

<https://introml.mit.edu/>

# 6.390 Intro to Machine Learning

## Lecture 5: Features, Neural Networks I

Shen Shen

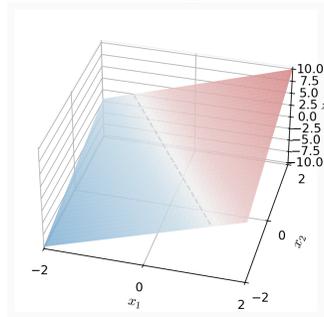
Mar 2, 2026

3pm, Room 10-250

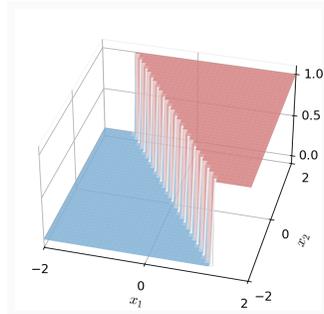
[Slides and Lecture Recording](#)

# Recall:

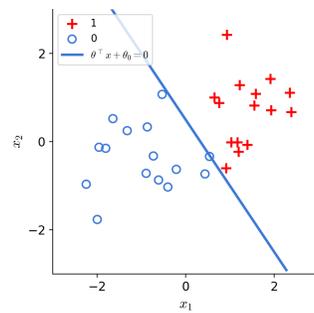
## step



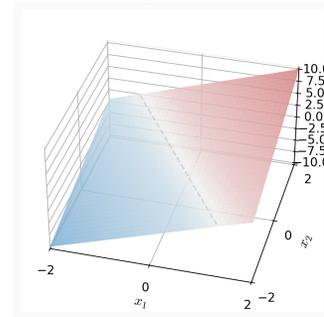
$$z = \theta^T x + \theta_0$$



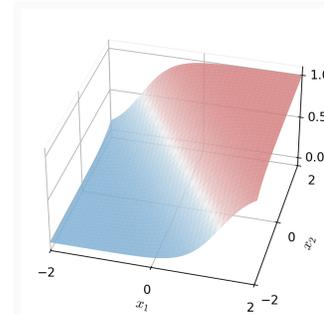
$$g = \text{step}(z)$$



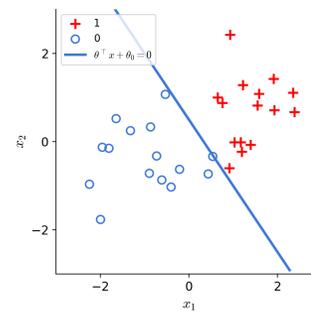
## logistic



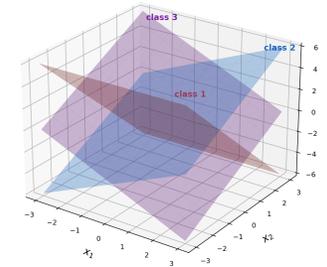
$$z = \theta^T x + \theta_0$$



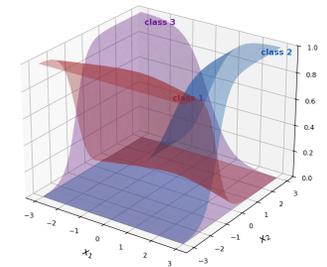
$$g = \sigma(z)$$



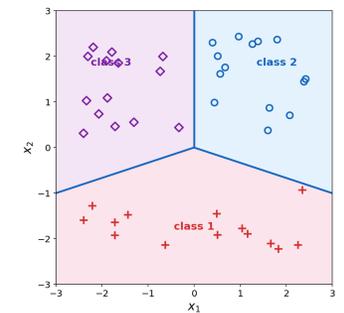
## multi-class



$$z_k = \theta_k^T x + \theta_{0,k}$$

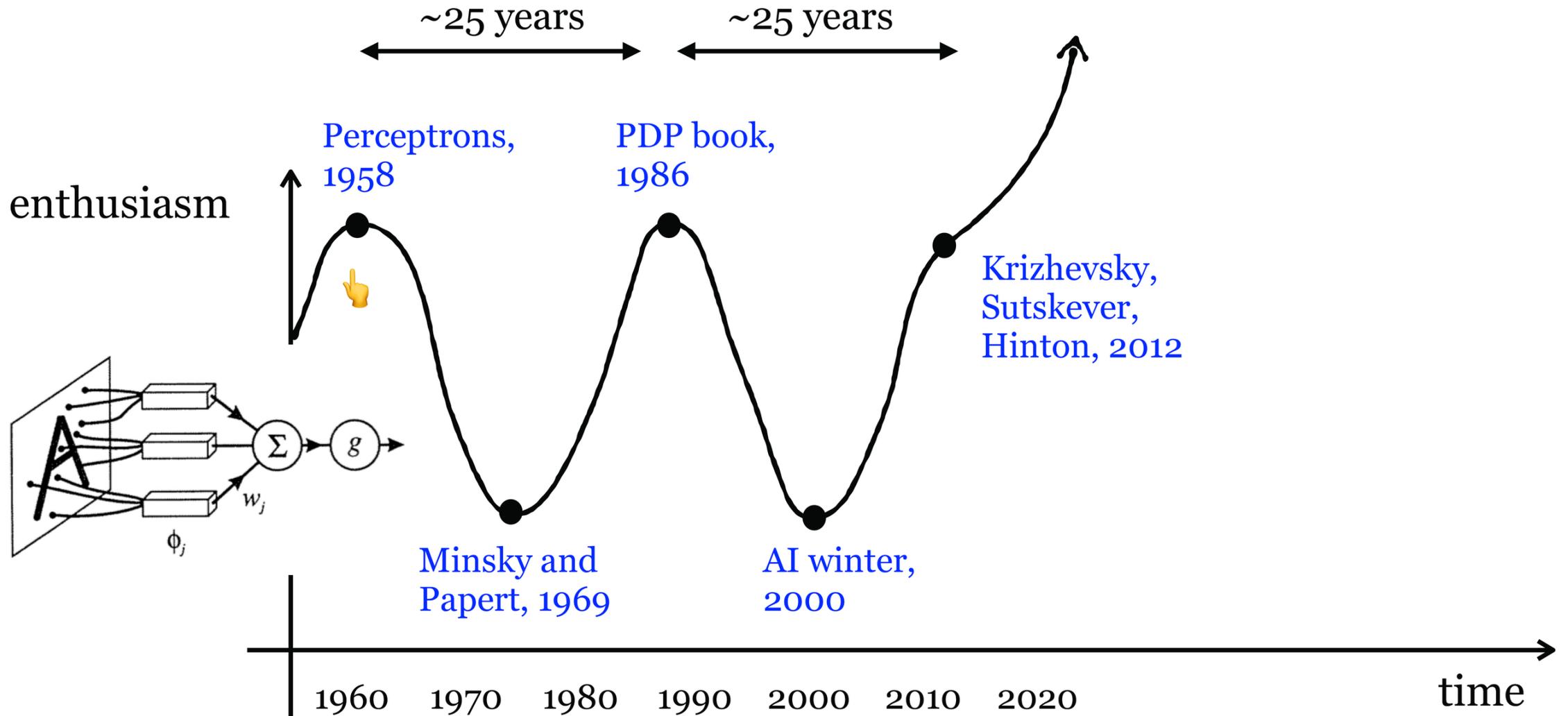


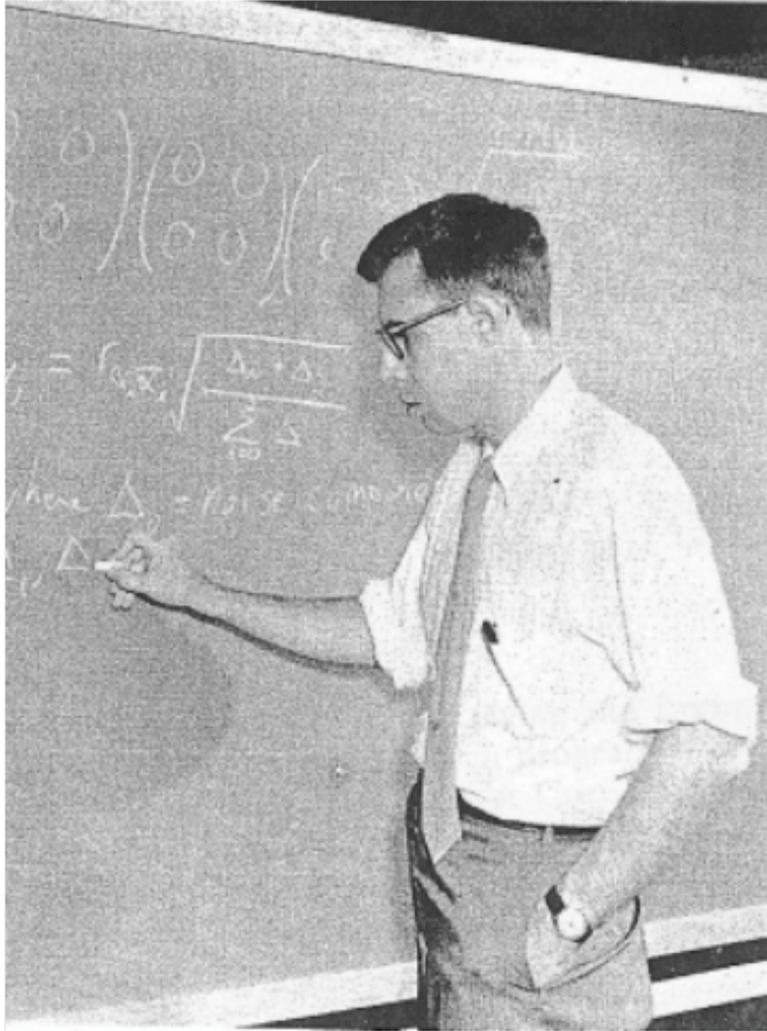
$$g = \text{softmax}(z_1, \dots, z_K)$$



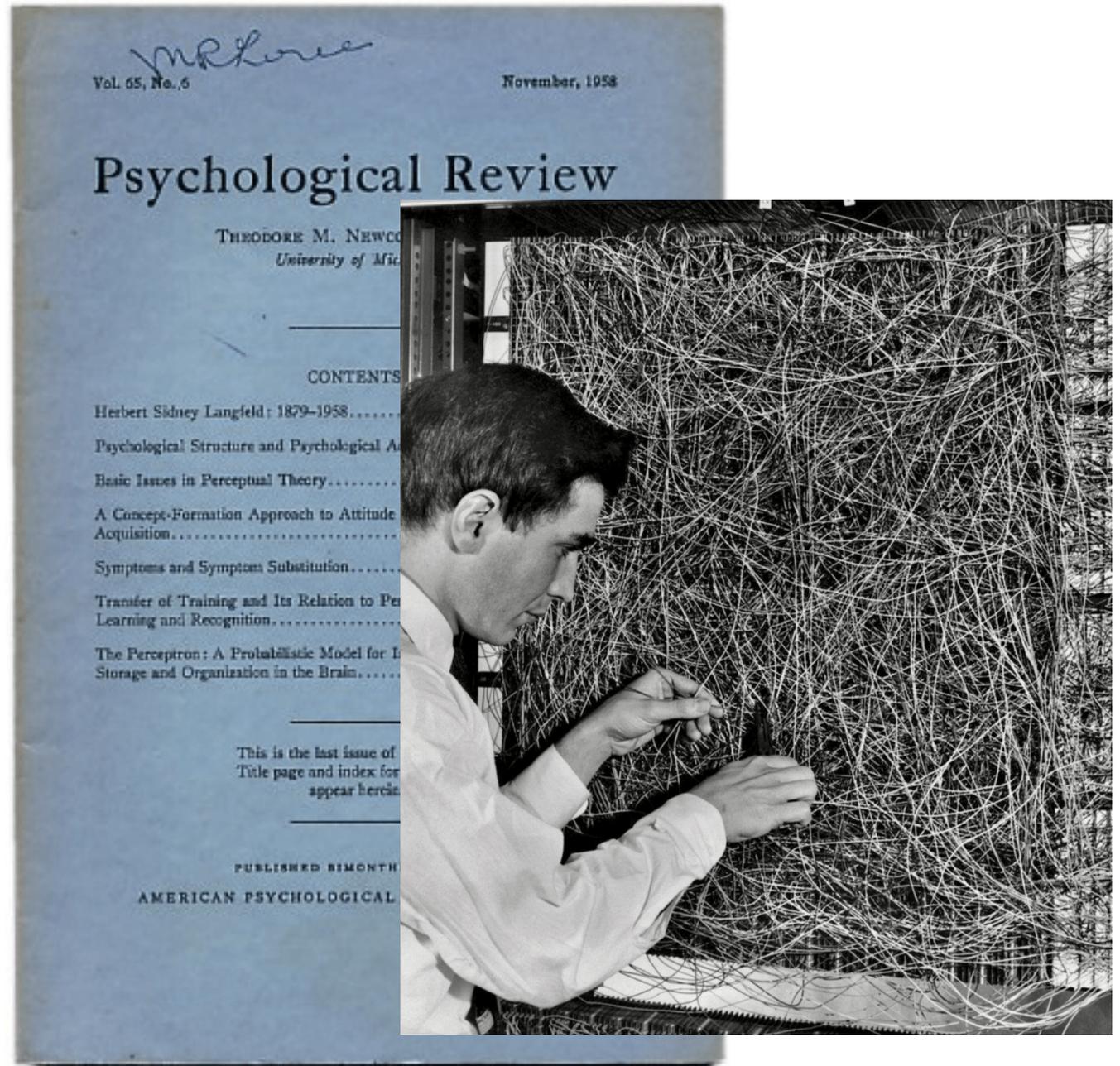
decision boundary is *linear* in feature  $x$

Linear classification played a pivotal role in kicking off the first wave of AI enthusiasm.





