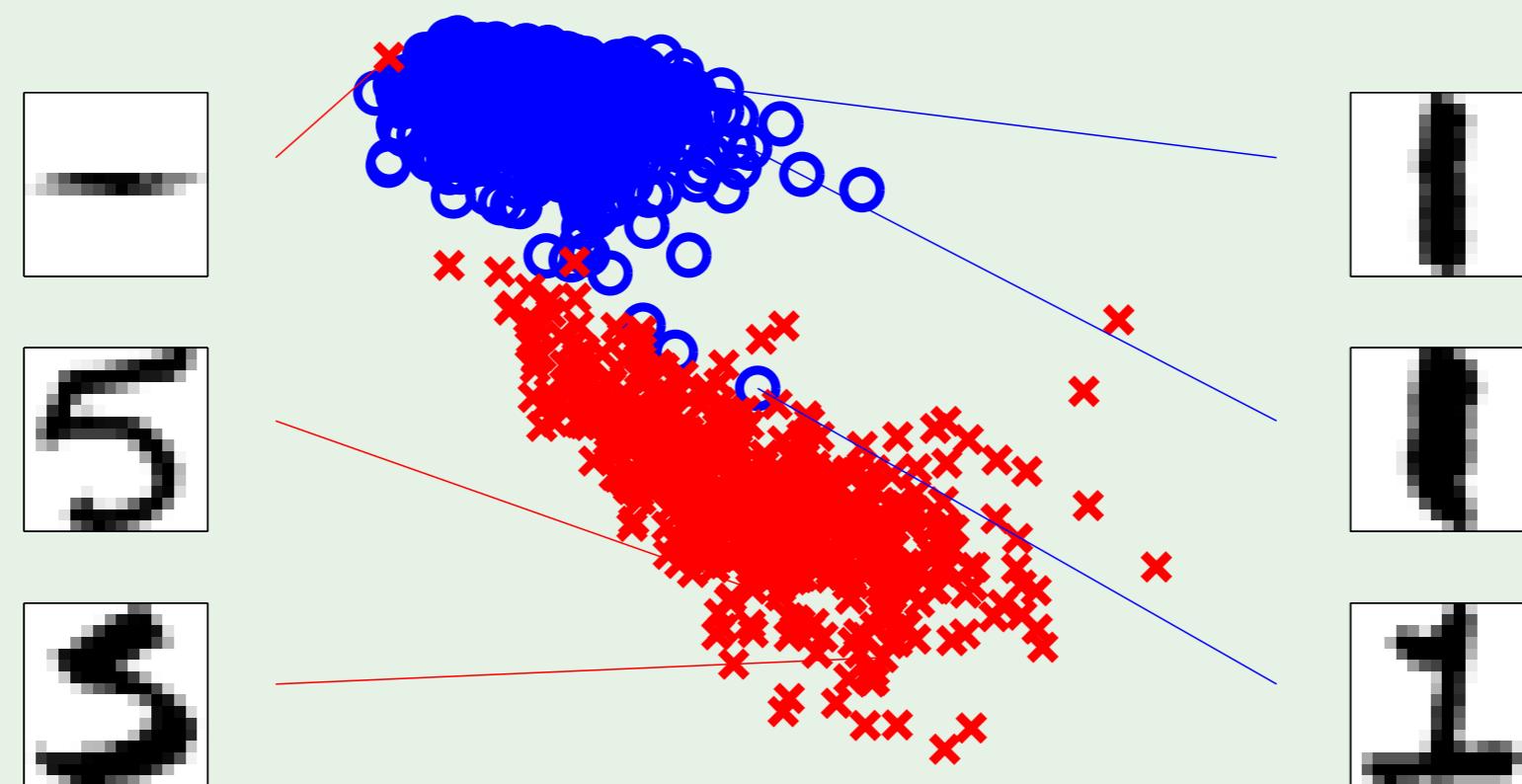


## Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

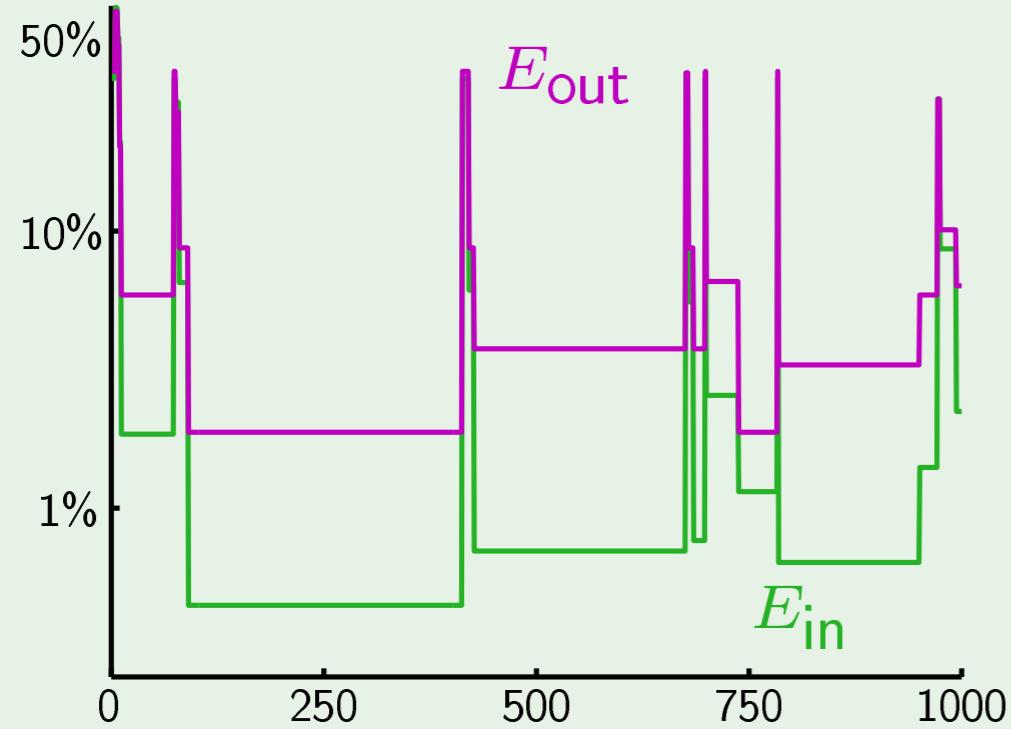
$x_1$ : intensity

$x_2$ : symmetry

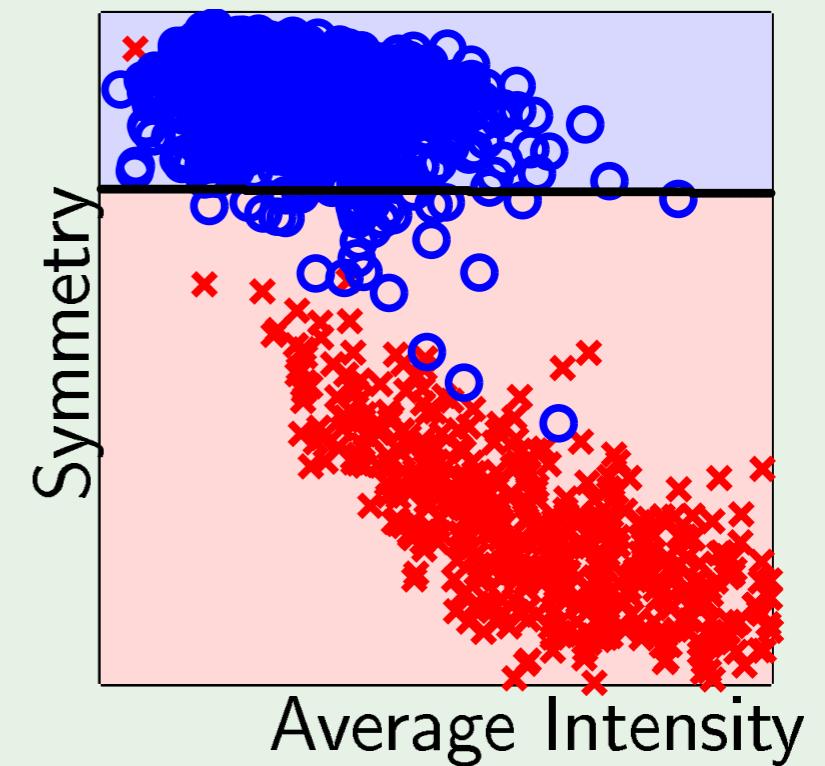


## What PLA does

Evolution of  $E_{\text{in}}$  and  $E_{\text{out}}$

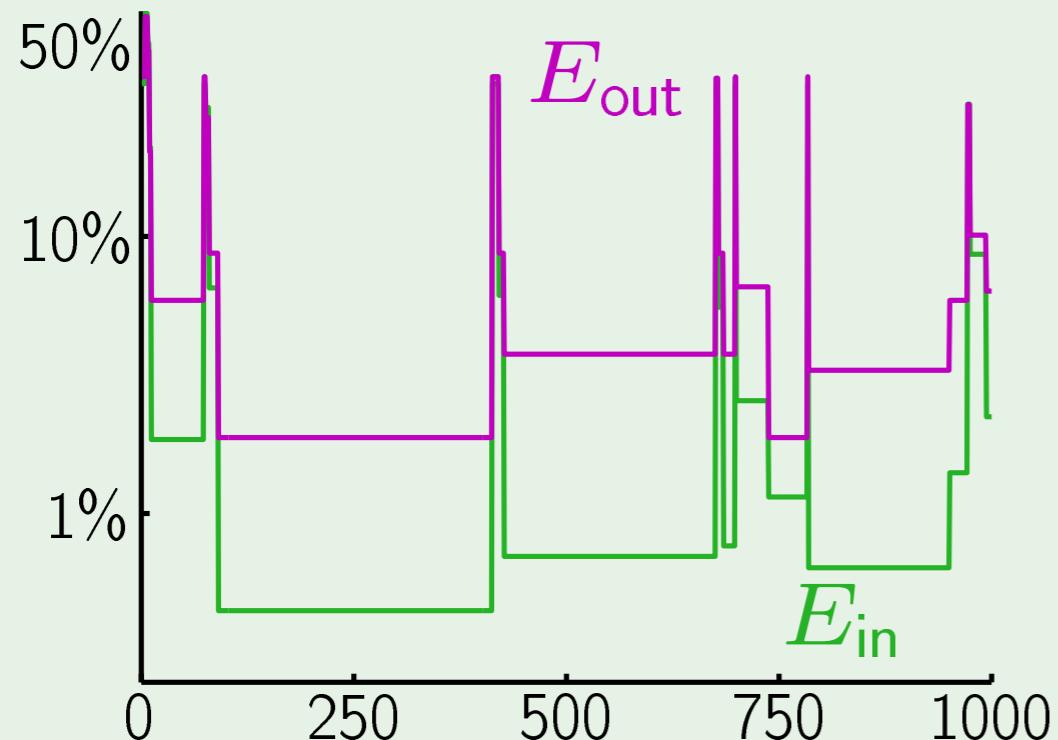


Final perceptron boundary

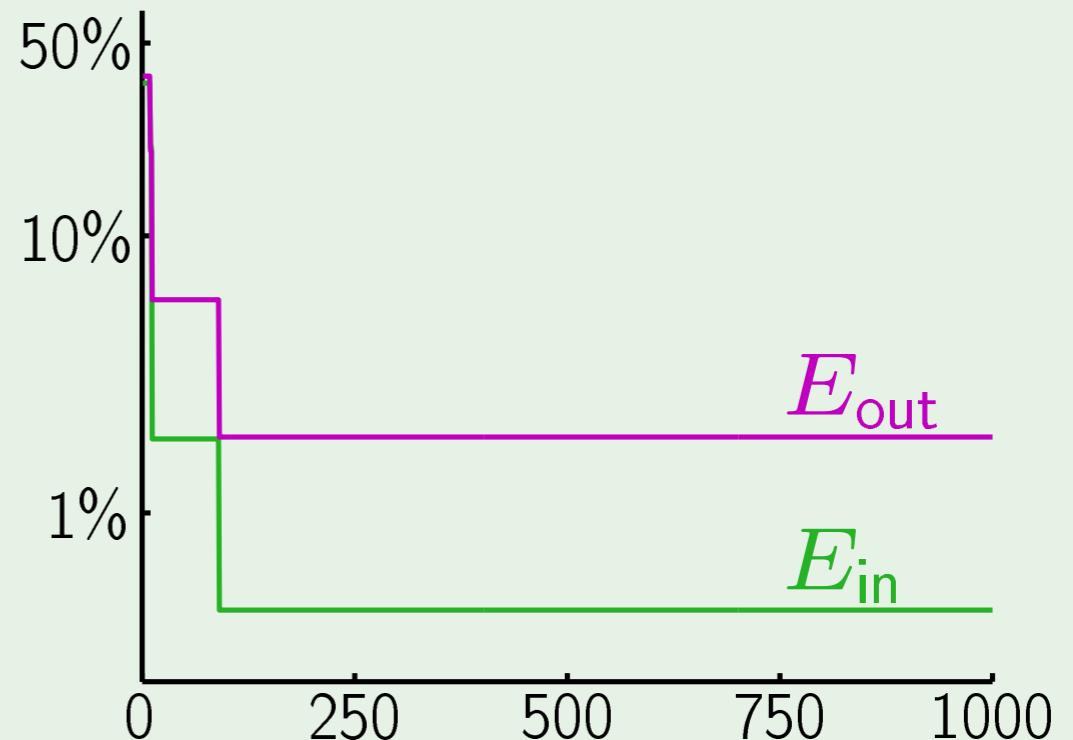


## The ‘pocket’ algorithm

PLA:

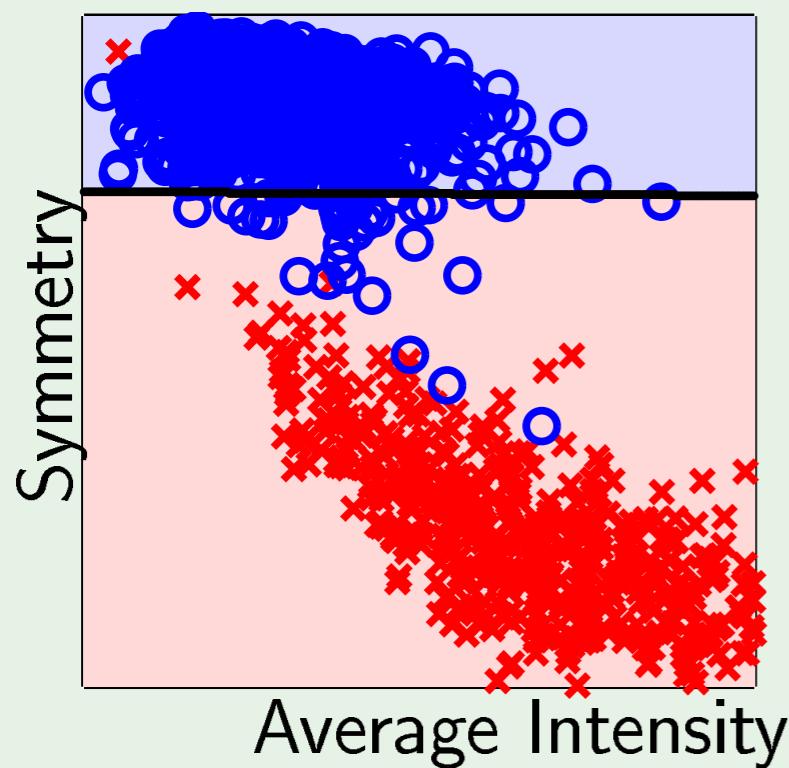


Pocket:

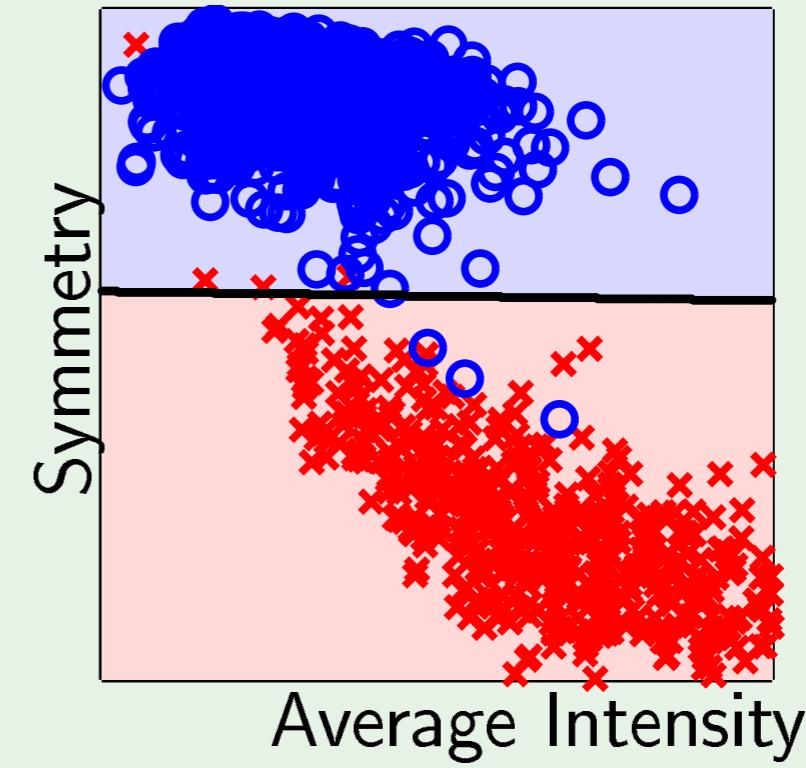


## Classification boundary - PLA versus Pocket

PLA:



Pocket:



## Outline

- Input representation
- Linear Classification
- Linear Regression **regression  $\equiv$  real-valued output**
- Nonlinear Transformation

## Credit again

**Classification:** Credit approval (yes/no)

**Regression:** Credit line (dollar amount)

Input:  $\mathbf{x} =$

age	23 years
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

Linear regression output:  $h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^\top \mathbf{x}$

## The data set

Credit officers decide on credit lines:

$$(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$$

$\mathbf{y}_n \in \mathbb{R}$  is the credit line for customer  $\mathbf{x}_n$ .

Linear regression tries to replicate that.

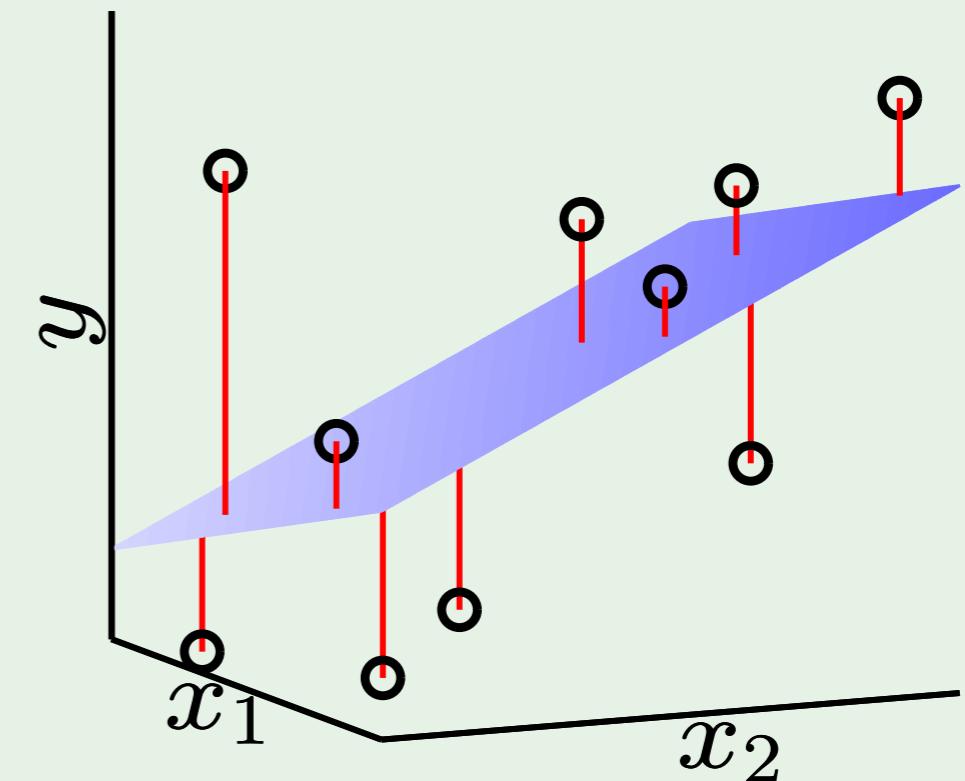
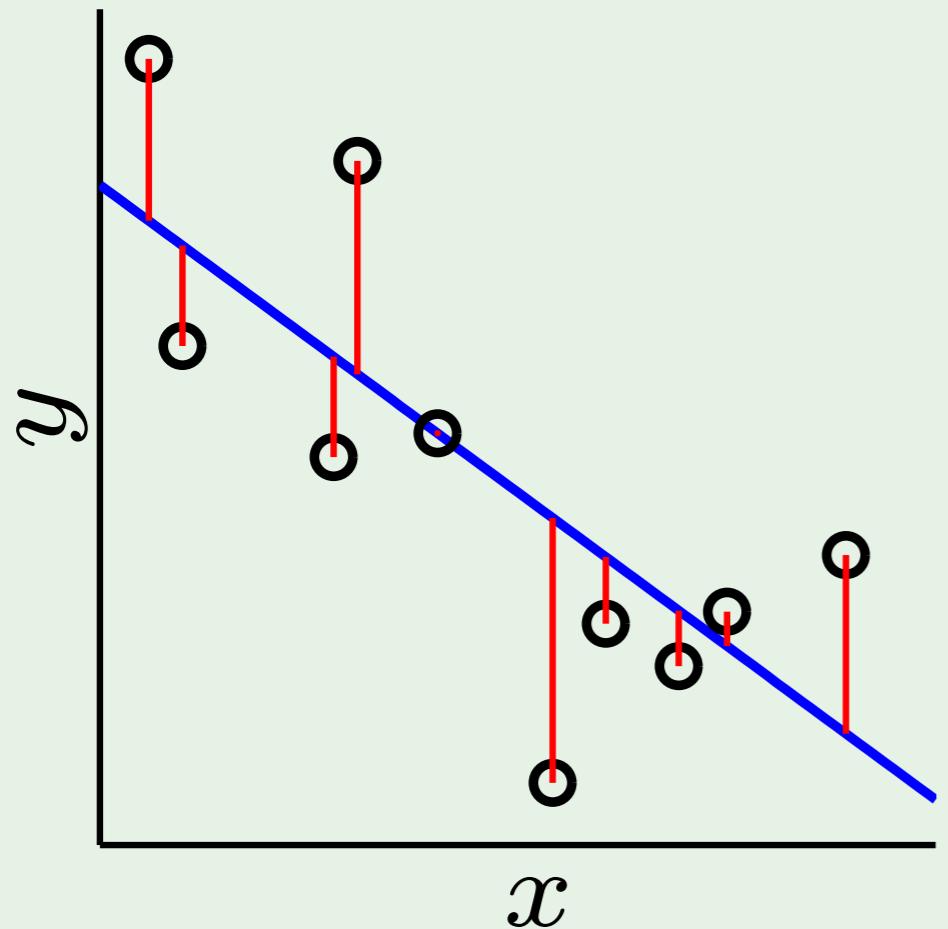
## How to measure the error

How well does  $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  approximate  $f(\mathbf{x})$ ?

In linear regression, we use squared error  $(h(\mathbf{x}) - f(\mathbf{x}))^2$

**in-sample error:**  $E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$

## Illustration of linear regression



## The expression for $E_{\text{in}}$

$$\begin{aligned} E_{\text{in}}(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 \\ &= \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

where

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

## Minimizing $E_{\text{in}}$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} \quad \text{where} \quad \mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$\mathbf{X}^\dagger$  is the '**pseudo-inverse**' of  $\mathbf{X}$

## The pseudo-inverse

$$X^\dagger = (X^T X)^{-1} X^T$$

$$\left( \underbrace{\begin{bmatrix} & \\ & d+1 \times d+1 \end{bmatrix}}_{d+1 \times N} \right)^{-1} \underbrace{\begin{bmatrix} & \\ & d+1 \times N \end{bmatrix}}_{d+1 \times N}$$

## The linear regression algorithm

- 1: Construct the matrix  $\mathbf{X}$  and the vector  $\mathbf{y}$  from the data set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  as follows

$$\mathbf{X} = \underbrace{\begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}}_{\text{input data matrix}}, \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse  $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .
- 3: Return  $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$ .

# Linear regression for classification

Linear regression learns a real-valued function  $y = f(\mathbf{x}) \in \mathbb{R}$

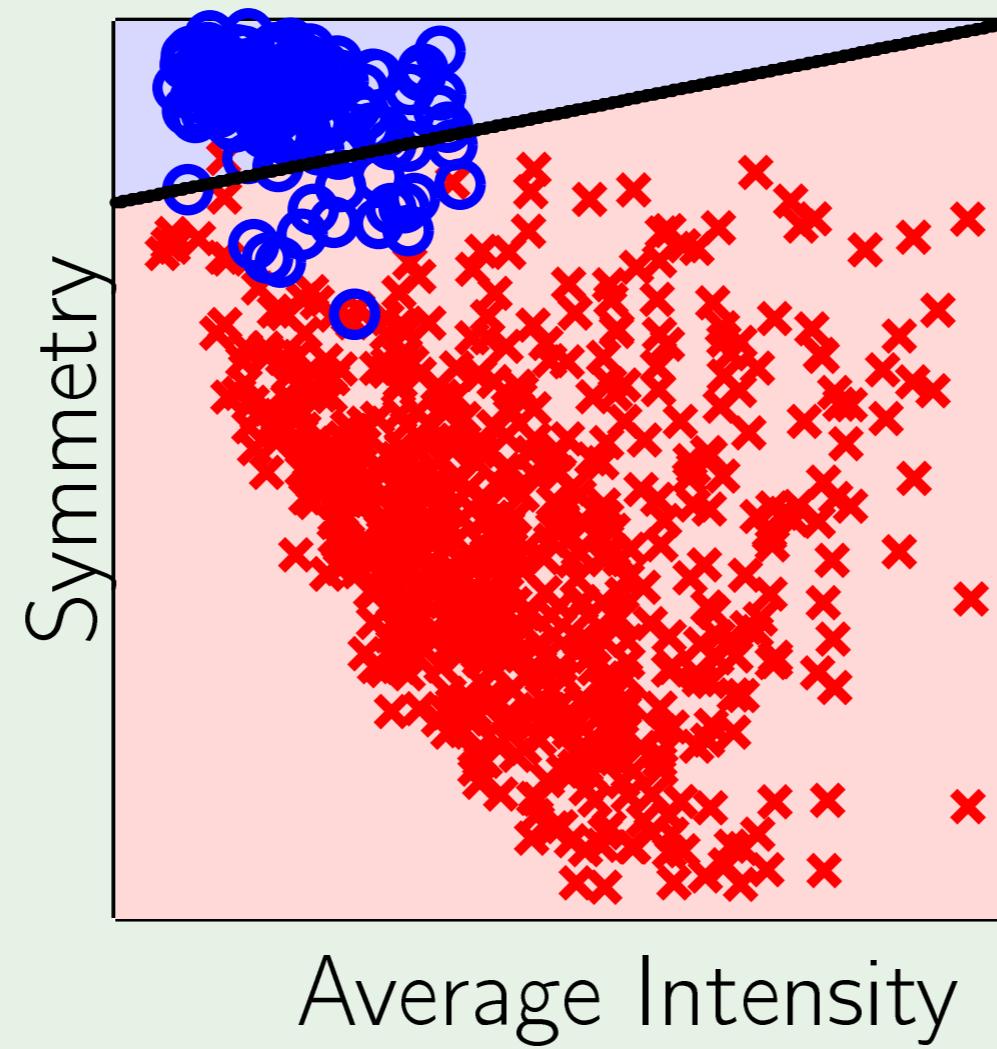
Binary-valued functions are also real-valued!  $\pm 1 \in \mathbb{R}$

Use linear regression to get  $\mathbf{w}$  where  $\mathbf{w}^\top \mathbf{x}_n \approx y_n = \pm 1$

In this case,  $\text{sign}(\mathbf{w}^\top \mathbf{x}_n)$  is likely to agree with  $y_n = \pm 1$