[http://www.ecse.rpi.edu/homepages/nagy/PDF\\_chrono/2011\\_Nagy\\_Pace\\_FR.pdf](http://www.ecse.rpi.edu/homepages/nagy/PDF_chrono/2011_Nagy_Pace_FR.pdf). Photo by George Nagy



<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.335.3398&rep=rep1&type=pdf>

# NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo  
of Computer Designed to  
Read and Grow Wiser

WASHINGTON, July 7 (UPI)  
—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the "Perceptron"  
costs \$2,000,000. "704" com-

1958 New York  
Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

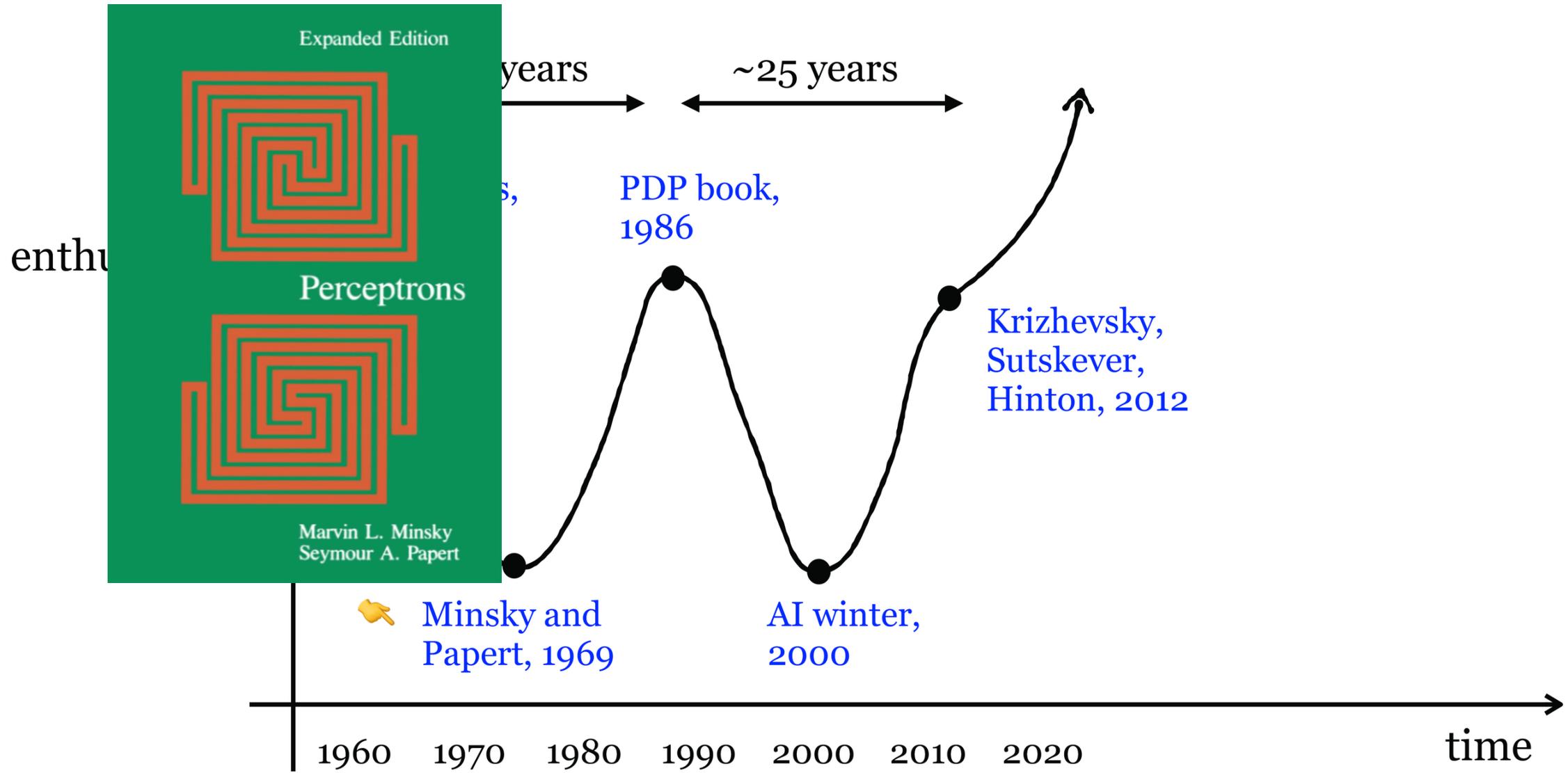
## Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

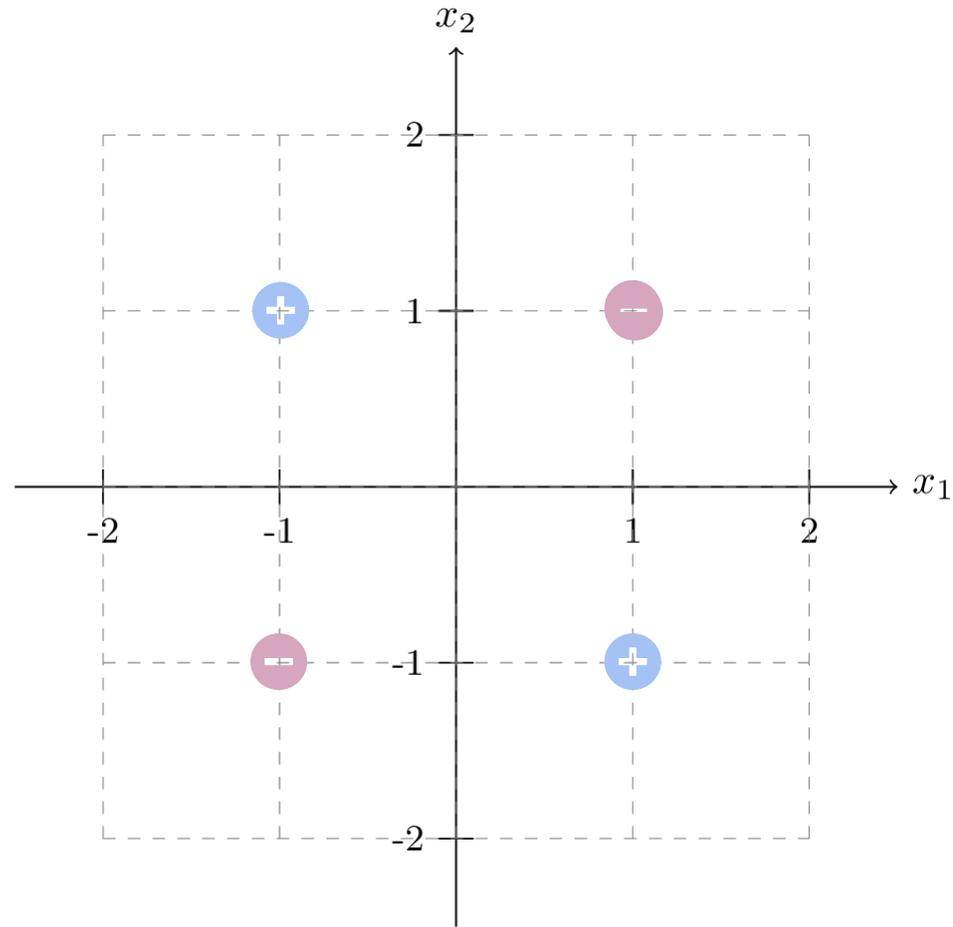
Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

nis-  
row  
he  
in  
ta  
rch  
nell  
w  
3uf-  
se  
be  
ani-  
p  
ron  
ring  
A  
eiv-  
ing  
fi  
any  
c  
to  
ma-  
rdi-  
only  
nch  
able  
w  
out  
ins-  
m  
to  
her  
ir  
rin-  
it  
to  
h  
pro-  
ibly  
con-  
re



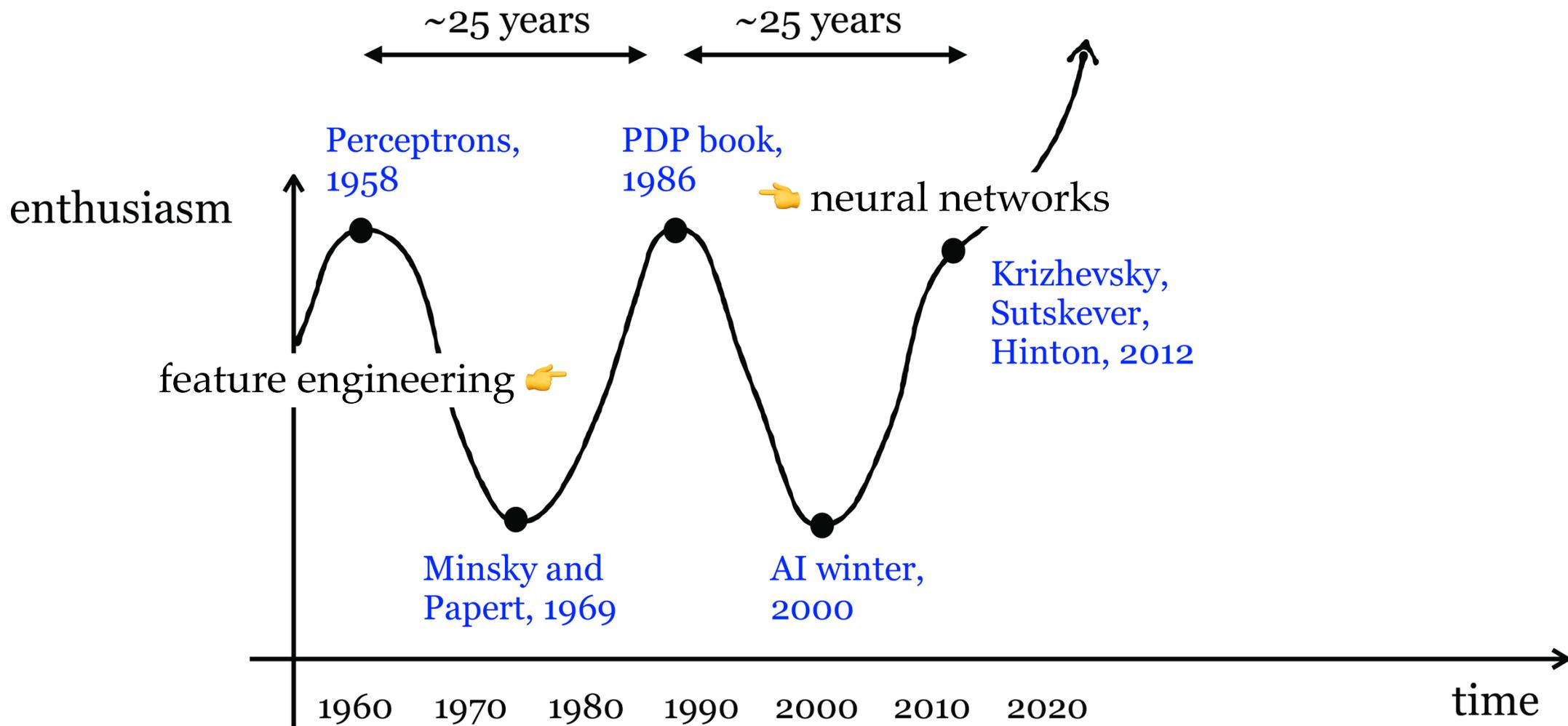
## XOR dataset



Not **linearly** separable.

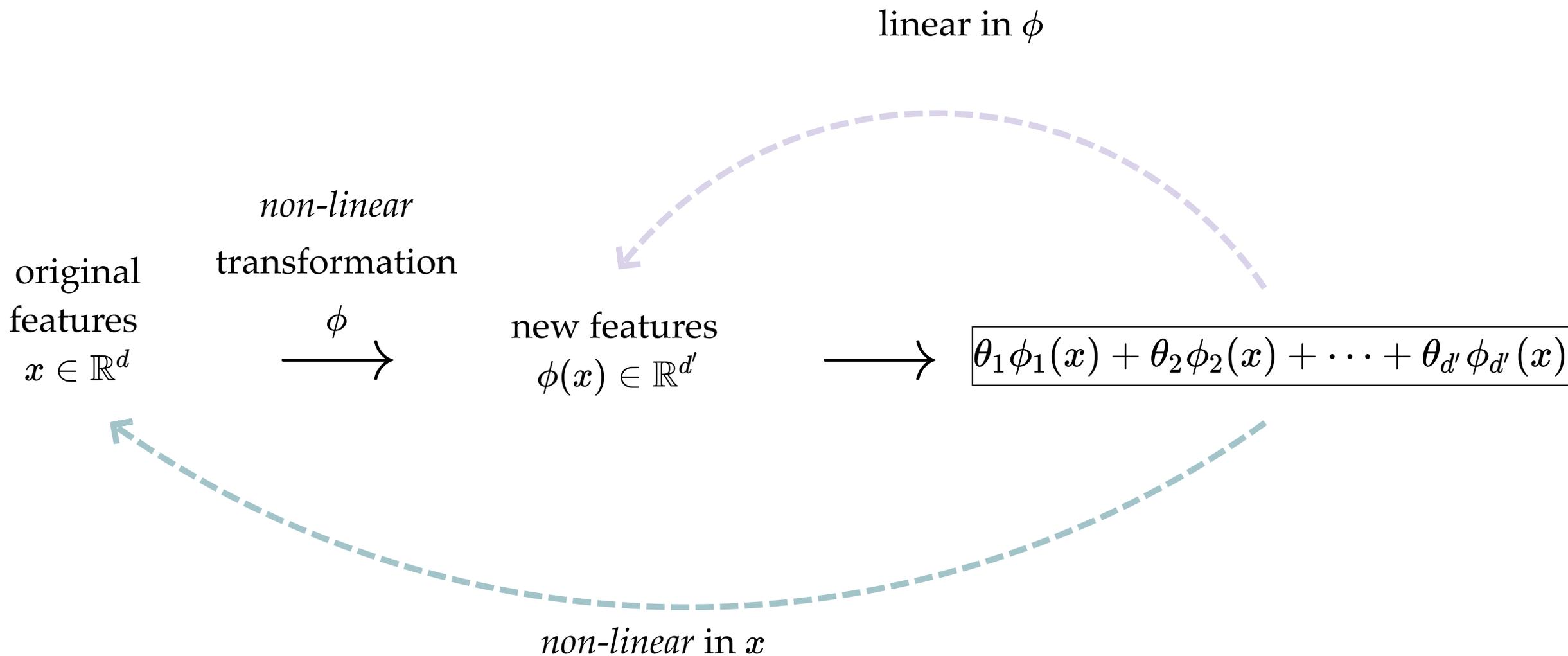
~~Linear tools cannot solve interesting tasks.~~

Linear tools cannot, *by themselves*, solve interesting tasks.