Good initial weights for classification

## Linear regression boundary

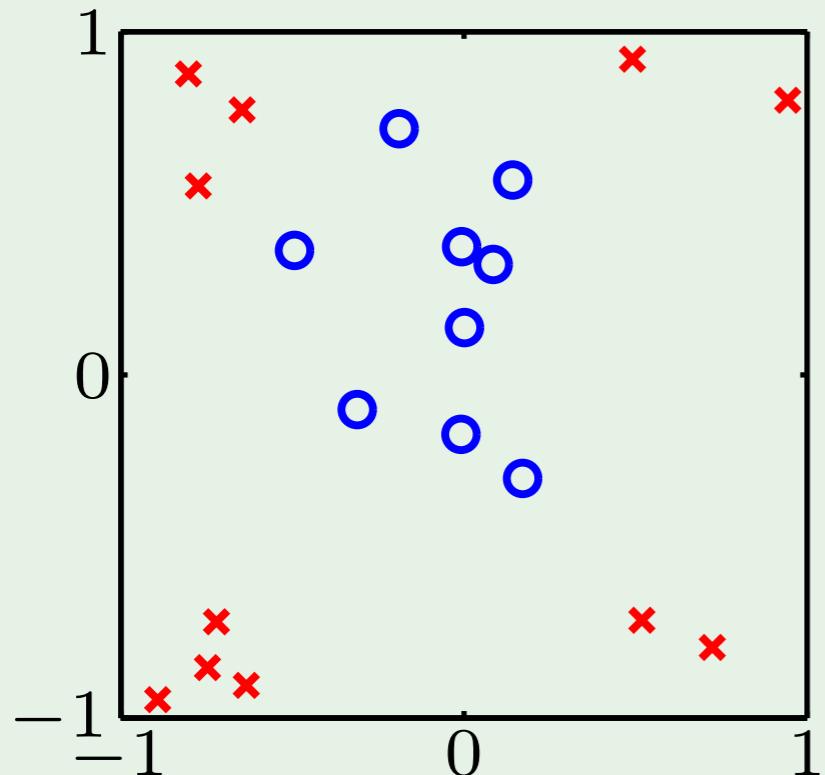


## Outline

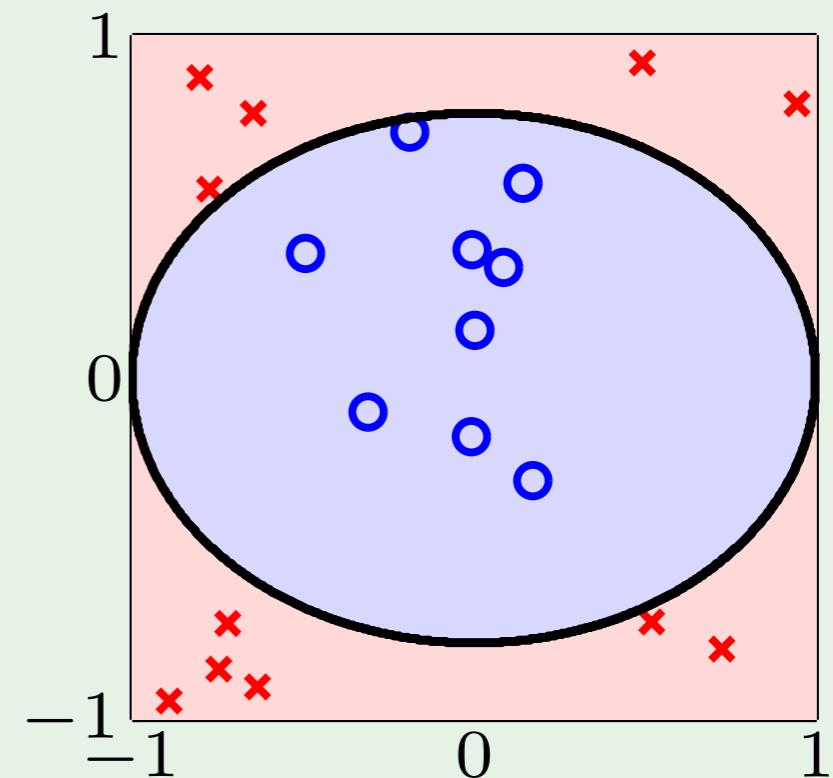
- Input representation
- Linear Classification
- Linear Regression
- Nonlinear Transformation

## Linear is limited

Data:



Hypothesis:



## Another example

Credit line is affected by ‘years in residence’

but **not** in a linear way!

Nonlinear  $[[x_i < 1]]$  and  $[[x_i > 5]]$  are better.

Can we do that with linear models?

# Linear in what?

Linear regression implements

$$\sum_{i=0}^d \textcolor{red}{w}_i x_i$$

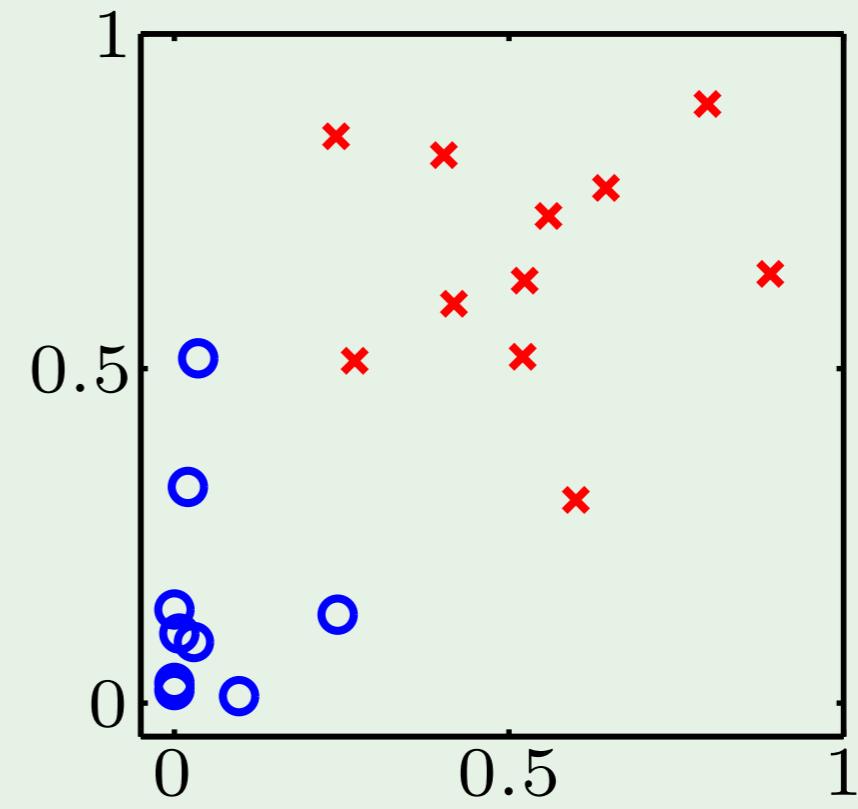
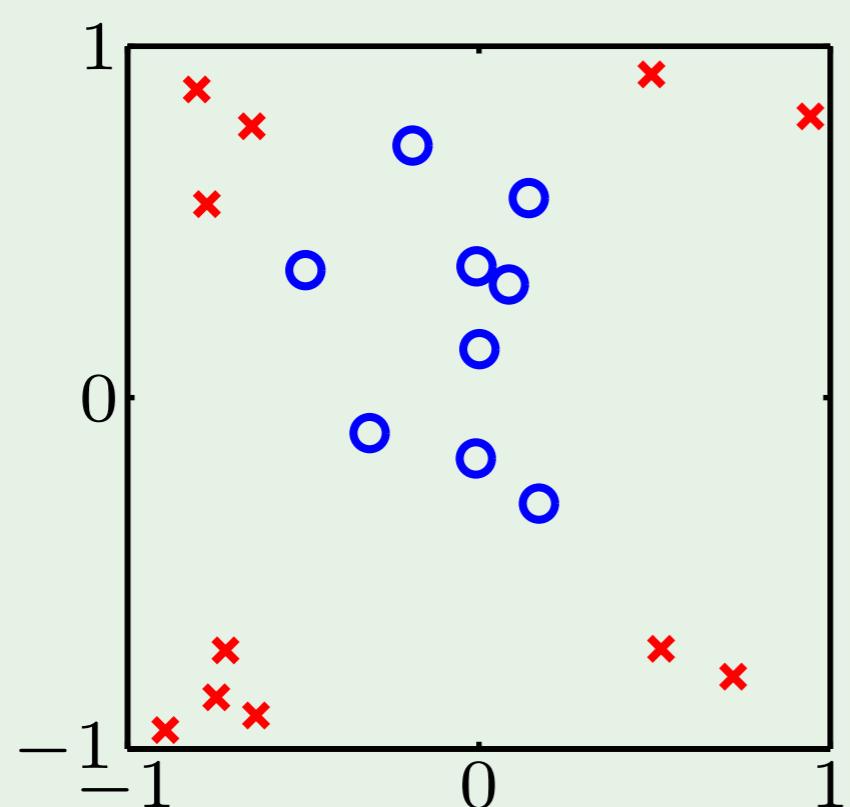
Linear classification implements

$$\text{sign} \left( \sum_{i=0}^d \textcolor{red}{w}_i x_i \right)$$

Algorithms work because of **linearity in the weights**

## Transform the data nonlinearly

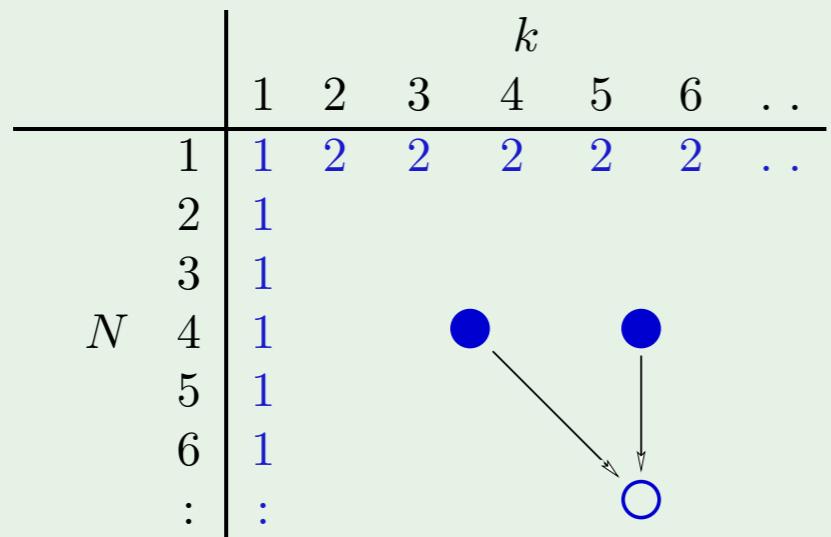
$$(x_1, x_2) \xrightarrow{\Phi} (x_1^2, x_2^2)$$



## Review of Lecture 6

- $m_{\mathcal{H}}(N)$  is polynomial

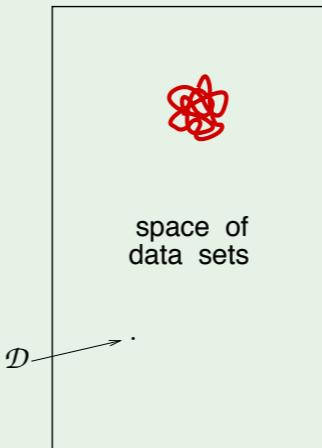
if  $\mathcal{H}$  has a break point  $k$



$$m_{\mathcal{H}}(N) \leq \underbrace{\sum_{i=0}^{k-1} \binom{N}{i}}_{\text{maximum power is } N^{k-1}}$$

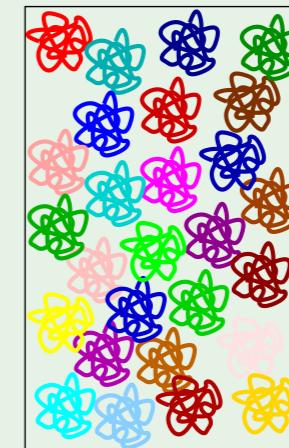
- The VC Inequality

Hoeffding Inequality



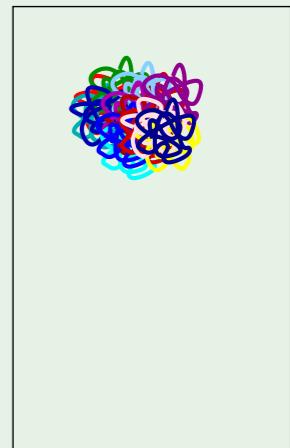
(a)

Union Bound



(b)

VC Bound



(c)

$$\mathbb{P} [ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon ] \leq 2M e^{-2\epsilon^2 N}$$



$$\mathbb{P} [ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon ] \leq 4 \boxed{m_{\mathcal{H}}(2N)} e^{-\frac{1}{8}\epsilon^2 N}$$

## 2. Number of data points needed

Two small quantities in the VC inequality:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq \underbrace{4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\epsilon^2 N}}_{\delta}$$

If we want certain  $\epsilon$  and  $\delta$ , how does  $N$  depend on  $d_{\text{VC}}$ ?

Let us look at

$$N^{\delta} e^{-N}$$

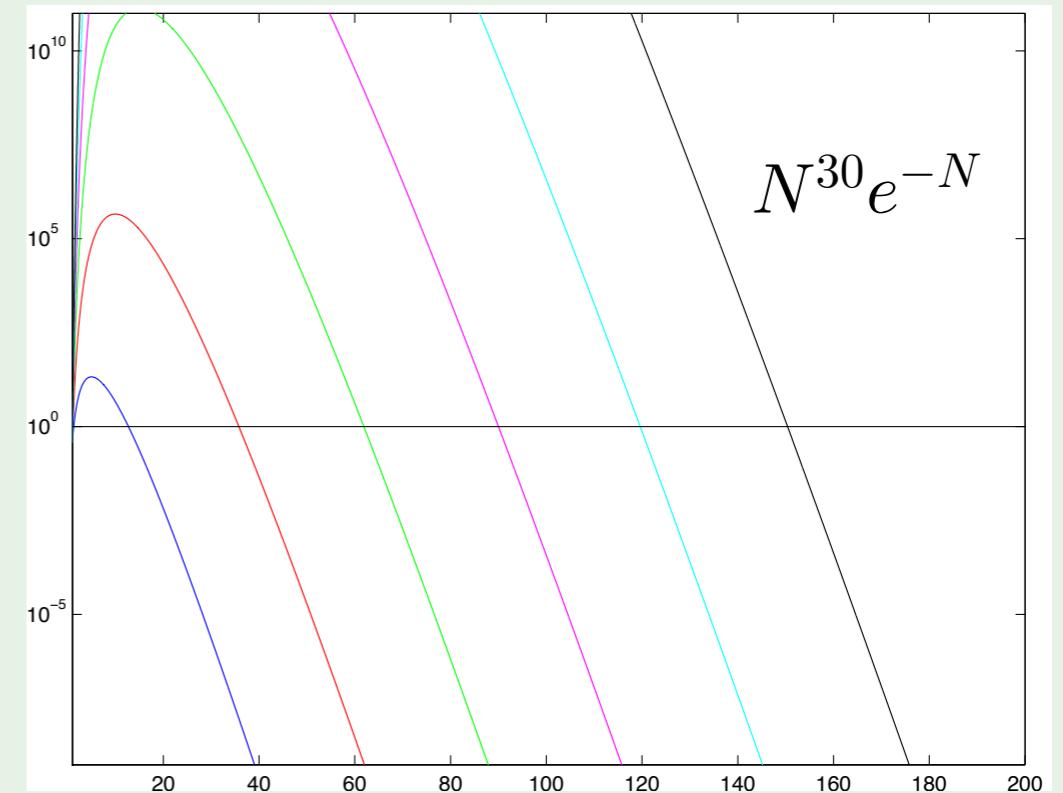
$$N^{\textcolor{red}{d}} e^{-N}$$

Fix  $N^{\textcolor{red}{d}} e^{-N} = \text{small value}$

How does  $N$  change with  $\textcolor{red}{d}$ ?

**Rule of thumb:**

$$N \geq 10 \textcolor{red}{d}_{\text{VC}}$$

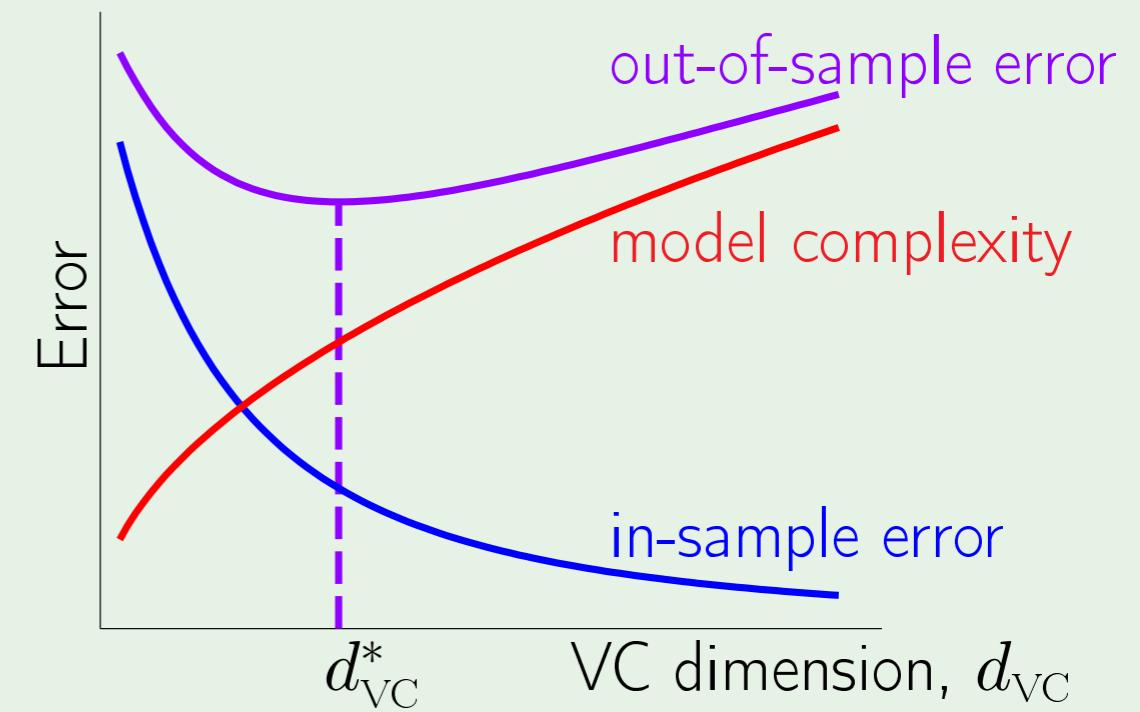


## What the theory will achieve

Characterizing the feasibility of learning for infinite  $M$

Characterizing the tradeoff:

Model complexity	$\uparrow$	$E_{\text{in}}$	$\downarrow$
Model complexity	$\uparrow$	$E_{\text{out}} - E_{\text{in}}$	$\uparrow$



# Where we are

- Linear classification ✓
- Linear regression ✓
- Logistic regression
- Nonlinear transforms ✎

## Nonlinear transforms

$$\mathbf{x} = (x_0, x_1, \dots, x_d) \xrightarrow{\Phi} \mathbf{z} = (z_0, z_1, \dots, \dots, z_{\tilde{d}})$$

Each  $z_i = \phi_i(\mathbf{x})$        $\mathbf{z} = \Phi(\mathbf{x})$

Example:  $\mathbf{z} = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$

Final hypothesis  $g(\mathbf{x})$  in  $\mathcal{X}$  space:

$$\text{sign} (\tilde{\mathbf{w}}^\top \Phi(\mathbf{x})) \quad \text{or} \quad \tilde{\mathbf{w}}^\top \Phi(\mathbf{x})$$

## The price we pay

$$\mathbf{x} = (x_0, x_1, \dots, x_d) \xrightarrow{\Phi} \mathbf{z} = (z_0, z_1, \dots, \dots, z_{\tilde{d}})$$



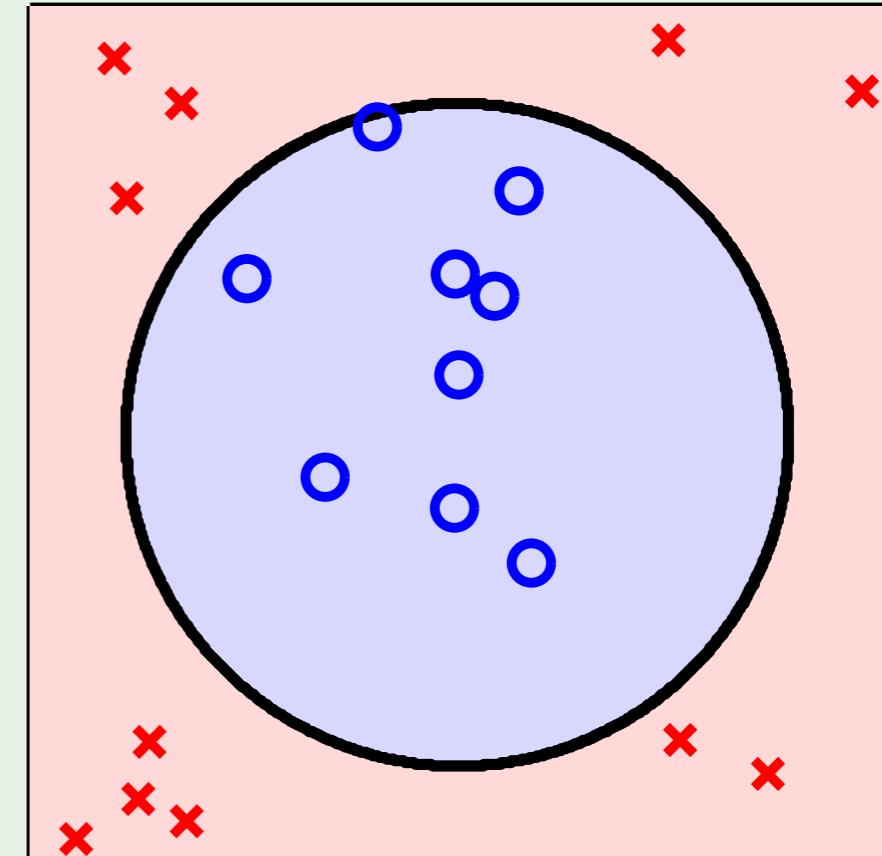
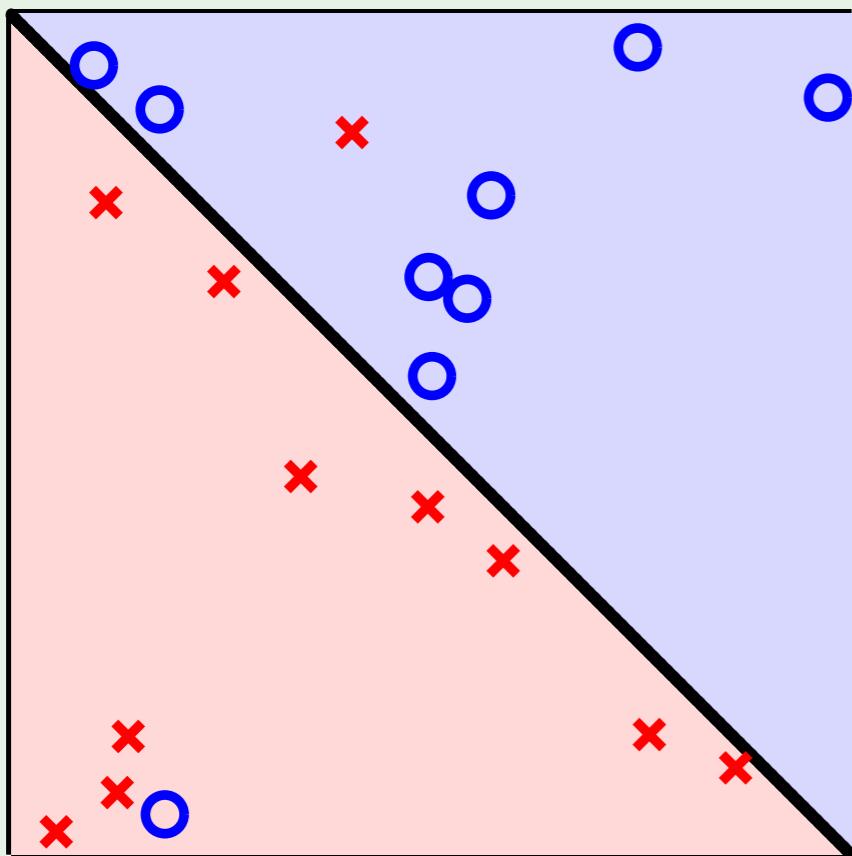
**w**

**w̃**

$$d_{\text{VC}} = d + 1$$

$$d_{\text{VC}} \leq \tilde{d} + 1$$

## Two non-separable cases

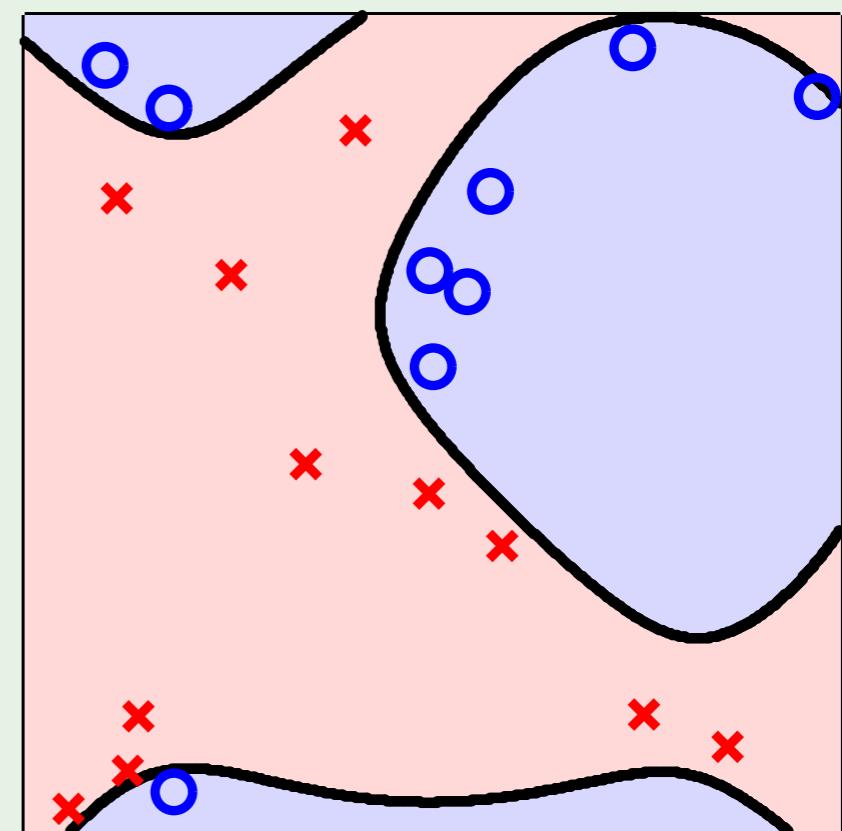


## First case

Use a linear model in  $\mathcal{X}$ ; accept  $E_{\text{in}} > 0$

or

Insist on  $E_{\text{in}} = 0$ ; go to high-dimensional  $\mathcal{Z}$



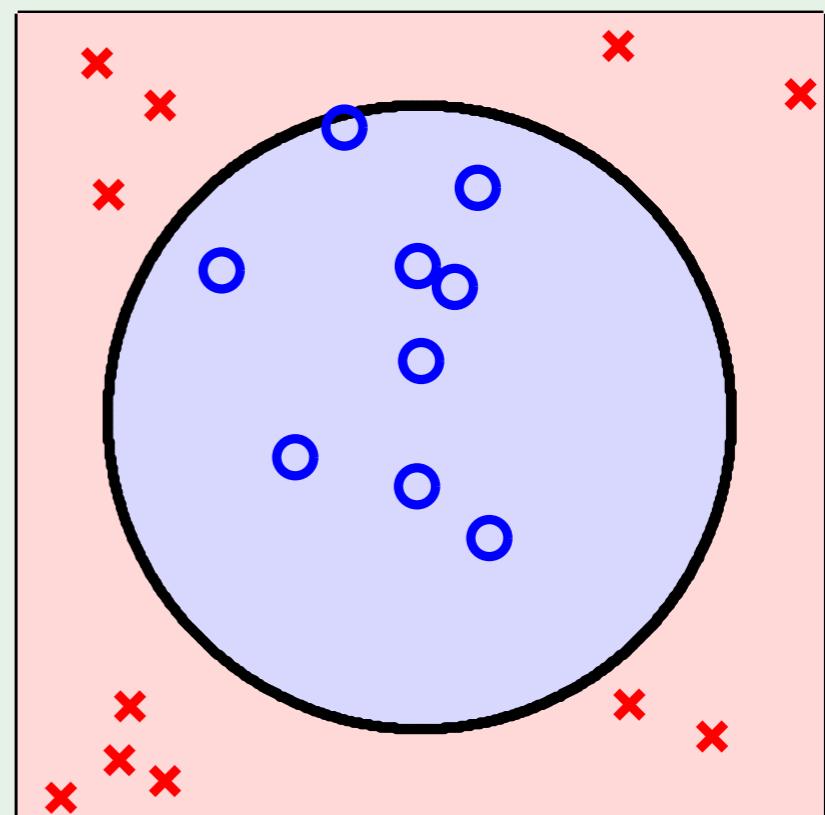
## Second case

$$\mathbf{z} = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$$

Why not:  $\mathbf{z} = (1, x_1^2, x_2^2)$

or better yet:  $\mathbf{z} = (1, x_1^2 + x_2^2)$

or even:  $\mathbf{z} = (x_1^2 + x_2^2 - 0.6)$



## Lesson learned

Looking at the data *before* choosing the model can be hazardous to your  $E_{\text{out}}$