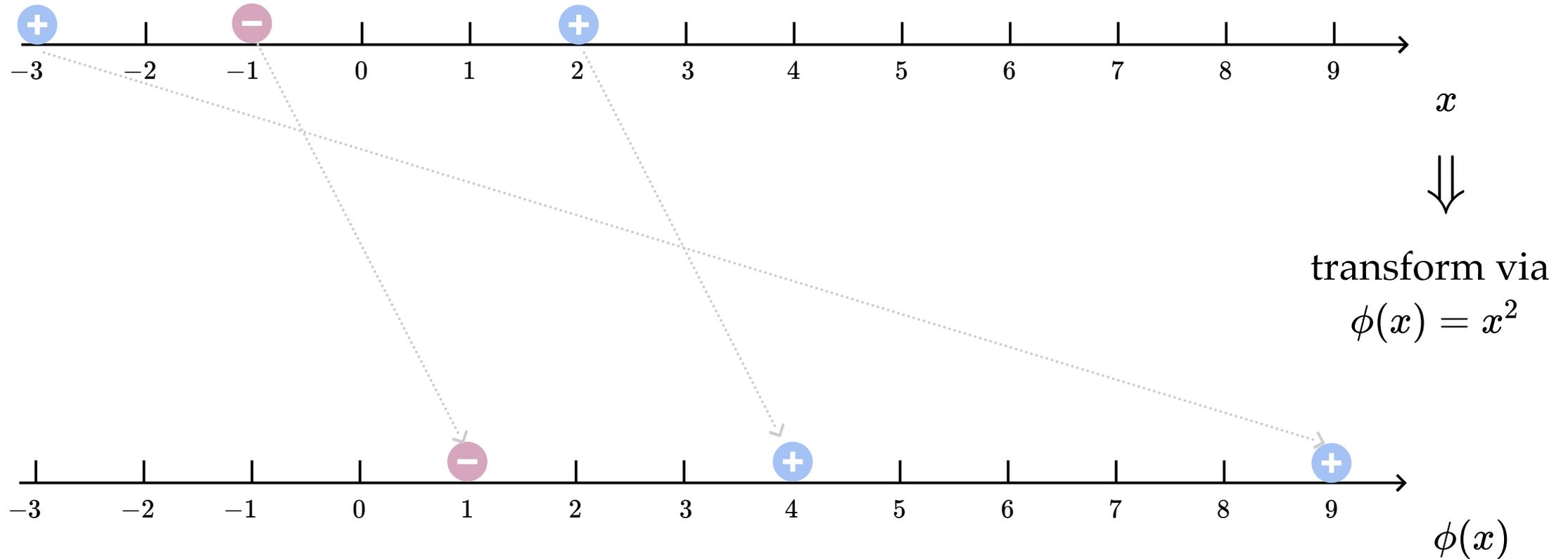


# Outline

- Systematic feature transformations
  - Engineered features
  - Polynomial features
  - Expressive power
- Neural networks

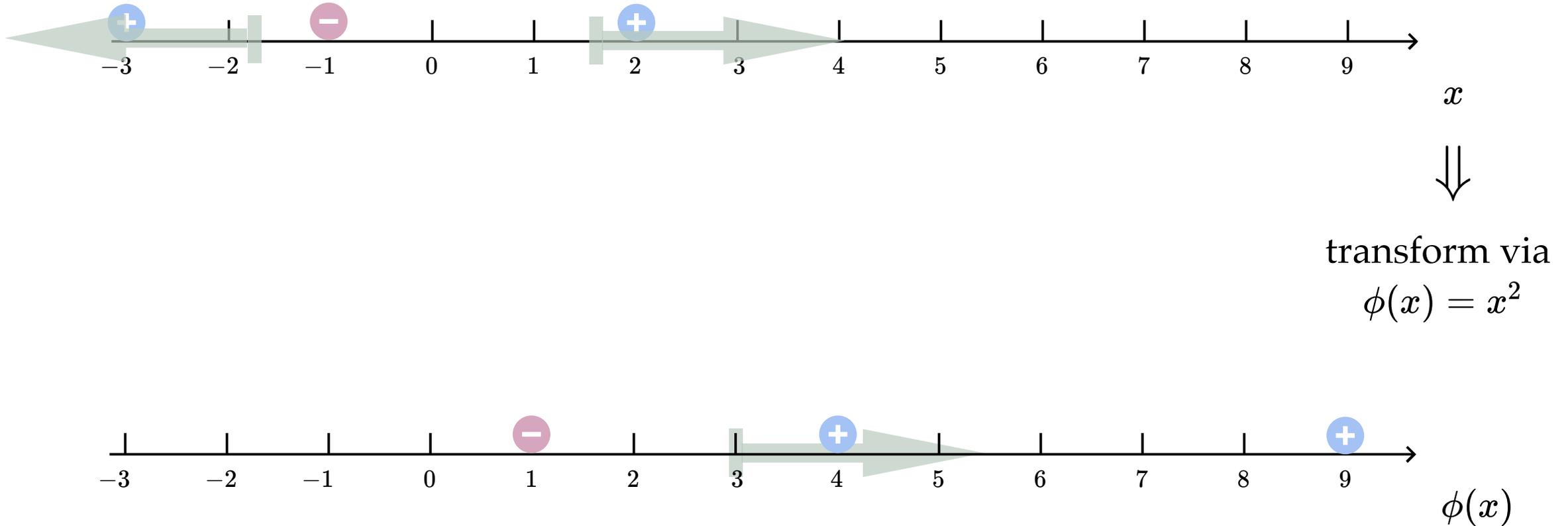


Not linearly separable in  $x$  space



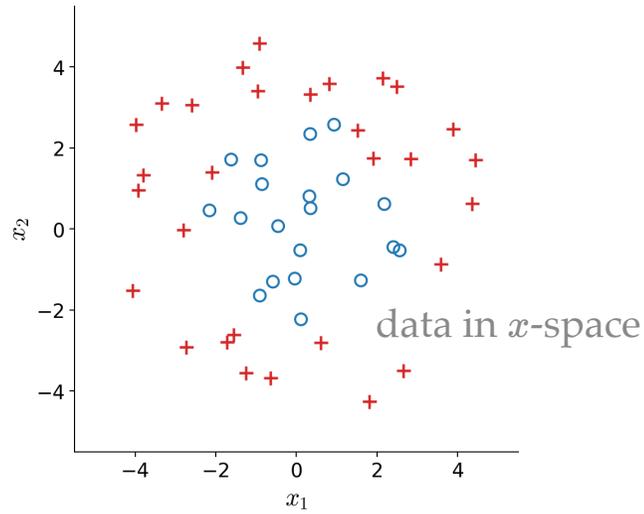
Linearly separable in  $\phi(x) = x^2$  space

Non-linearly separated in  $x$  space, e.g. predict positive if  $x^2 \geq 3$

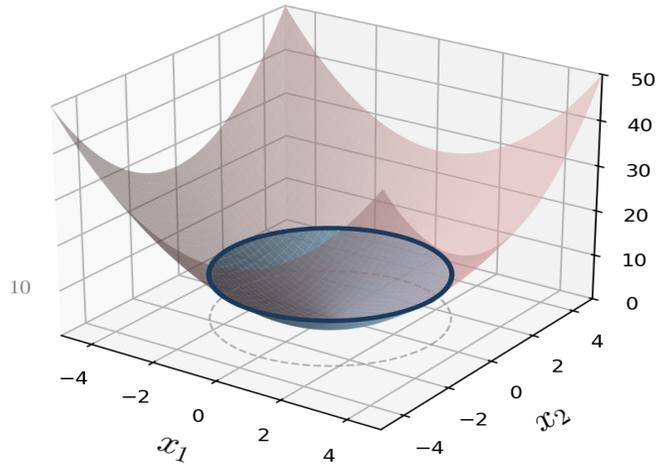


Linearly separable in  $\phi(x) = x^2$  space, e.g. predict positive if  $\phi \geq 3$

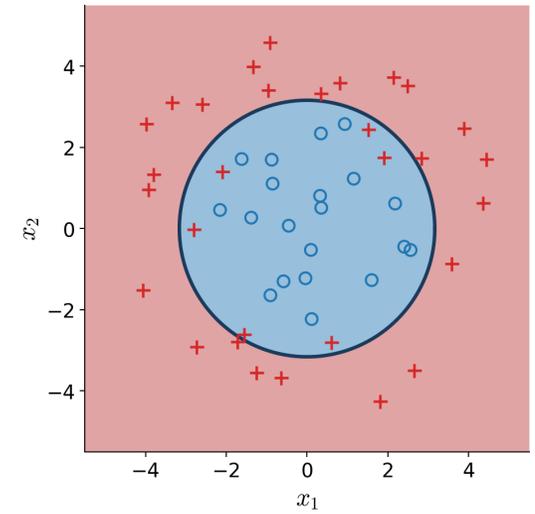
# non-linear classification



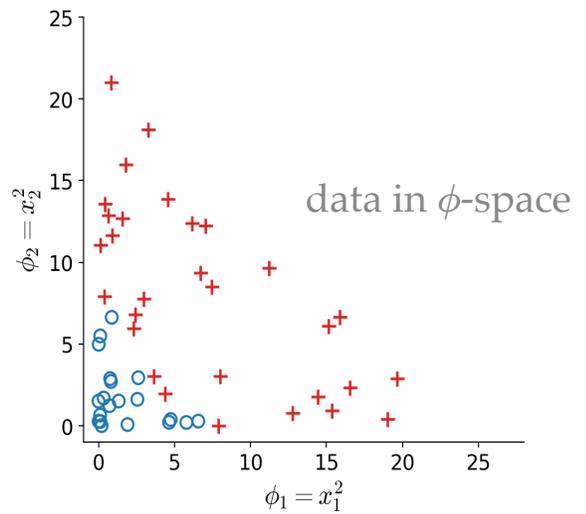
$$z = x_1^2 + x_2^2, \text{ threshold at } z = 10$$



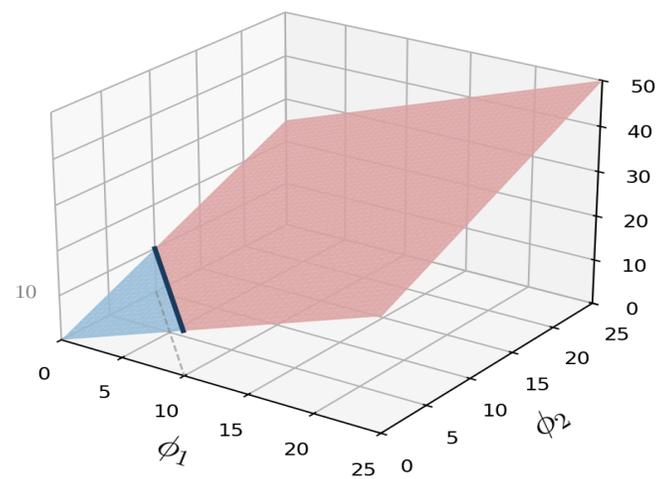
$$x_1^2 + x_2^2 = 10: \text{ non-linear in } x$$



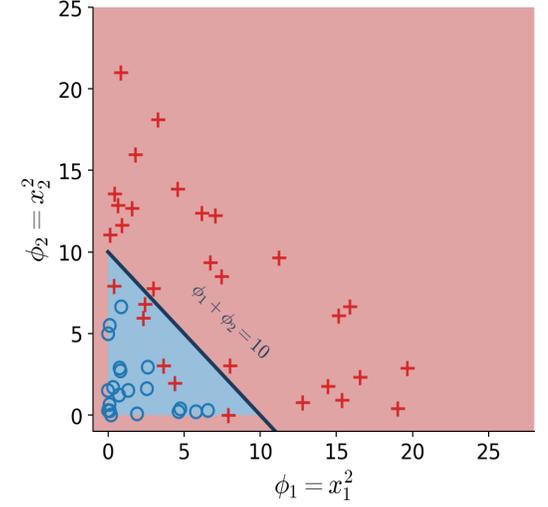
transform via  $\phi(x) = (x_1^2, x_2^2)$



$$z = \phi_1 + \phi_2, \text{ threshold at } z = 10$$

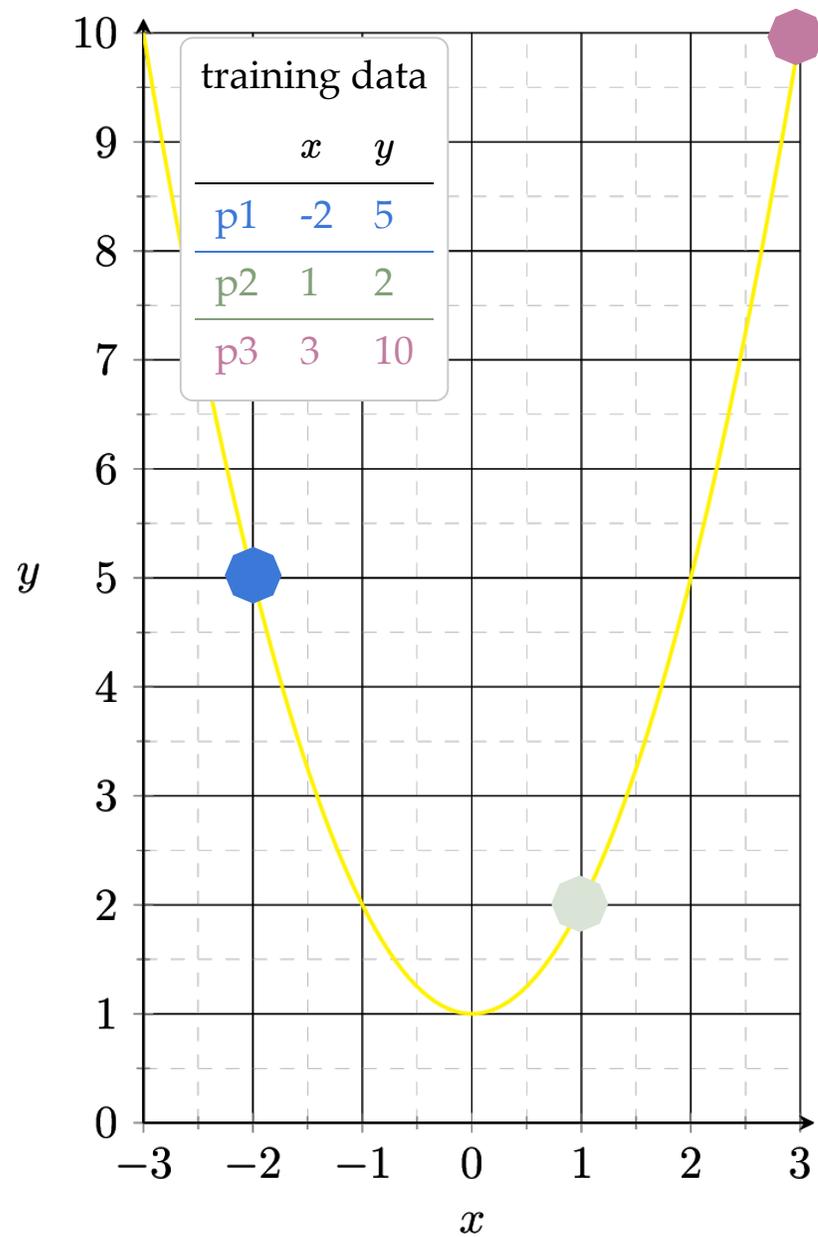


$$\phi_1 + \phi_2 = 10: \text{ linear in } \phi$$



decision boundary is *linear* in  $\phi$ , *nonlinear* in  $x$

# non-linear regression

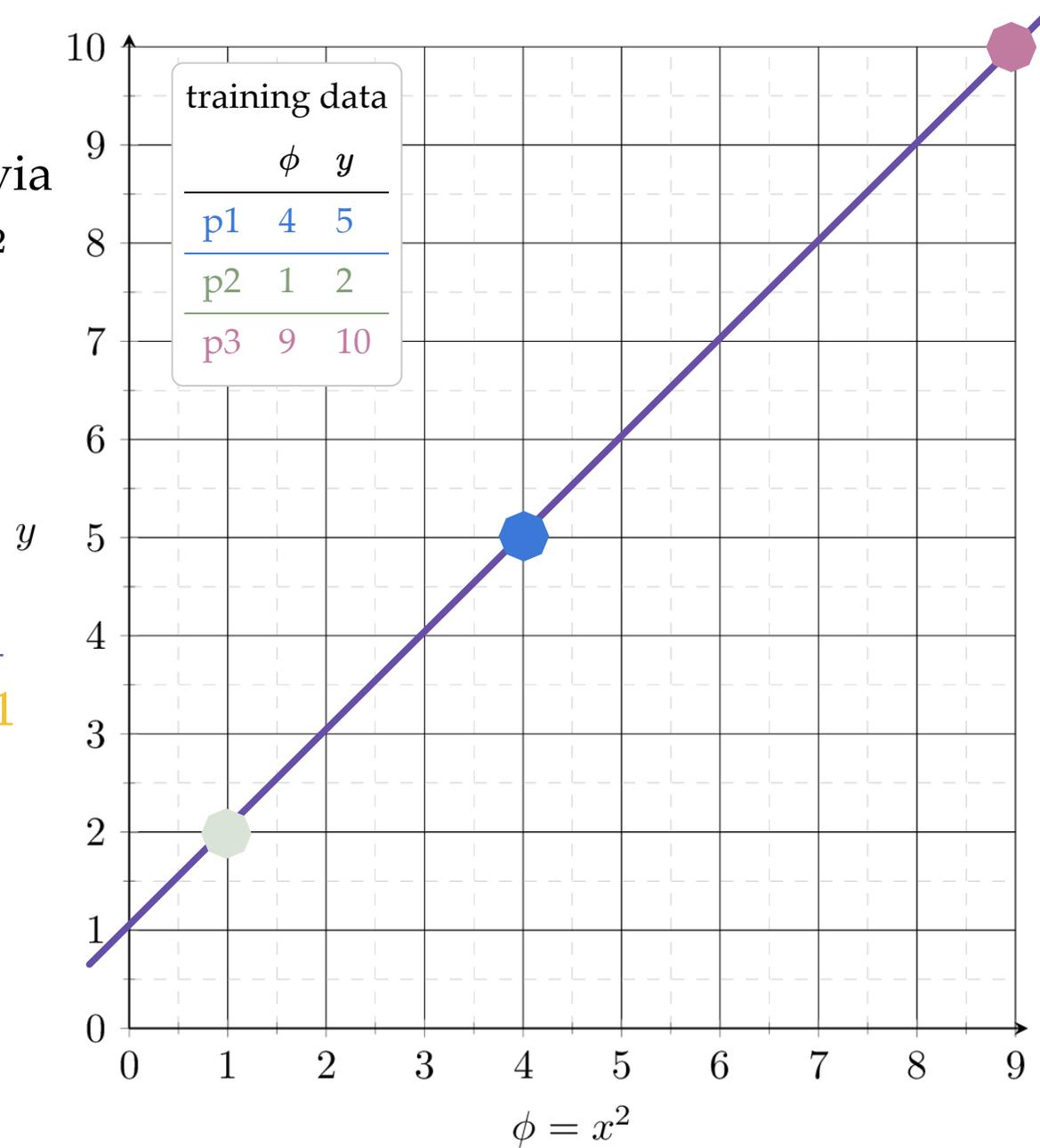


transform via

$$\phi(x) = x^2$$



$$g = \phi + 1$$
$$g = x^2 + 1$$



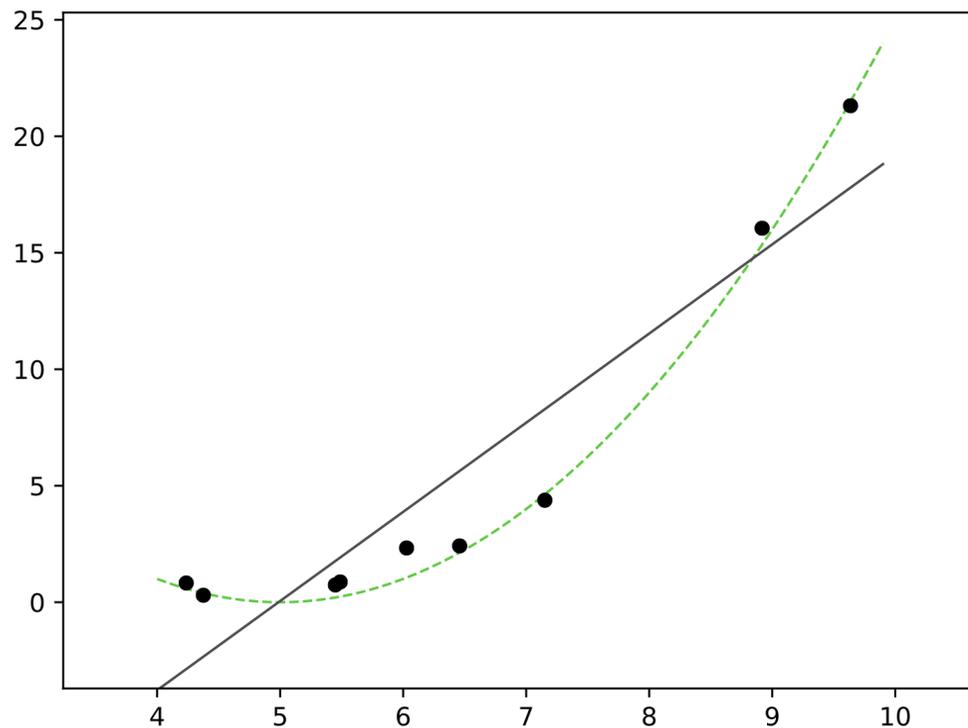
# systematic polynomial features construction

	$d = 1, \text{ features: } x_1$	$d = 2, \text{ features: } x_1, x_2$
$k = 0$	1	1
$k = 1$	1, $x_1$	1, $x_1, x_2$
$k = 2$	1, $x_1$ , $x_1^2$	1, $x_1, x_2$ , $x_1^2, x_1x_2, x_2^2$
$k = 3$	1, $x_1$ , $x_1^2$ , $x_1^3$	1, $x_1, x_2$ , $x_1^2, x_1x_2, x_2^2$ , $x_1^3, x_1^2x_2, x_1x_2^2, x_2^3$
$\vdots$	$\vdots$	$\vdots$

- Elements in the basis are the monomials of original features raised up to power  $k$
- With a given  $d$  and a fixed  $k$ , the basis is **fixed**.

- $n = 9$  data points, each with feature  $x \in \mathbb{R}$  and label  $y \in \mathbb{R}$
- data generated from green dashed curve