### Data snooping



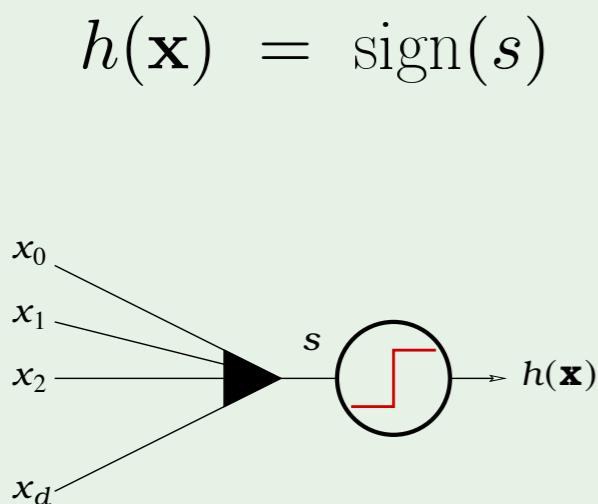
## Logistic regression - Outline

- The model
- Error measure
- Learning algorithm

## A third linear model

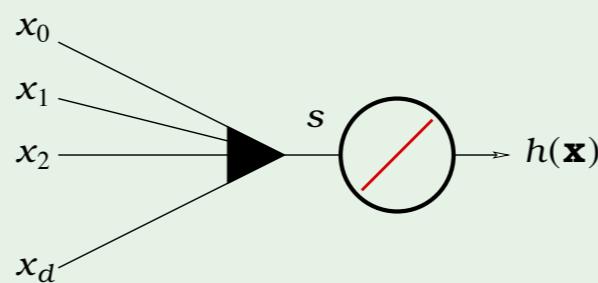
$$s = \sum_{i=0}^d w_i x_i$$

linear classification



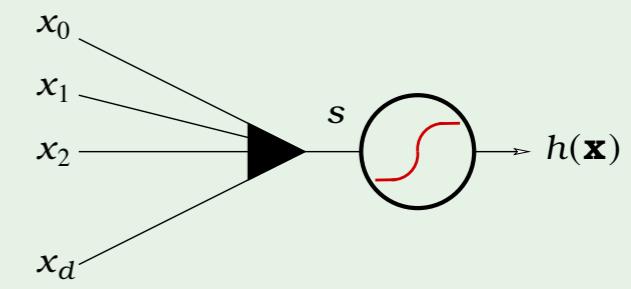
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

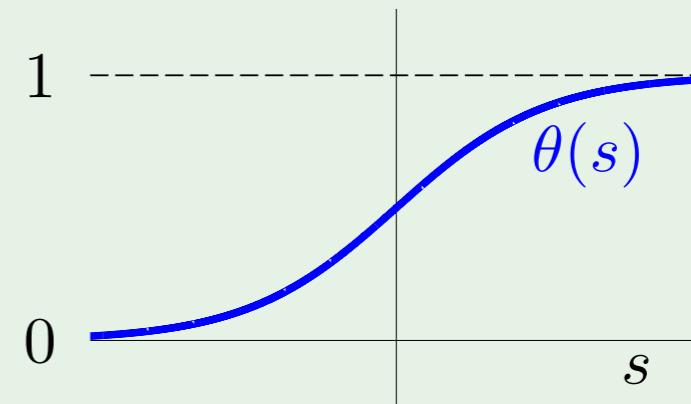
$$h(\mathbf{x}) = \theta(s)$$



## The logistic function $\theta$

The formula:

$$\theta(s) = \frac{e^s}{1 + e^s}$$



soft threshold: uncertainty

sigmoid: flattened out ‘s’

## Probability interpretation

$h(\mathbf{x}) = \theta(s)$  is interpreted as a probability

**Example.** Prediction of heart attacks

Input  $\mathbf{x}$ : cholesterol level, age, weight, etc.

$\theta(s)$ : probability of a heart attack

The signal  $s = \mathbf{w}^T \mathbf{x}$       "risk score"

## Genuine probability

Data  $(\mathbf{x}, y)$  with **binary  $y$** , generated by a noisy target:

$$P(y \mid \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1; \\ 1 - f(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

The target  $f : \mathbb{R}^d \rightarrow [0, 1]$  is the probability

$$\text{Learn } g(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x}) \approx f(\mathbf{x})$$

## Error measure

For each  $(\mathbf{x}, y)$ ,  $y$  is generated by probability  $f(\mathbf{x})$

Plausible error measure based on **likelihood**:

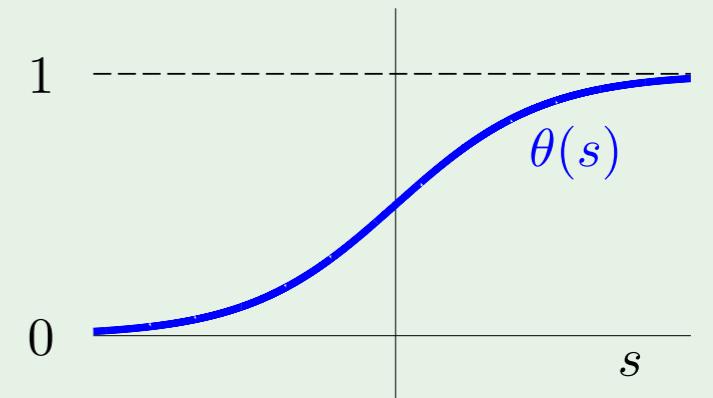
If  $h = f$ , how likely to get  $y$  from  $\mathbf{x}$ ?

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

## Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute  $h(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x})$ , noting  $\theta(-s) = 1 - \theta(s)$



$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^\top \mathbf{x})$$

Likelihood of  $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

## Maximizing the likelihood

Minimize

$$-\frac{1}{N} \ln \left( \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \ln \left( \frac{1}{\theta(y_n \mathbf{w}^\top \mathbf{x}_n)} \right)$$

$$\left[ \theta(s) = \frac{1}{1 + e^{-s}} \right]$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)}_{\text{e}(h(\mathbf{x}_n), y_n)}$$

“cross-entropy” error

## Logistic regression - Outline

- The model
- Error measure
- Learning algorithm

## How to minimize $E_{\text{in}}$

For logistic regression,

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right) \quad \longleftarrow \text{iterative solution}$$

Compare to linear regression:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 \quad \longleftarrow \text{closed-form solution}$$

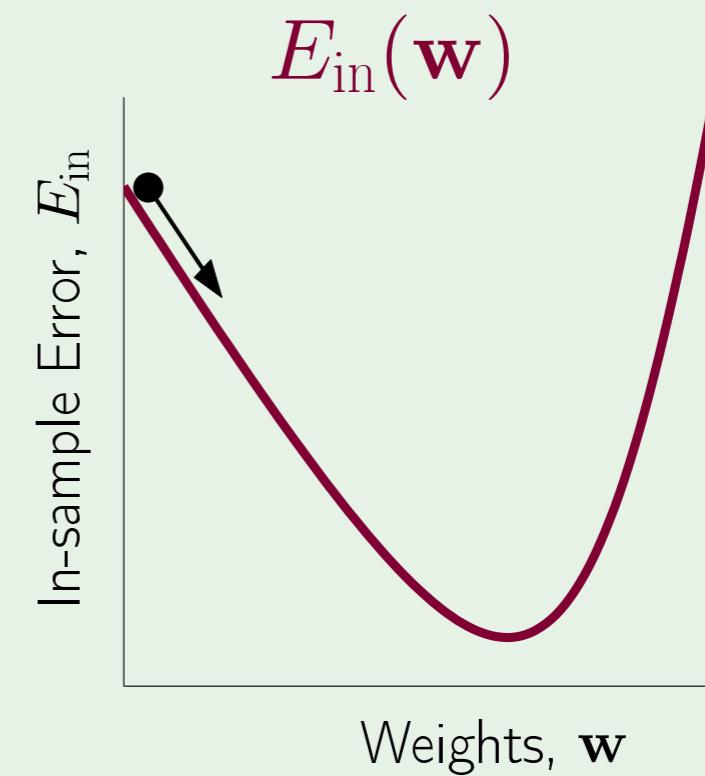
## Iterative method: gradient descent

General method for nonlinear optimization

Start at  $\mathbf{w}(0)$ ; take a step along steepest slope

Fixed step size:  $\mathbf{w}(1) = \mathbf{w}(0) + \eta \hat{\mathbf{v}}$

What is the direction  $\hat{\mathbf{v}}$ ?



## Formula for the direction $\hat{\mathbf{v}}$

$$\Delta E_{\text{in}} = E_{\text{in}}(\mathbf{w}(0) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(0))$$

$$= \eta \nabla E_{\text{in}}(\mathbf{w}(0))^T \hat{\mathbf{v}} + O(\eta^2)$$

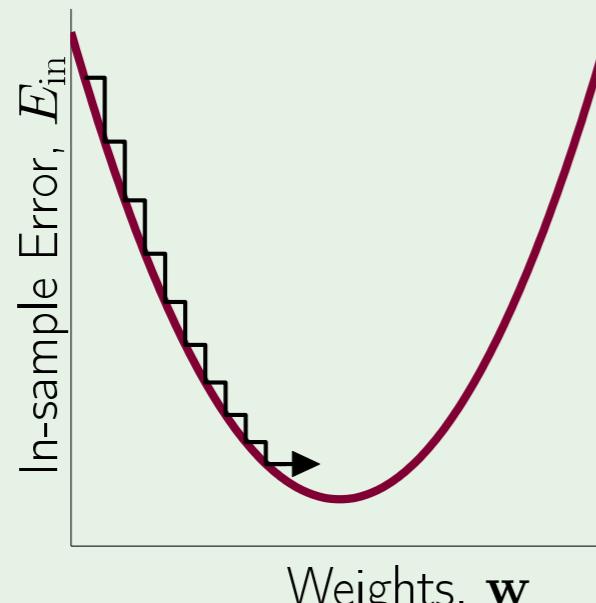
$$\geq -\eta \|\nabla E_{\text{in}}(\mathbf{w}(0))\|$$

Since  $\hat{\mathbf{v}}$  is a unit vector,

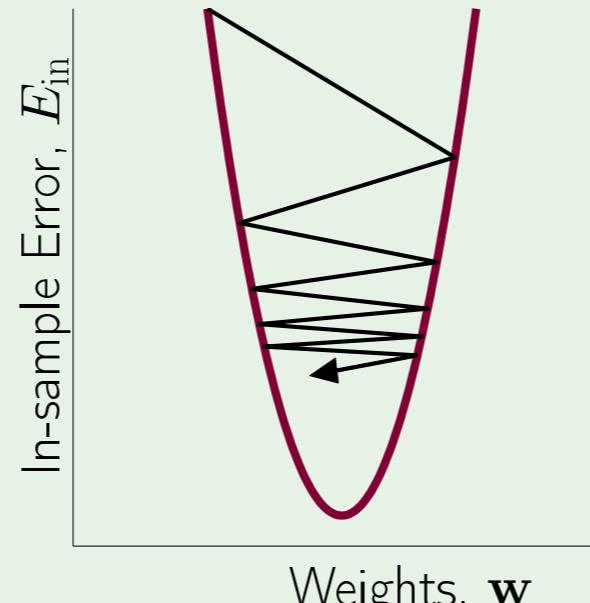
$$\hat{\mathbf{v}} = - \frac{\nabla E_{\text{in}}(\mathbf{w}(0))}{\|\nabla E_{\text{in}}(\mathbf{w}(0))\|}$$

## Fixed-size step?

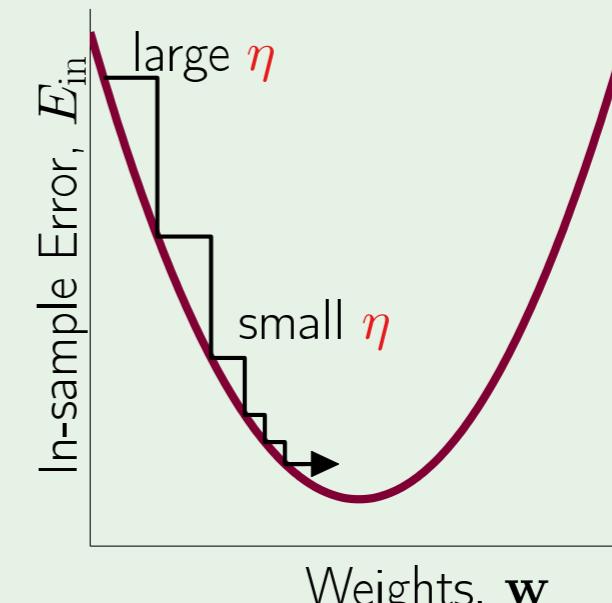
How  $\eta$  affects the algorithm:



$\eta$  too small



$\eta$  too large



variable  $\eta$  – just right

$\eta$  should increase with the slope

## Easy implementation

Instead of

$$\Delta \mathbf{w} = \eta \hat{\mathbf{v}}$$

$$= -\eta \frac{\nabla E_{\text{in}}(\mathbf{w}(0))}{\|\nabla E_{\text{in}}(\mathbf{w}(0))\|}$$

Have

$$\Delta \mathbf{w} = -\eta \nabla E_{\text{in}}(\mathbf{w}(0))$$

Fixed learning rate  $\eta$

## Logistic regression algorithm

1: Initialize the weights at  $t = 0$  to  $\mathbf{w}(0)$

2: **for**  $t = 0, 1, 2, \dots$  **do**

3:   Compute the gradient

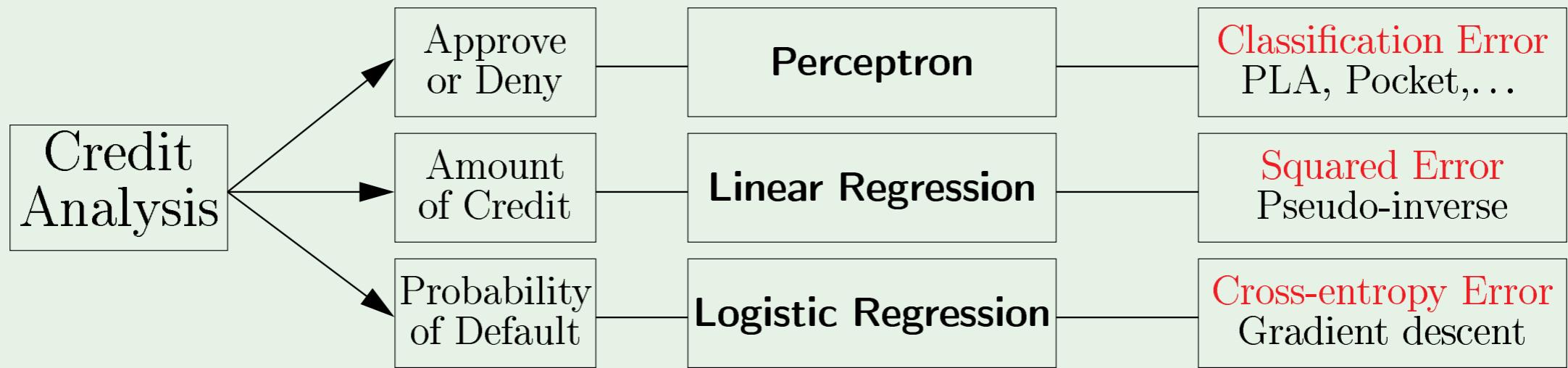
$$\nabla E_{\text{in}} = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^\top(t) \mathbf{x}_n}}$$

4:   Update the weights:  $\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}$