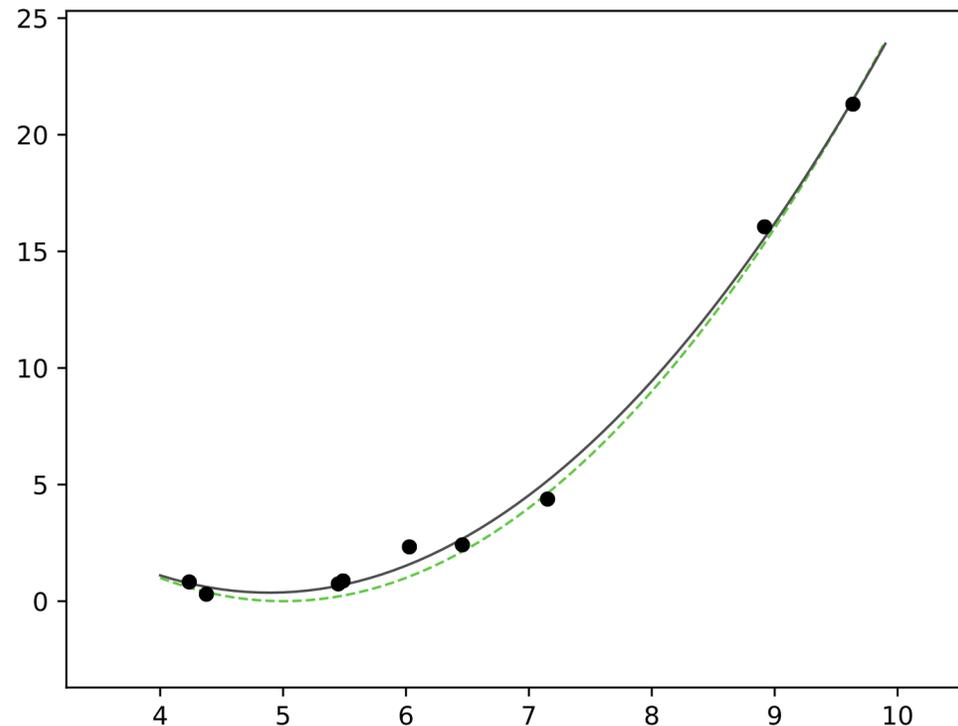
$k = 1$



$$h = \theta_0 + \theta_1 x$$

Learn 2 parameters — degree-1

$k = 2$

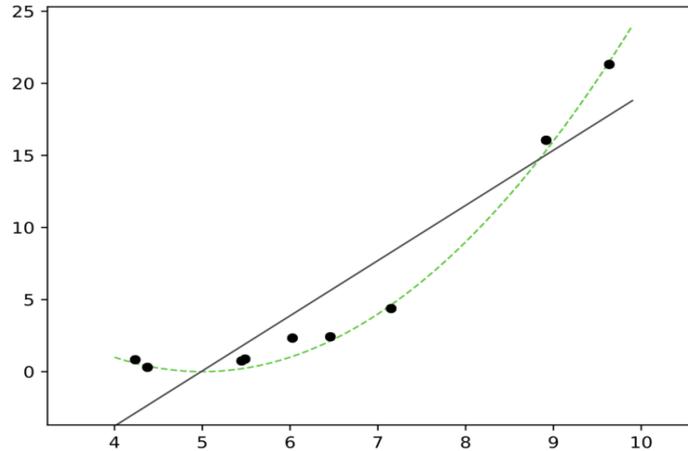


$$h = \theta_0 + \theta_1 x + \theta_2 x^2$$

Learn 3 parameters — degree-2

## Underfitting

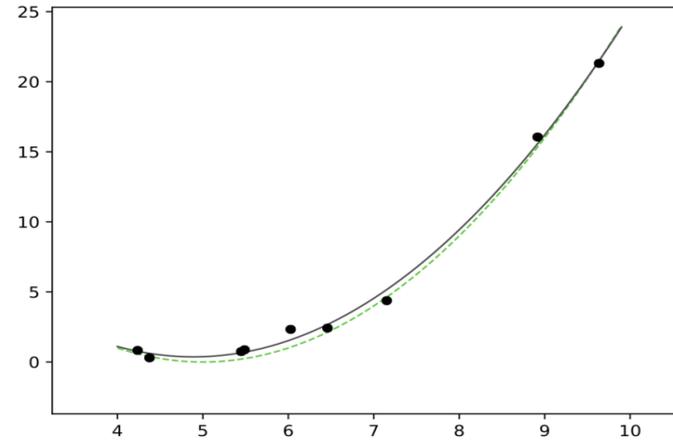
$$k = 1$$



high error on train set  
high error on test set

## Appropriate model

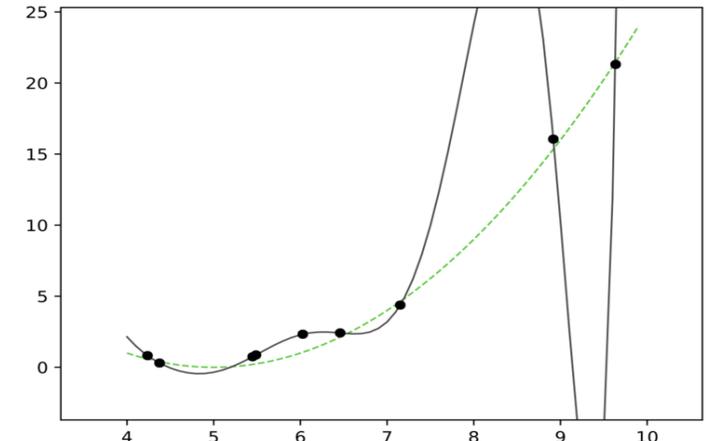
$$k = 2$$



low error on train set  
low error on test set

## Overfitting

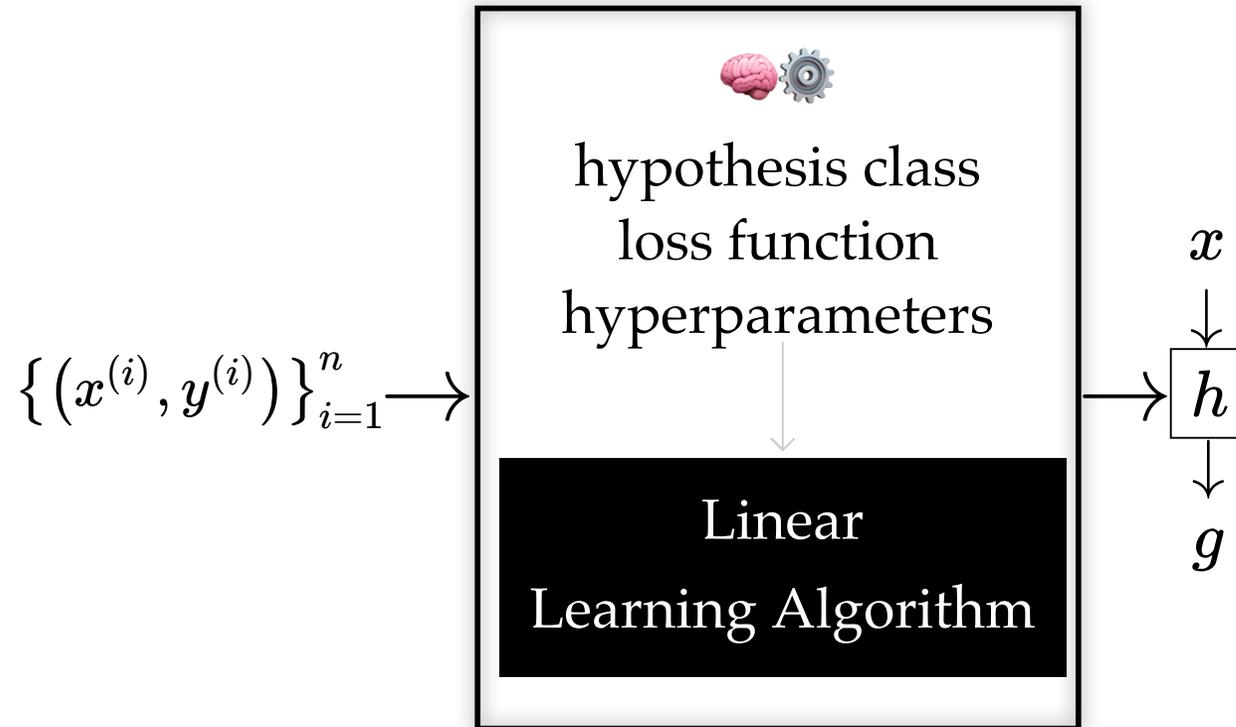
$$k = 10$$



low error on train set  
high error on test set

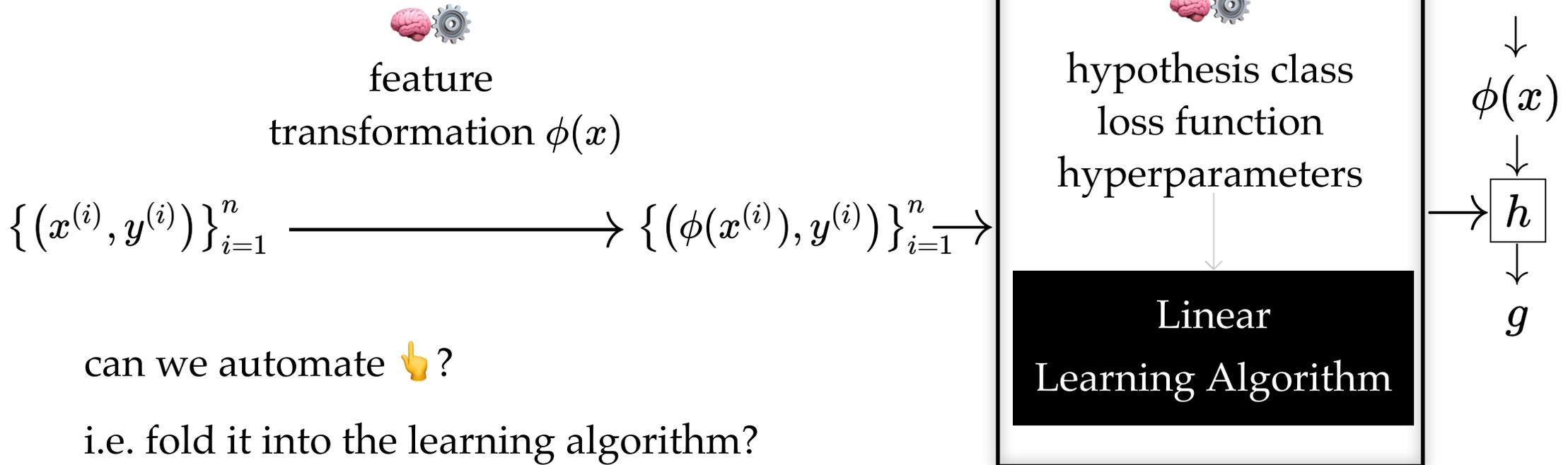
- $k$  : a hyperparameter that determines the capacity (expressiveness) of the hypothesis class.
- Models with many rich features and free parameters tend to have high capacity but also greater risk of overfitting.
- How to choose  $k$ ? Validation/cross-validation.

## Previously:

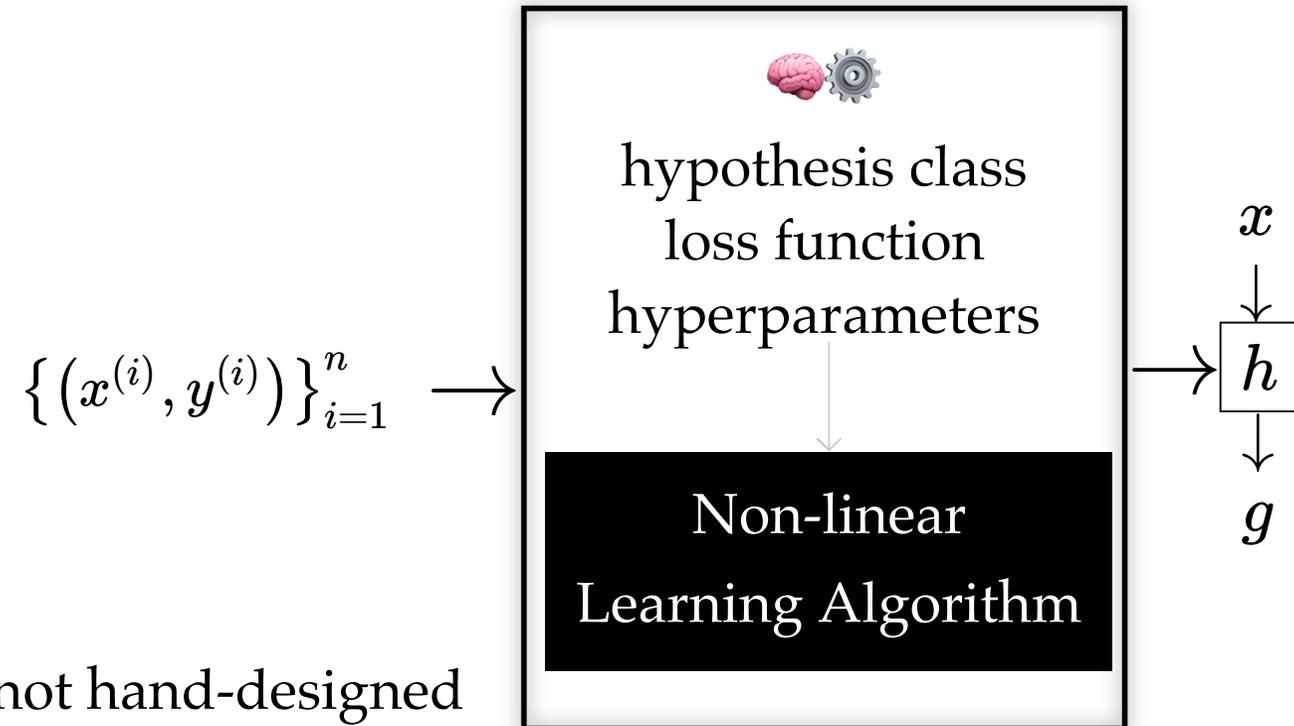


- Linear models are convenient but lack expressiveness for real-world tasks.
- We applied **fixed** nonlinear feature transformations, then used linear methods.
- The essence of ML was **feature engineering** — manually designing useful representations.

today, so far:



## neural networks:



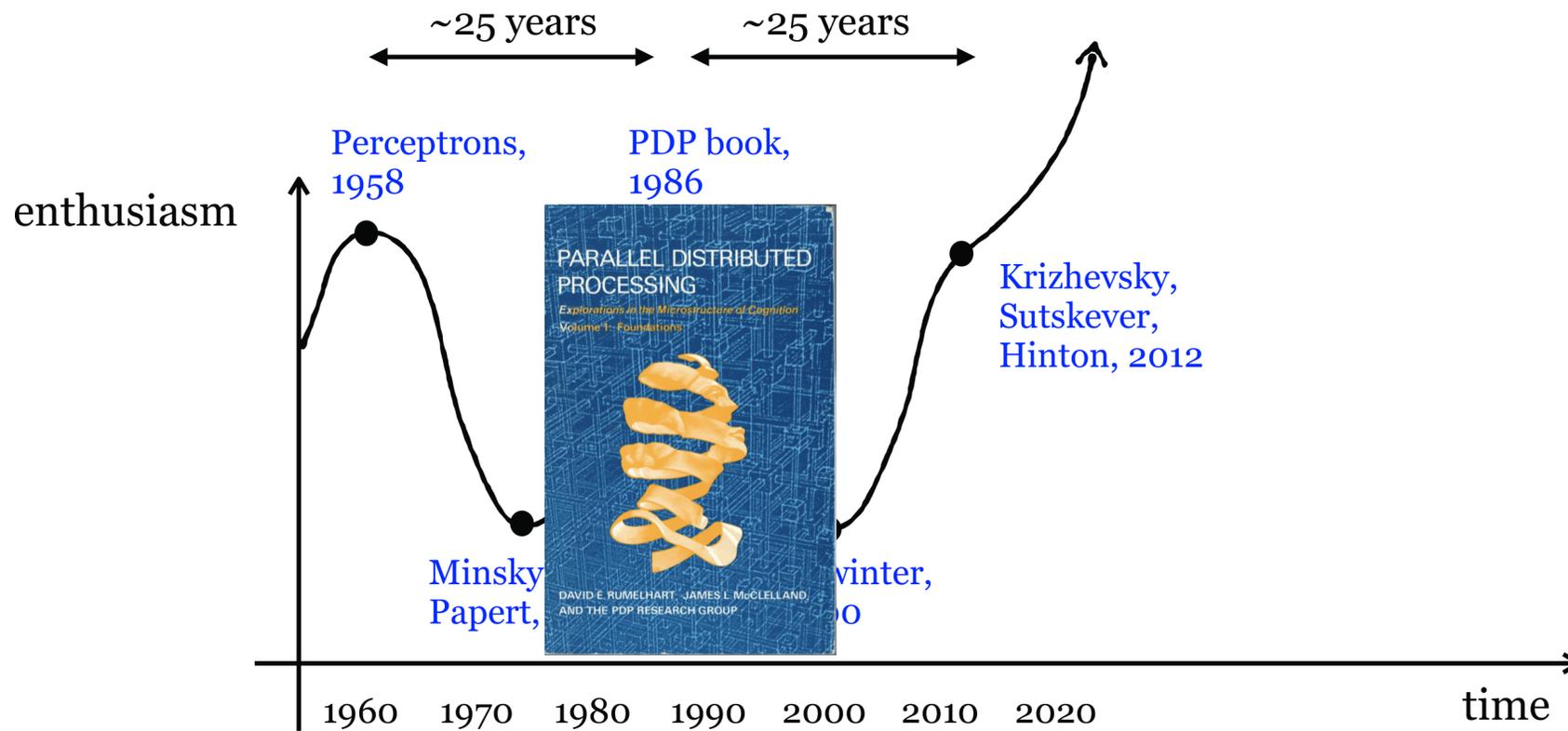
the nonlinearity  $\phi$  is learned not hand-designed