5:   Iterate to the next step until it is time to stop

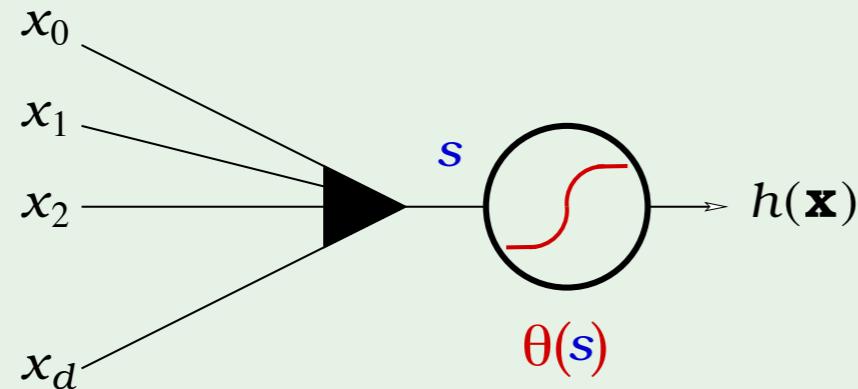
6:   Return the final weights  $\mathbf{w}$

# Summary of Linear Models



## Review of Lecture 9

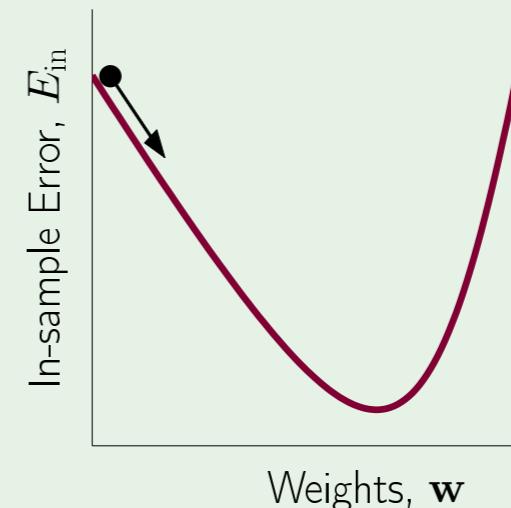
- Logistic regression



- Likelihood measure

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

- Gradient descent



- Initialize  $\mathbf{w}(0)$
- For  $t = 0, 1, 2, \dots$  [to termination]

$$\mathbf{w}(t + 1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}(\mathbf{w}(t))$$

- Return final  $\mathbf{w}$

# Learning From Data

Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 10: Neural Networks

Sponsored by Caltech's Provost Office, E&AS Division, and IST • Thursday, May 3, 2012

# Outline

- Stochastic gradient descent
- Neural network model
- Backpropagation algorithm

## Stochastic gradient descent

GD minimizes:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbf{e}(\mathbf{h}(\mathbf{x}_n), y_n)}_{\ln(1+e^{-y_n \mathbf{w}^\top \mathbf{x}_n})} \quad \leftarrow \text{in logistic regression}$$

by iterative steps along  $-\nabla E_{\text{in}}$ :

$$\Delta \mathbf{w} = -\eta \nabla E_{\text{in}}(\mathbf{w})$$

$\nabla E_{\text{in}}$  is based on all examples  $(\mathbf{x}_n, y_n)$

“batch” GD

## The stochastic aspect

Pick one  $(\mathbf{x}_n, y_n)$  at a time. Apply GD to  $\mathbf{e}(h(\mathbf{x}_n), y_n)$

“Average” direction:

$$\begin{aligned}\mathbb{E}_{\mathbf{n}}[-\nabla \mathbf{e}(h(\mathbf{x}_n), y_n)] &= \frac{1}{N} \sum_{n=1}^N -\nabla \mathbf{e}(h(\mathbf{x}_n), y_n) \\ &= -\nabla E_{\text{in}}\end{aligned}$$

randomized version of GD

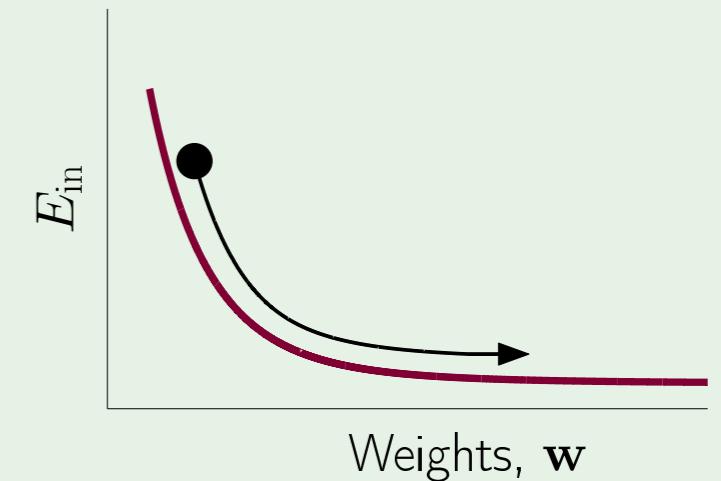
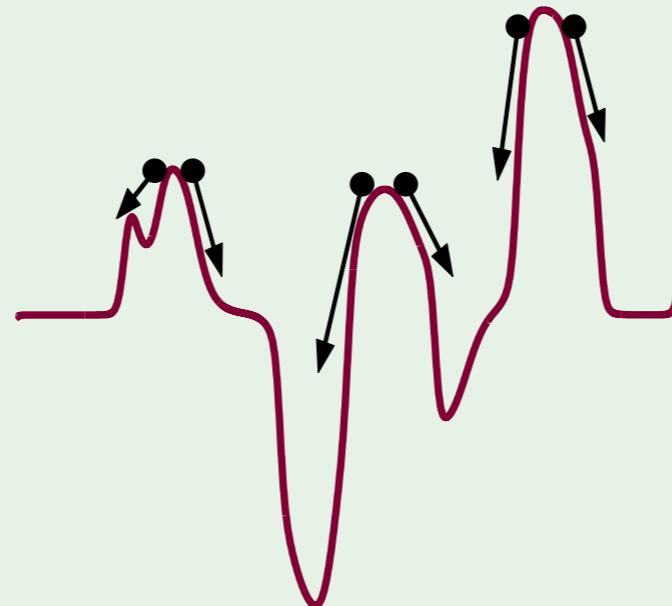
**stochastic** gradient descent (SGD)

## Benefits of SGD

1. cheaper computation
2. randomization
3. simple

**Rule of thumb:**

$\eta = 0.1$  works

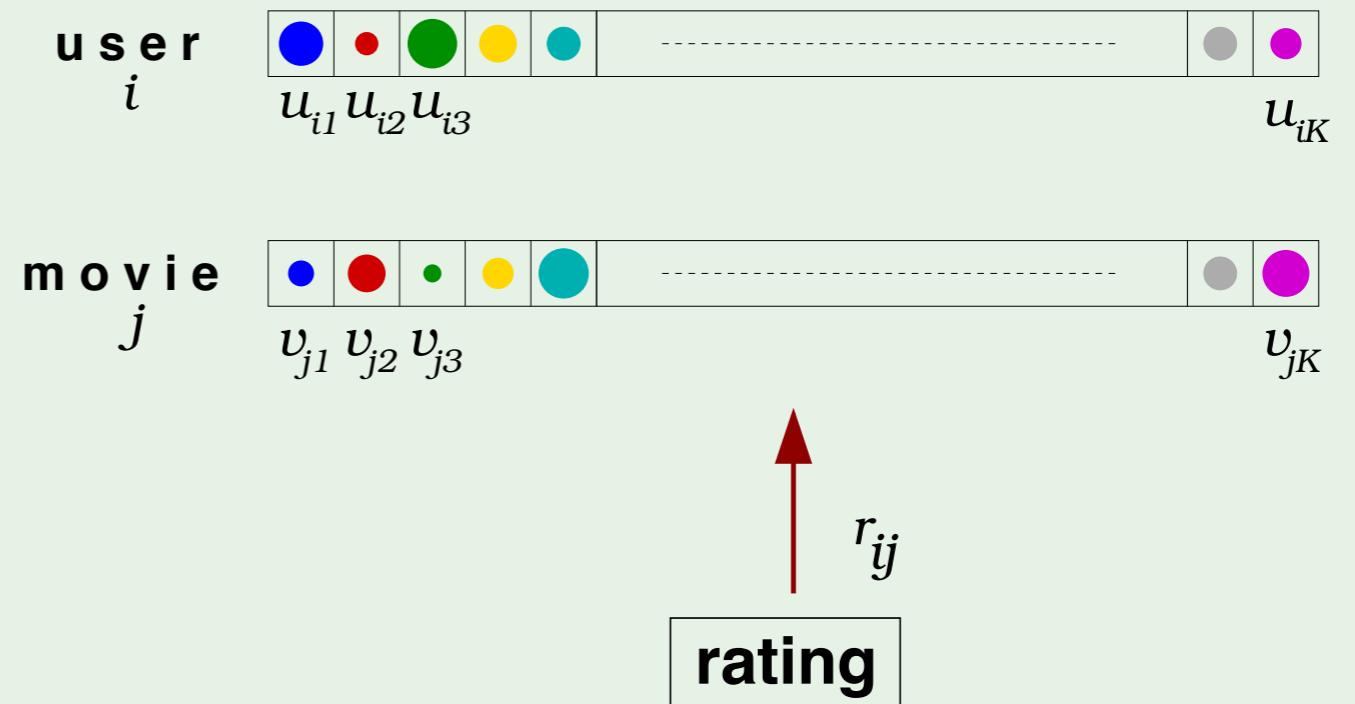


**randomization helps**

## SGD in action

Remember movie ratings?

$$\mathbf{e}_{ij} = \left( r_{ij} - \sum_{k=1}^K u_{ik} v_{jk} \right)^2$$

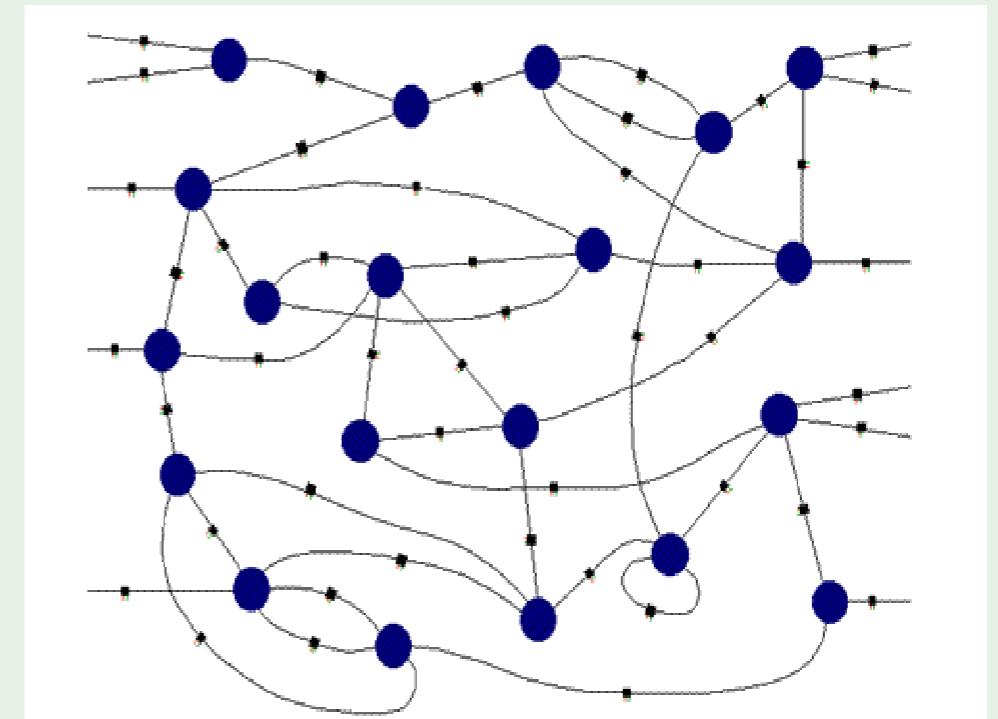
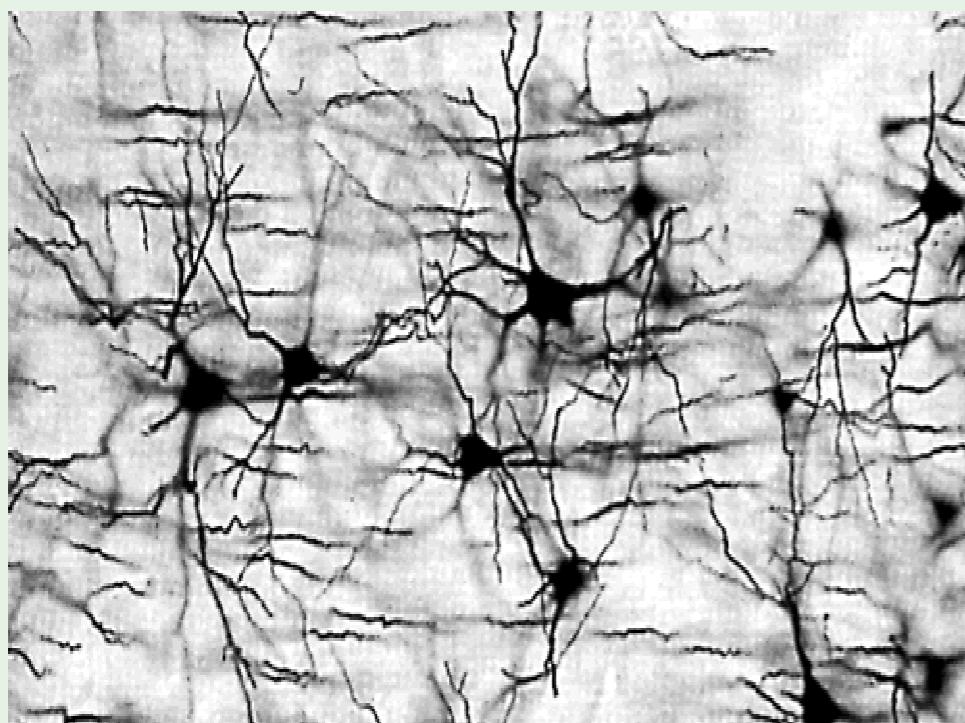