# Outline

- Systematic feature transformations
- Neural networks
  - Terminology
  - Design choices

👉 heads-up:

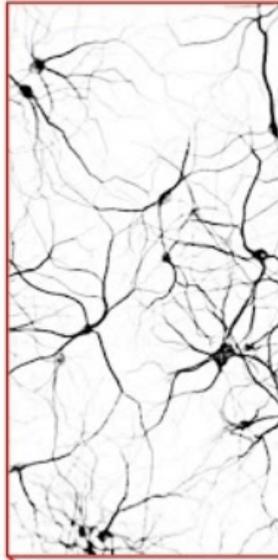
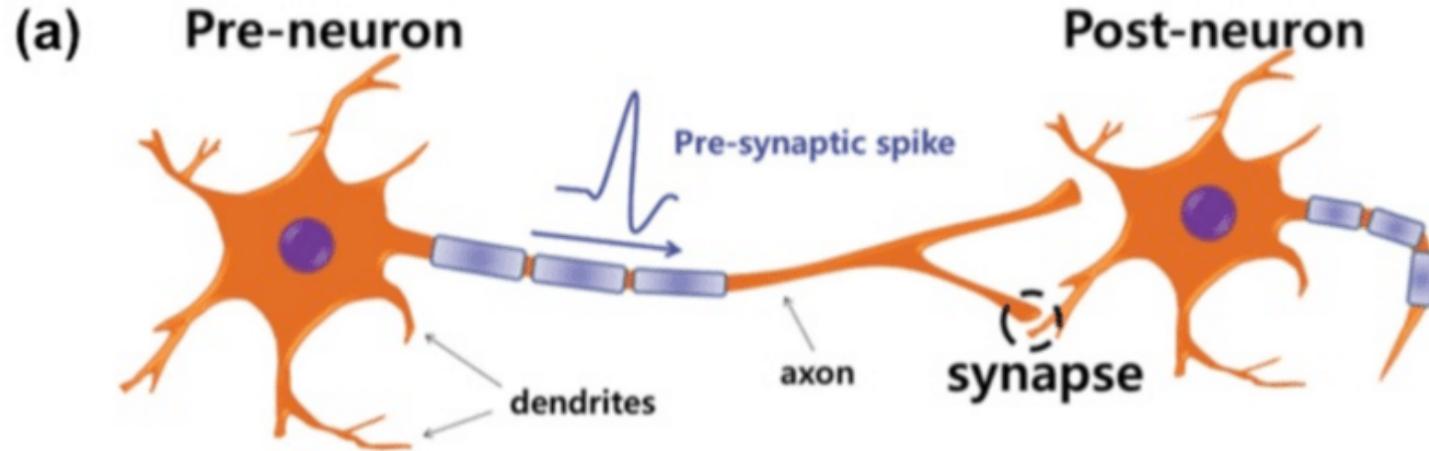
all neural network diagrams focus on a single data point



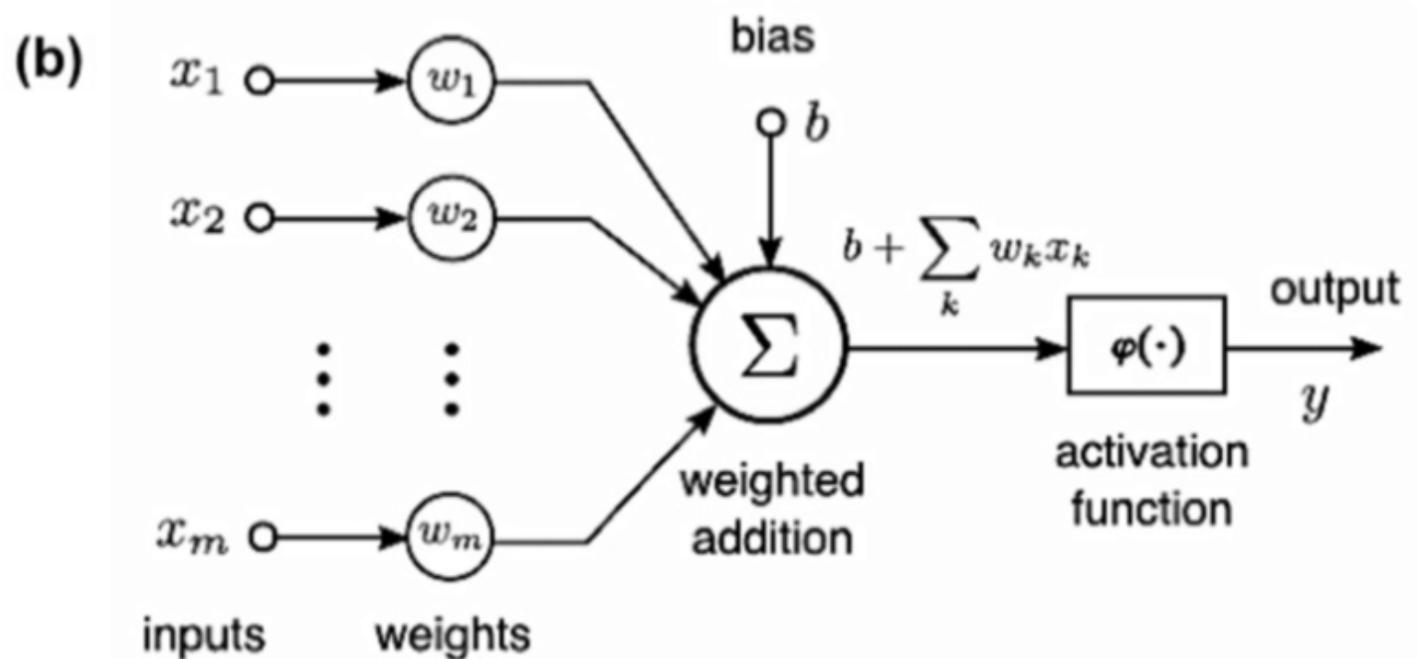
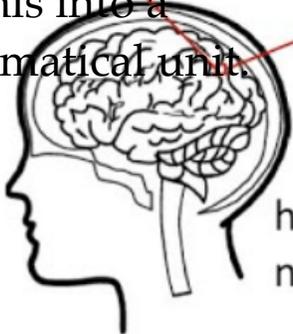
Outlined the fundamental concepts of neural networks:

- |                     |   |   |                |
|---------------------|---|---|----------------|
| layered structure { | <ul style="list-style-type: none"> <li>• Nonlinear feature transformation</li> </ul>                      | } | expressiveness |
|                     | <ul style="list-style-type: none"> <li>• Composing simple nonlinearities amplifies this effect</li> </ul> |   |                |
|                     | <ul style="list-style-type: none"> <li>• Backpropagation</li> </ul>                                       |   |                |

# Neurons and the brain

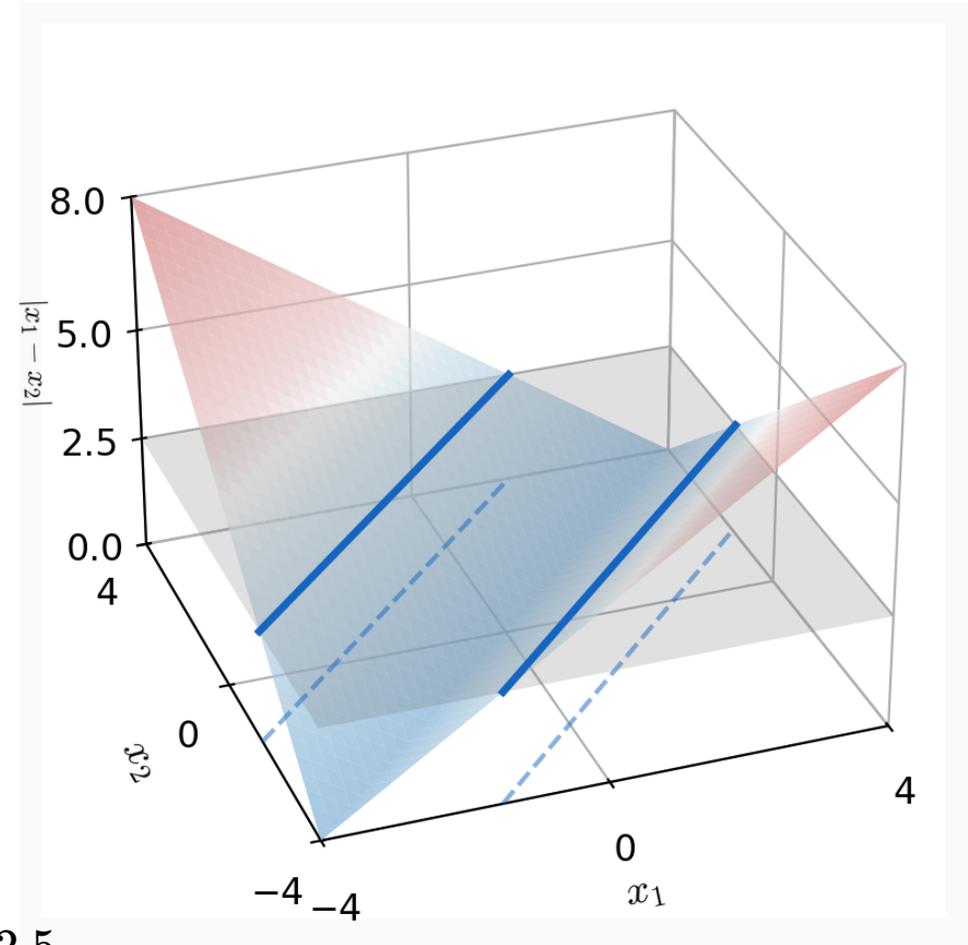
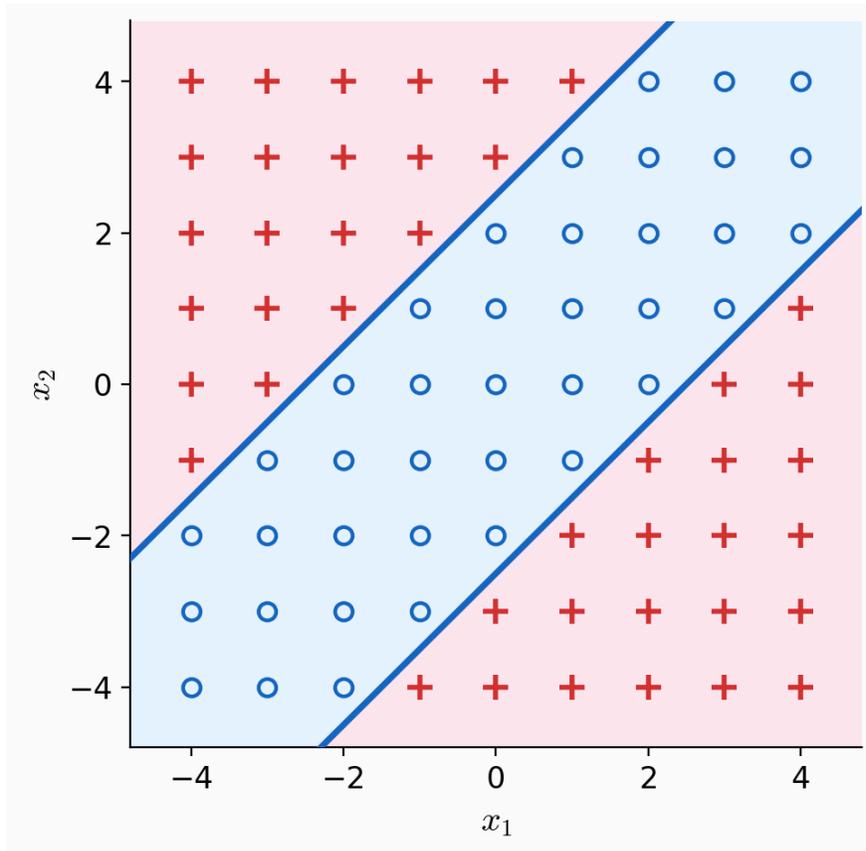


We abstract this into a simple mathematical unit.



- Nonlinear feature transformation

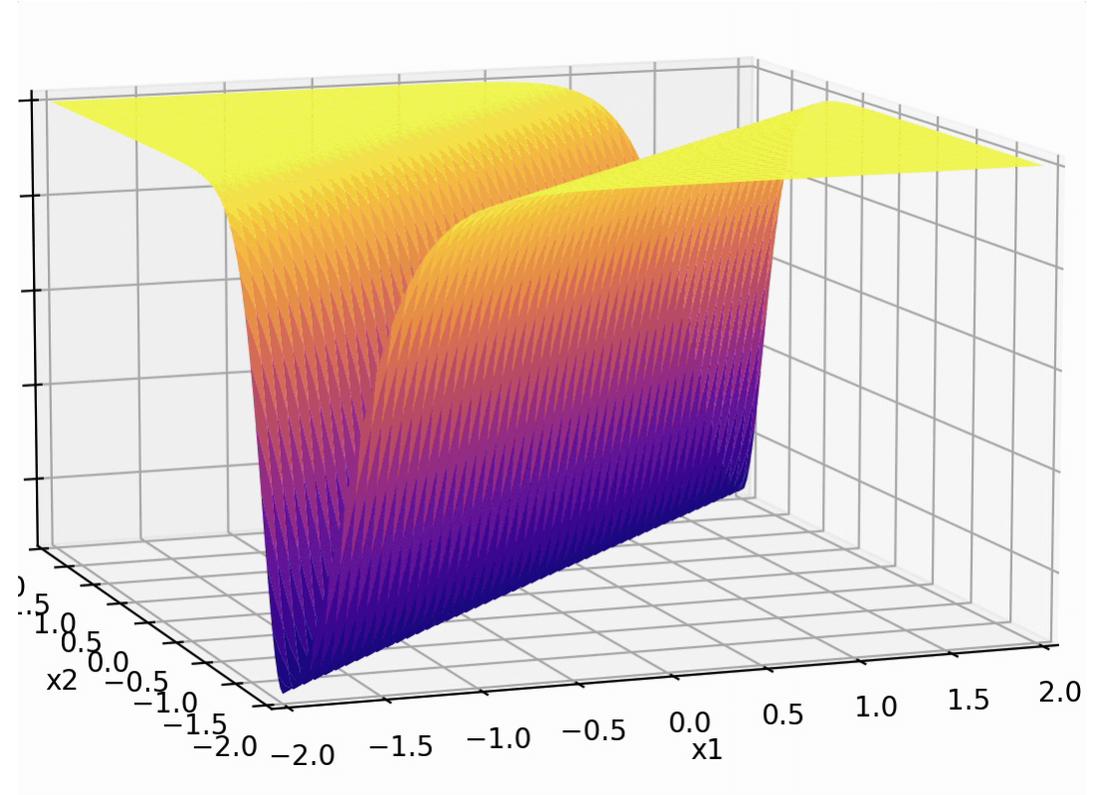
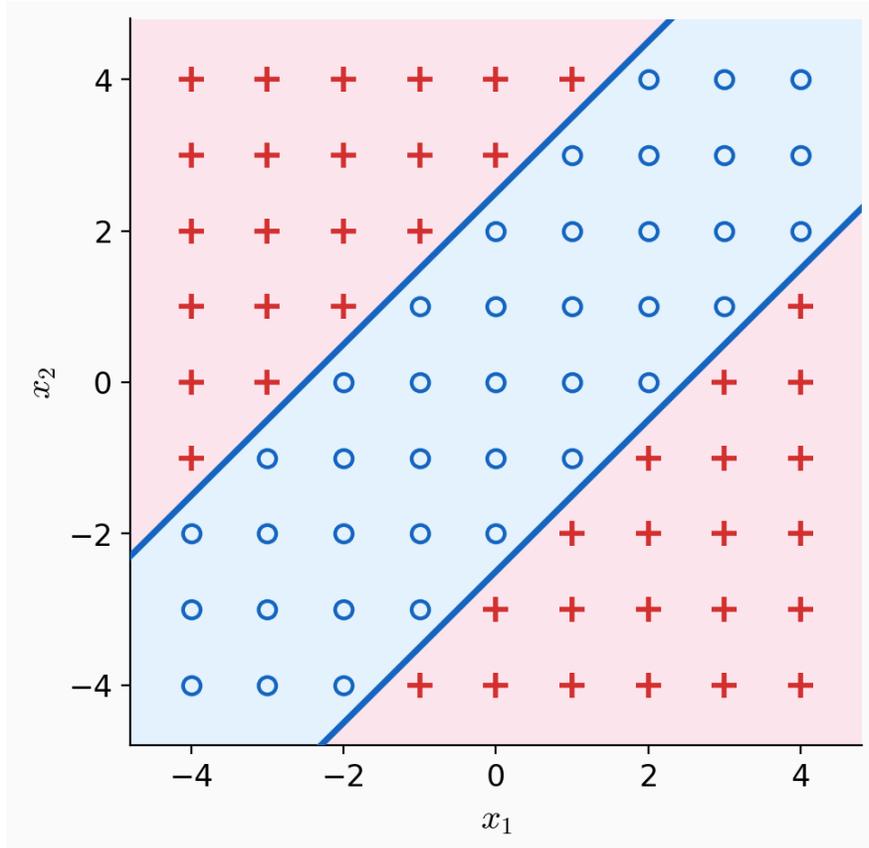
transform via  $\phi([x_1; x_2]) = [|x_1 - x_2|]$



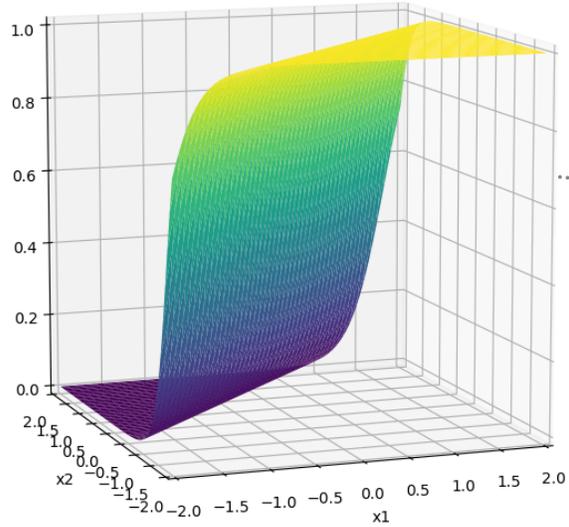
$$\phi([x_1; x_2]) = [|x_1 - x_2|] = 2.5$$

importantly, linear in  $\phi$ , non-linear in  $x$

- "Composing" simple transformations

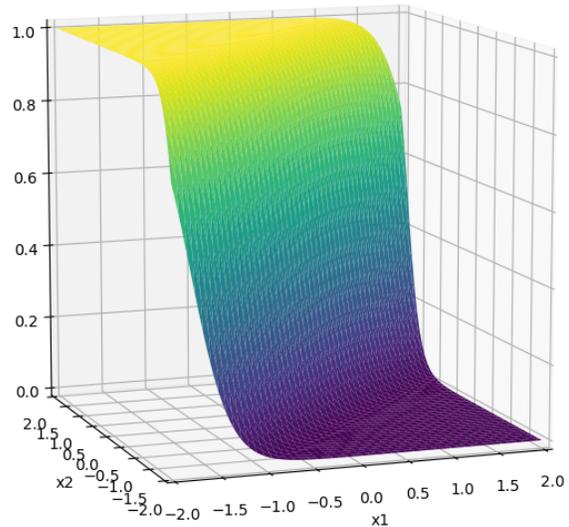


$$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$$



- "Composing" simple transformations

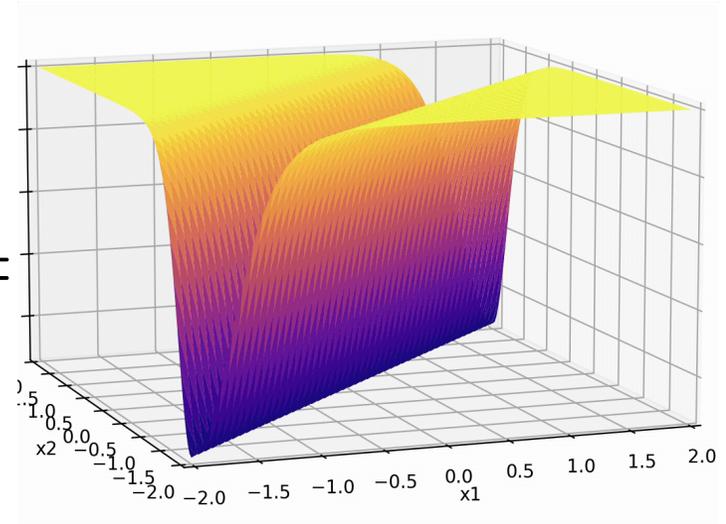
$$\sigma_2 = \sigma(-5x_1 + 5x_2 + 1)$$



some appropriately  
weighted sum

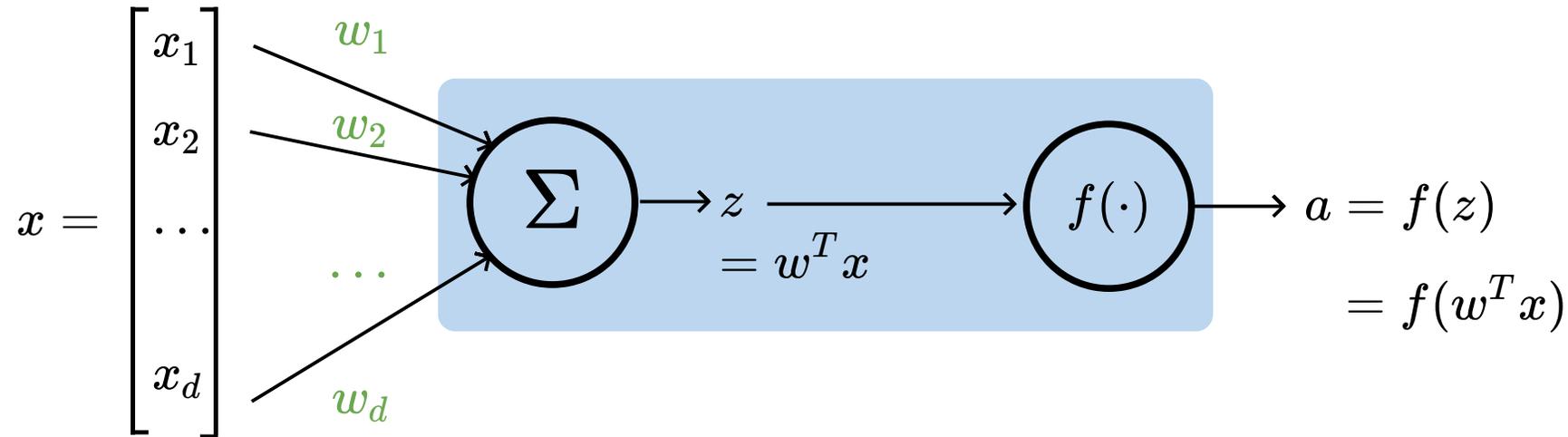


=



<https://shenshen.mit.edu/demos/2layers.html?embed=sigmoid>

## A neuron:



- $x$ :  $d$ -dimensional input
- $w$ : weights (i.e. parameters)
- $z$ : pre-activation output
- $f$ : activation function
- $a$ : post-activation output