# Outline

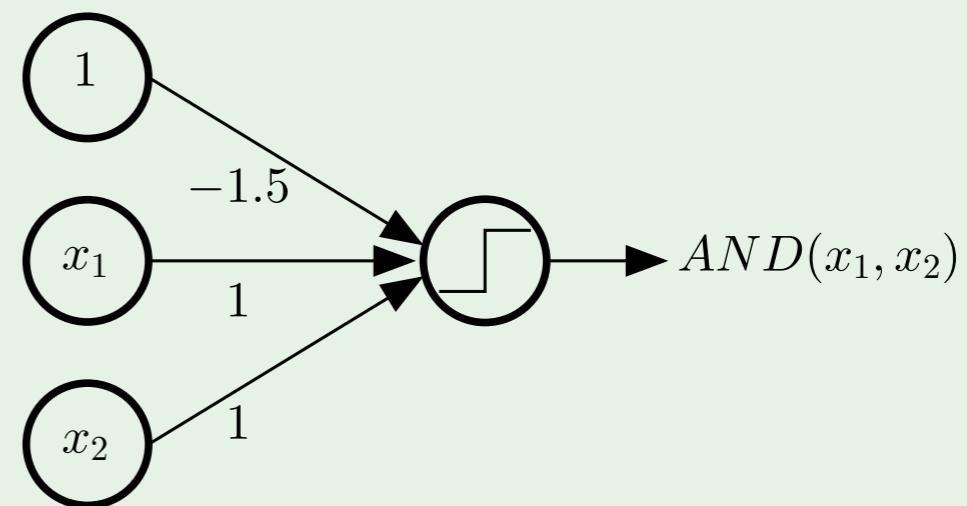
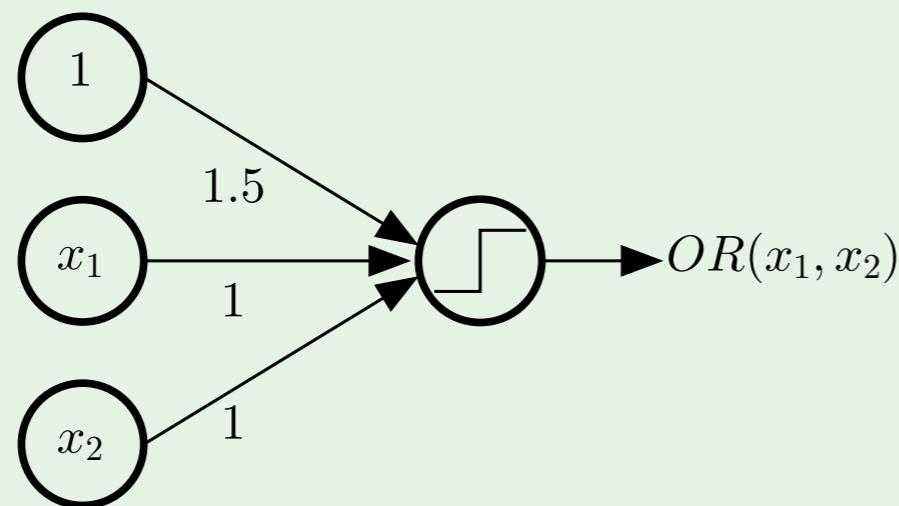
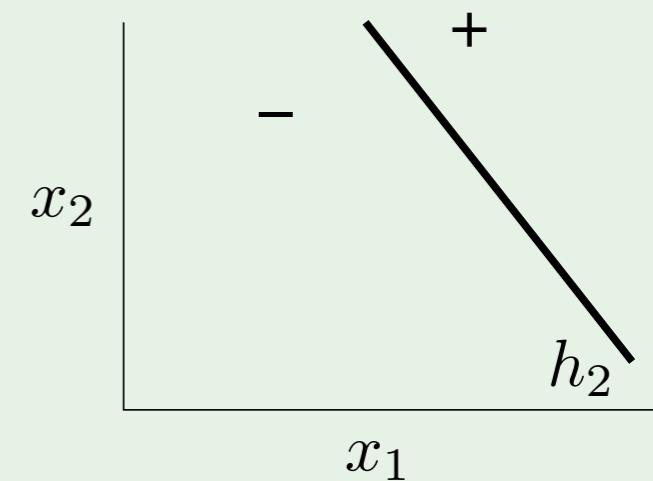
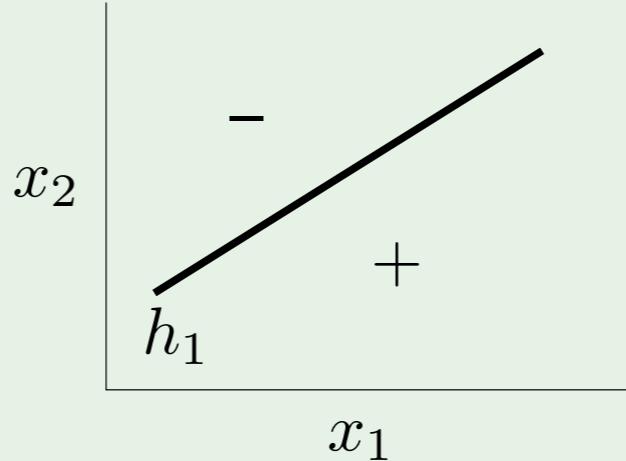
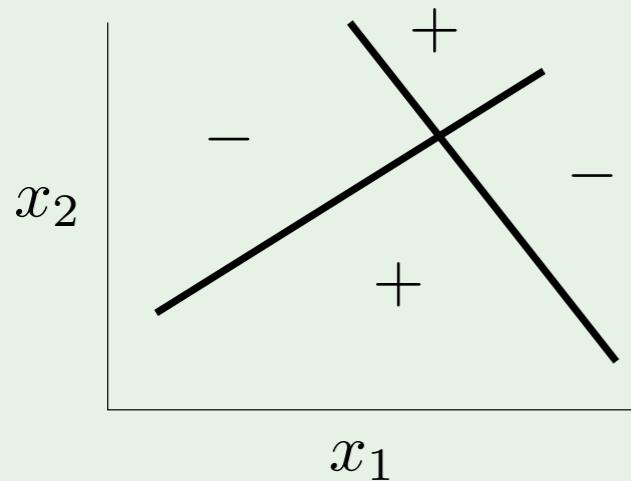
- Stochastic gradient descent
- Neural network model
- Backpropagation algorithm

# Biological inspiration

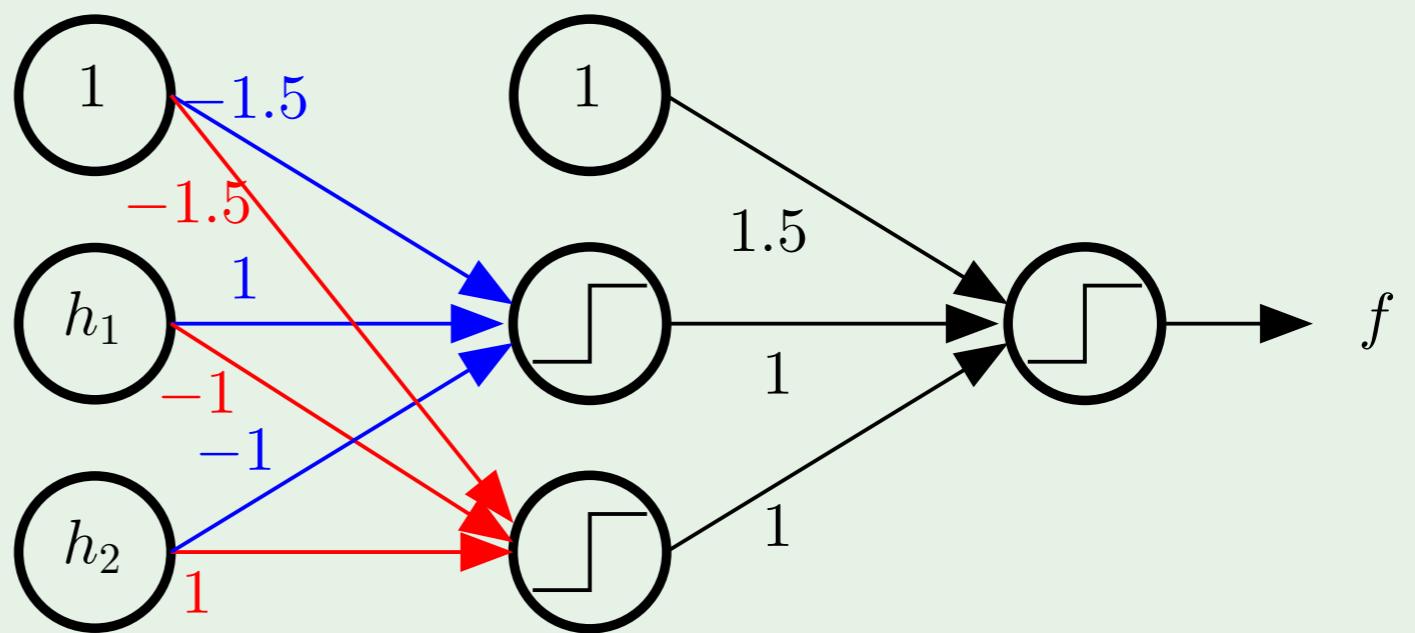
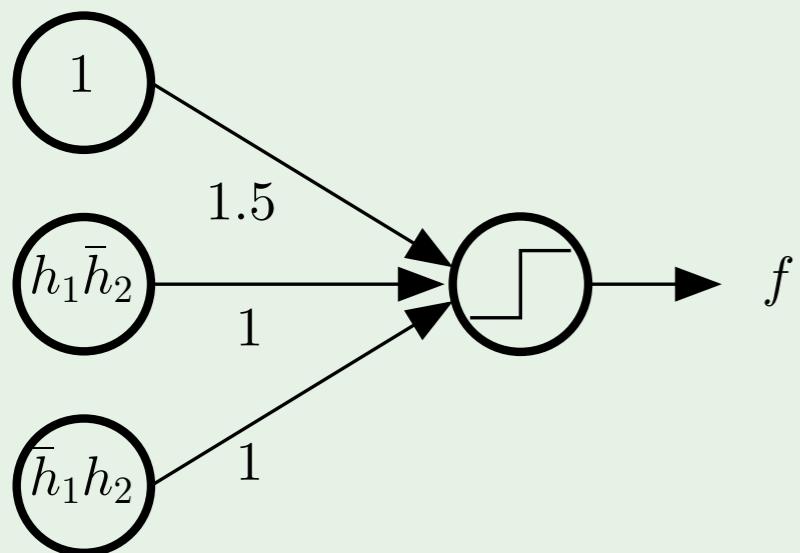
biological function → biological structure



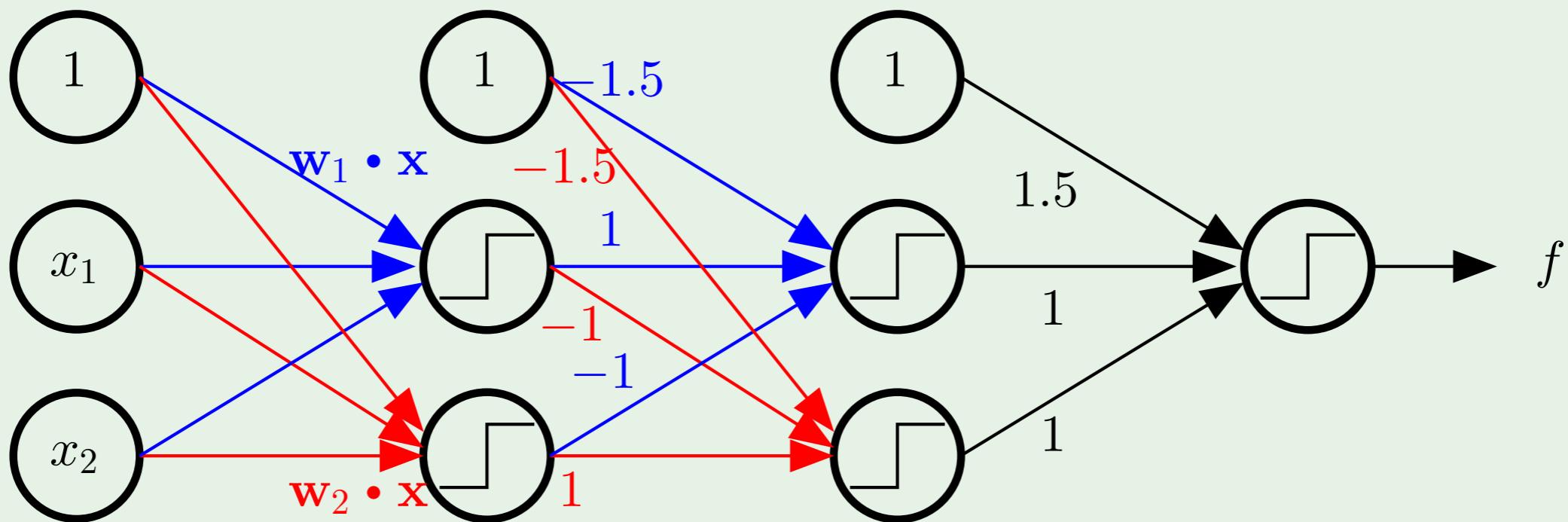
## Combining perceptrons



## Creating layers

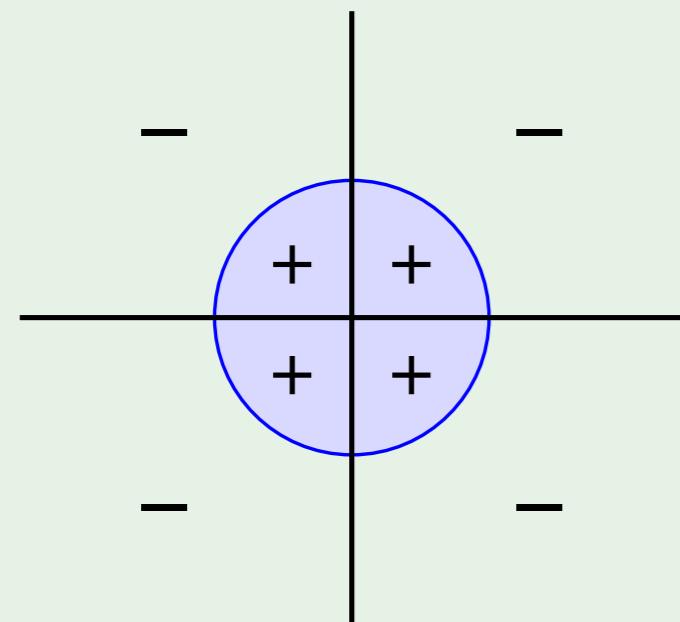


## The multilayer perceptron

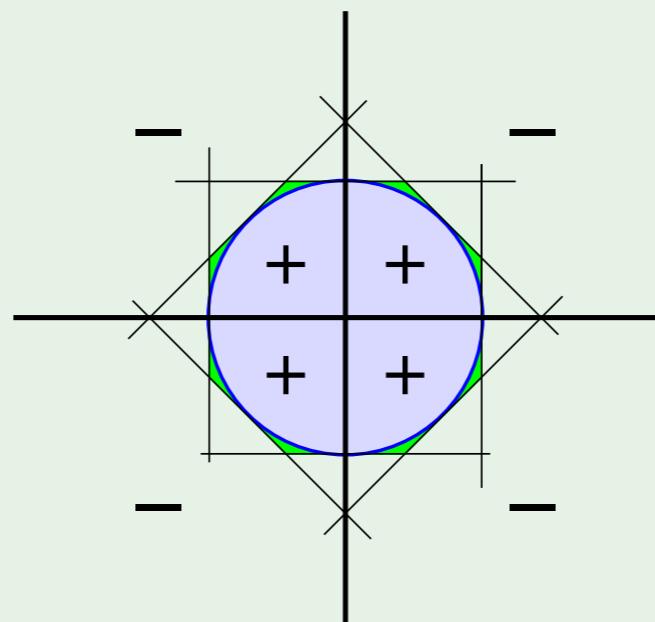


3 layers     “feedforward”

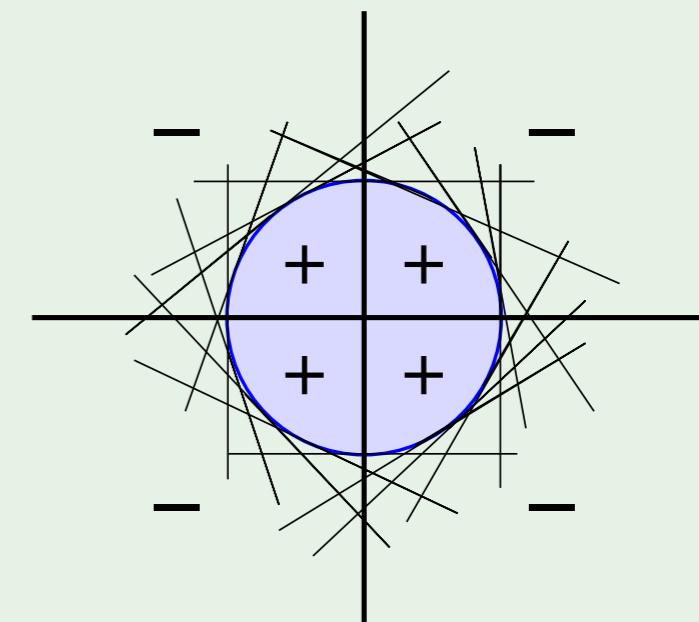
## A powerful model



Target



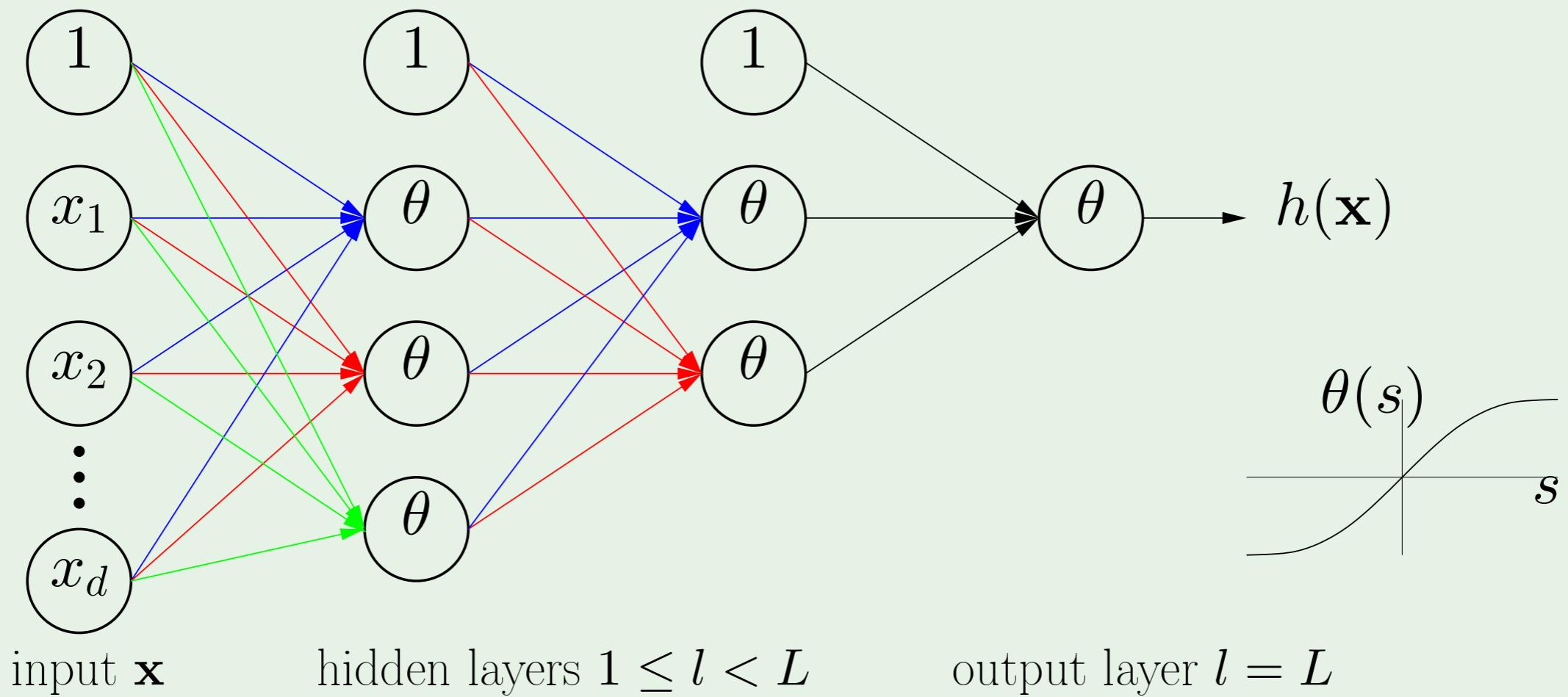
8 perceptrons



16 perceptrons

**2 red flags** for generalization and optimization

# The neural network

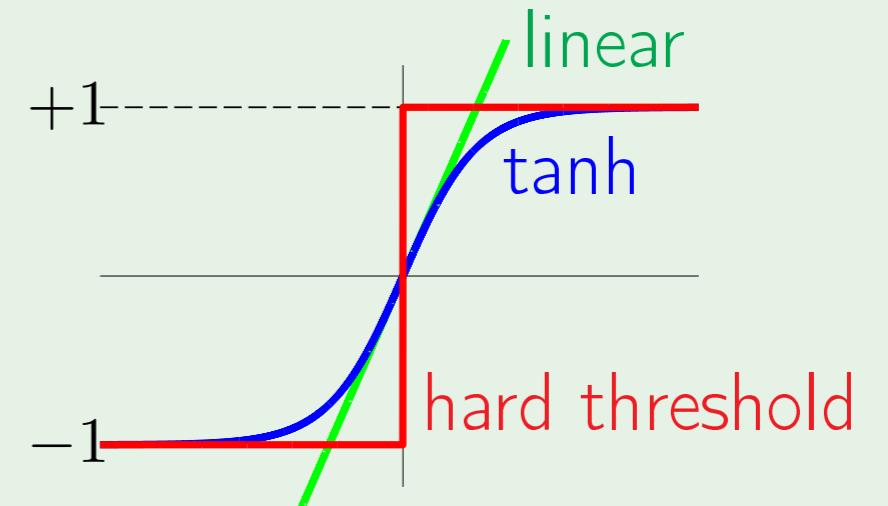


## How the network operates

$$w_{ij}^{(l)} \quad \begin{cases} 1 \leq l \leq L & \text{layers} \\ 0 \leq i \leq d^{(l-1)} & \text{inputs} \\ 1 \leq j \leq d^{(l)} & \text{outputs} \end{cases}$$

$$x_j^{(l)} = \theta(s_j^{(l)}) = \theta \left( \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)} \right)$$

Apply  $\mathbf{x}$  to  $x_1^{(0)} \cdots x_{d^{(0)}}^{(0)} \rightarrow \rightarrow x_1^{(L)} = h(\mathbf{x})$



$$\theta(s) = \tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

# Outline

- Stochastic gradient descent
- Neural network model
- Backpropagation algorithm

# Applying SGD

All the weights  $\mathbf{w} = \{\mathbf{w}_{ij}^{(l)}\}$  determine  $h(\mathbf{x})$

Error on example  $(\mathbf{x}_n, y_n)$  is

$$\mathbf{e}(h(\mathbf{x}_n), y_n) = \mathbf{e}(\mathbf{w})$$

To implement SGD, we need the gradient

$$\nabla \mathbf{e}(\mathbf{w}): \frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}} \quad \text{for all } i, j, l$$

## Computing $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$

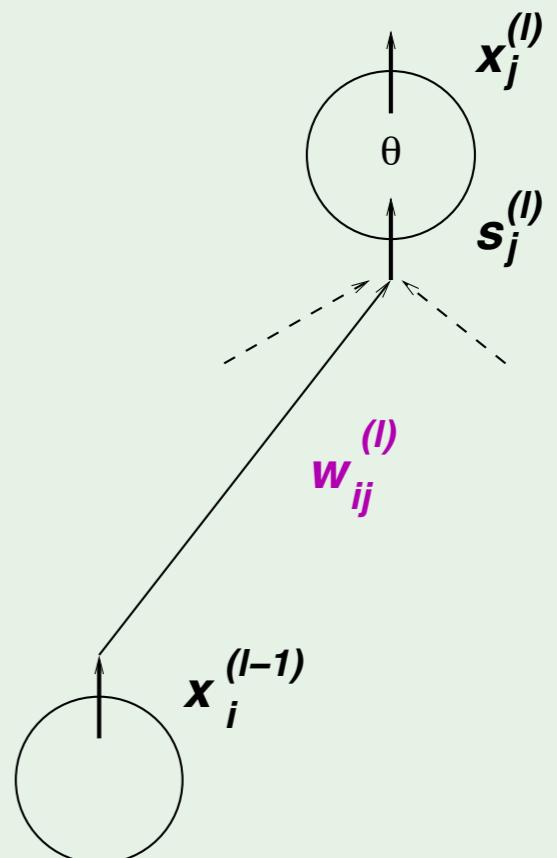
We can evaluate  $\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}}$  one by one: analytically or numerically

A trick for efficient computation:

$$\frac{\partial \mathbf{e}(\mathbf{w})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}$$

$$\text{We have } \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$$

$$\text{We only need: } \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)}$$



## $\delta$ for the final layer

$$\delta_j^{(l)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer  $l = L$  and  $j = 1$ :

$$\delta_1^{(L)} = \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_1^{(L)}}$$

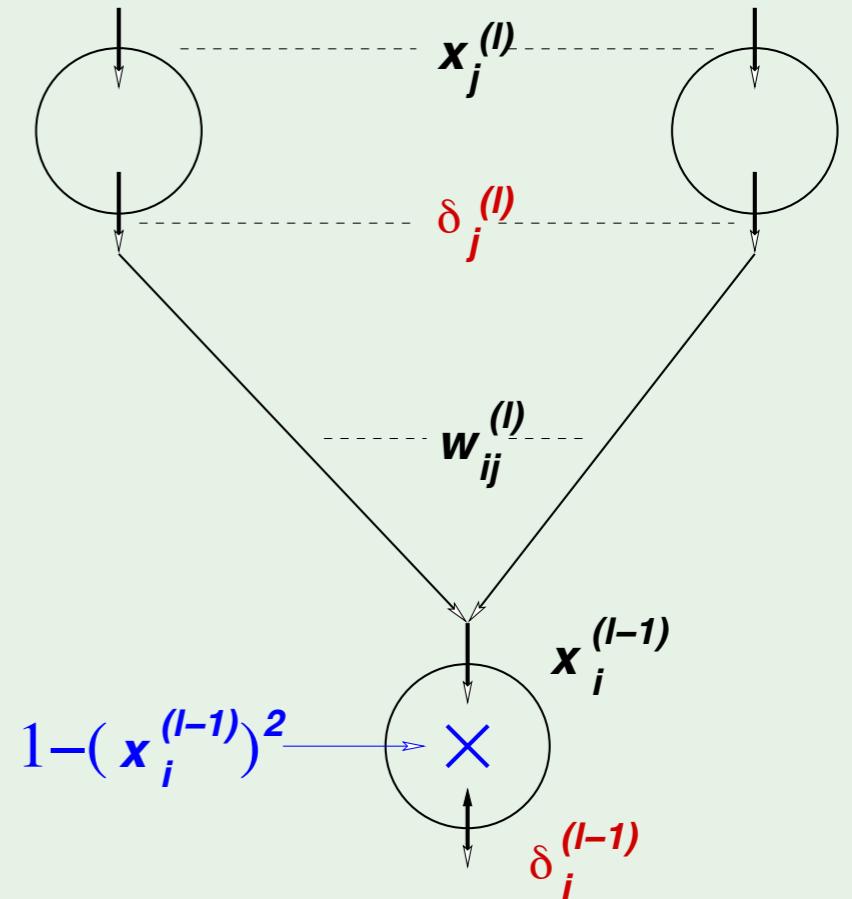
$$\mathbf{e}(\mathbf{w}) = (x_1^{(L)} - y_n)^2$$

$$x_1^{(L)} = \theta(s_1^{(L)})$$

$$\theta'(s) = 1 - \theta^2(s) \quad \text{for the tanh}$$

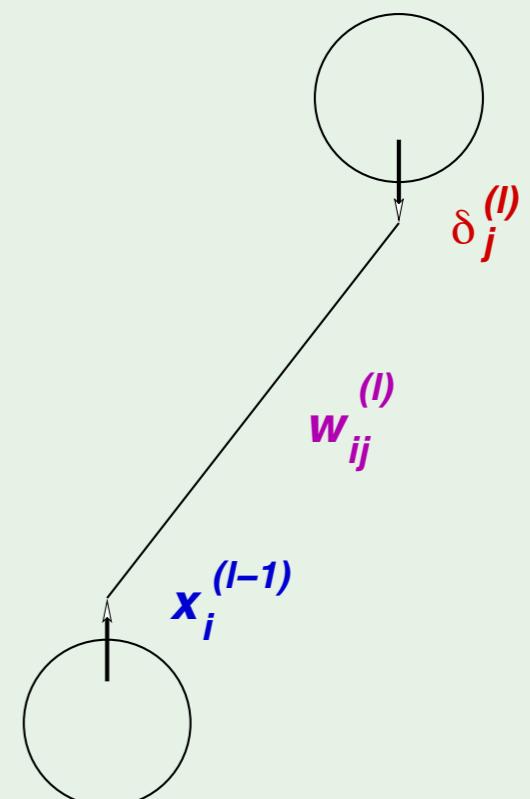
## Back propagation of $\delta$

$$\begin{aligned}
 \delta_i^{(l-1)} &= \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d(l)} \frac{\partial \mathbf{e}(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d(l)} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(s_i^{(l-1)}) \\
 \delta_i^{(l-1)} &= (1 - (x_i^{(l-1)})^2) \sum_{j=1}^{d(l)} w_{ij}^{(l)} \delta_j^{(l)}
 \end{aligned}$$



# Backpropagation algorithm

- 1: Initialize all weights  $w_{ij}^{(l)}$  **at random**
- 2: **for**  $t = 0, 1, 2, \dots$  **do**
- 3:   Pick  $n \in \{1, 2, \dots, N\}$
- 4:   *Forward:* Compute all  $x_j^{(l)}$
- 5:   *Backward:* Compute all  $\delta_j^{(l)}$
- 6:   Update the weights:  $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$
- 7:   Iterate to the next step until it is time to stop
- 8: Return the final weights  $w_{ij}^{(l)}$



## Final remark: hidden layers

learned nonlinear transform

interpretation?