$w$ : what the algorithm learns

$z$ : scalar

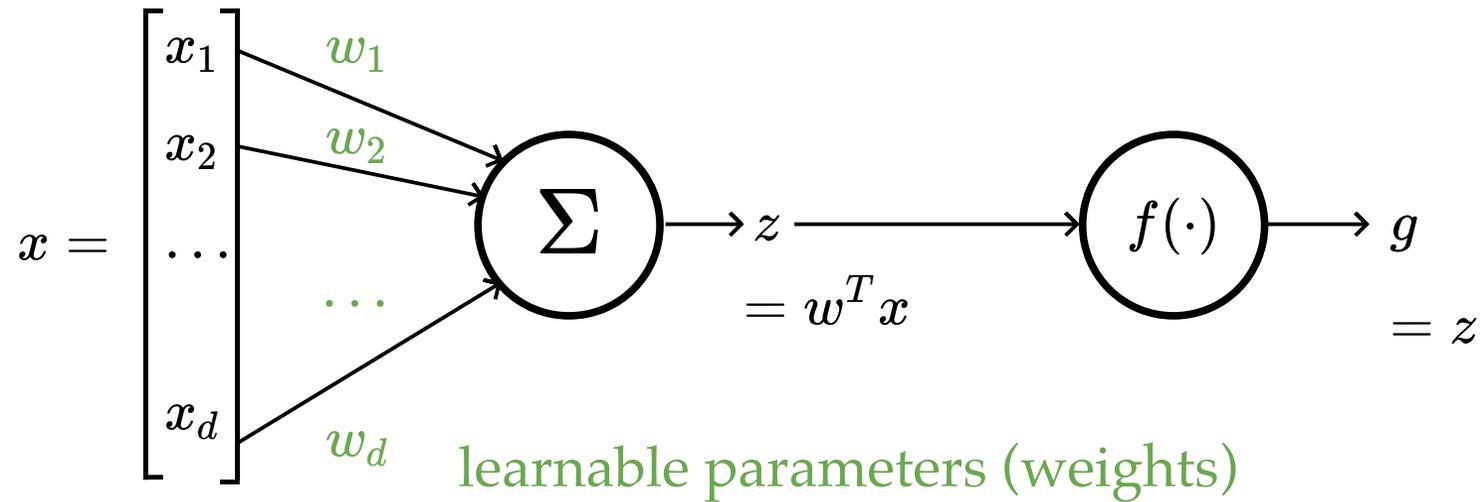
↓

$f$ : what we engineers choose

↓

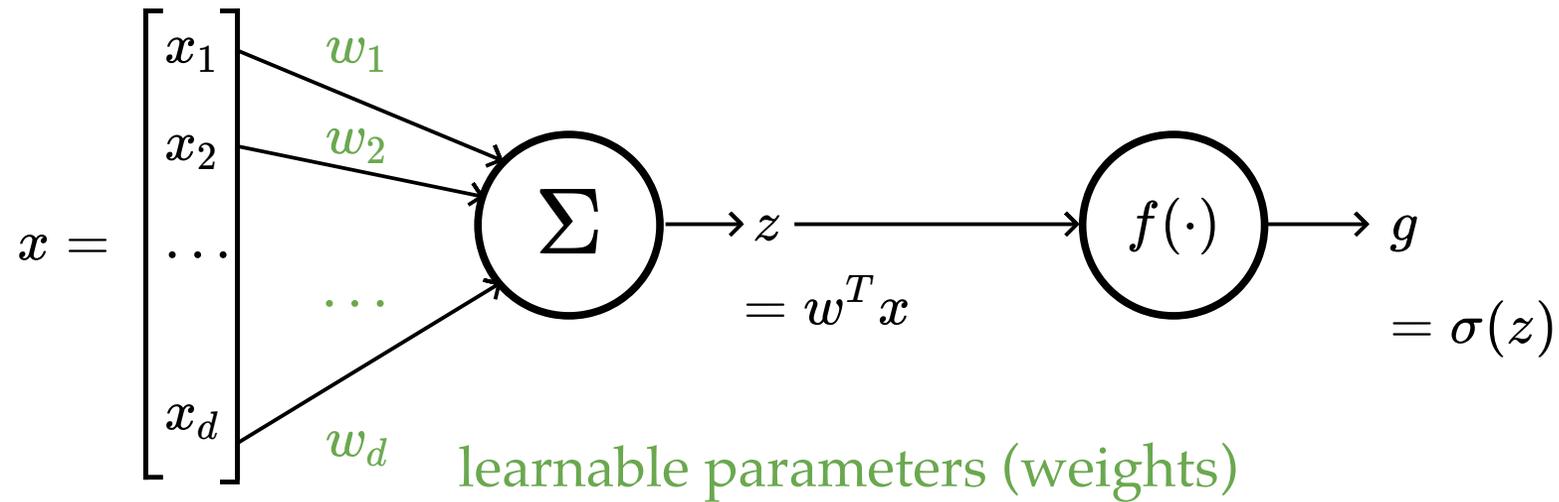
$a$ : scalar

e.g. linear regressor represented as a computation graph



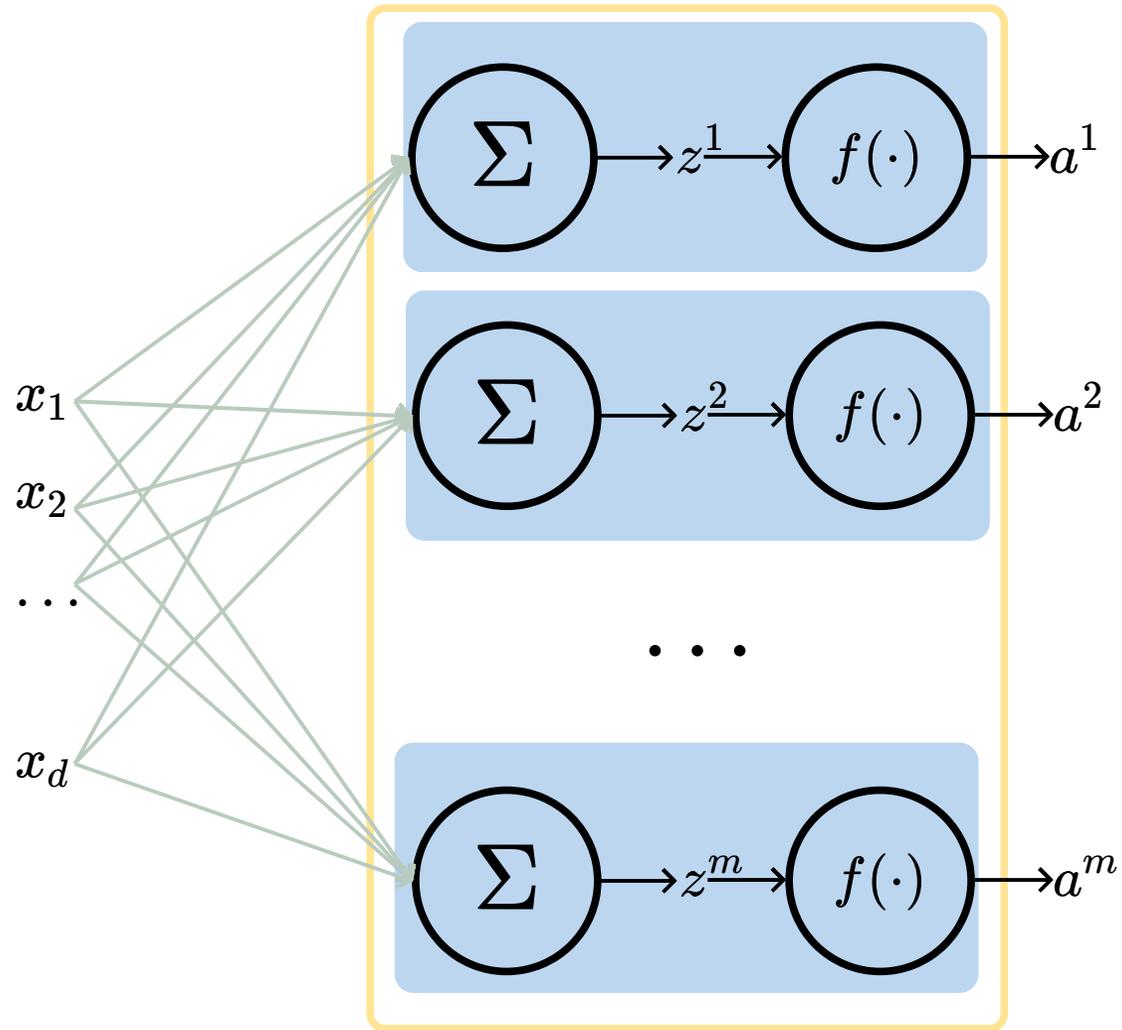
Choose activation  $f(z) = z$

e.g. linear logistic classifier represented as a computation graph



Choose activation  $f(z) = \sigma(z)$

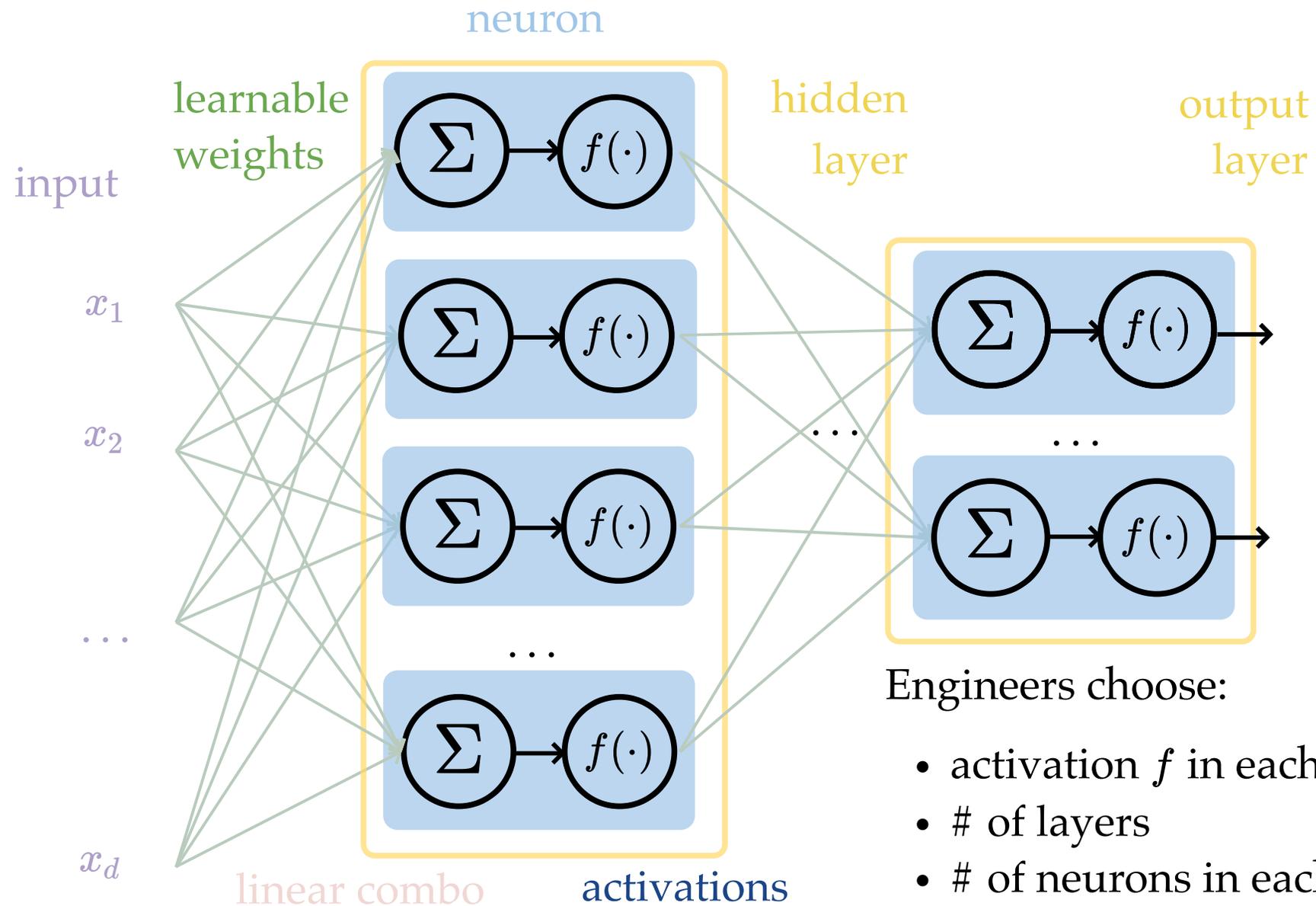
A layer:



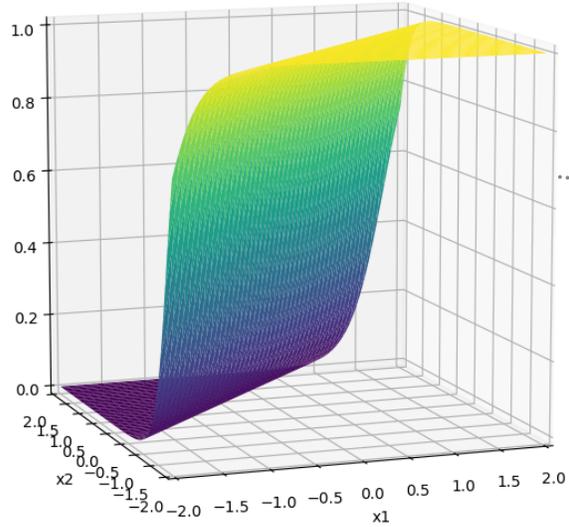
- (# of neurons) = layer's output dimension.
- activations applied element-wise, e.g.  $z^1$  won't affect  $a^2$ . (softmax in the output layer is a common exception)
- all neurons in a layer typically share the same  $f$  (mixed  $f$  is possible but complicates the math).
- layers are typically fully connected, each input  $x_i$  influences every  $a^j$ .

learnable weights

A (fully-connected, feed-forward) neural network:

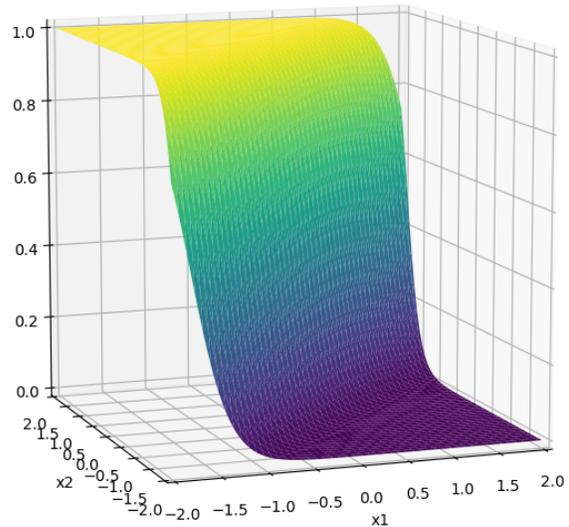


$$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$$



recall this example

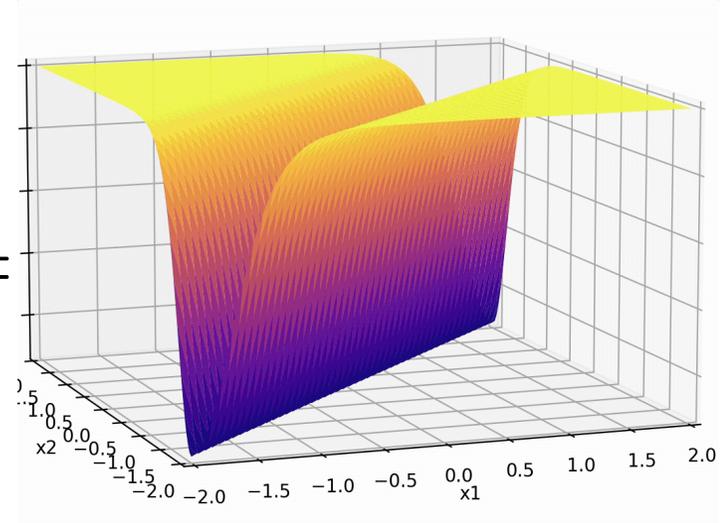
$$\sigma_2 = \sigma(-5x_1 + 5x_2 + 1)$$



some appropriately  
weighted sum

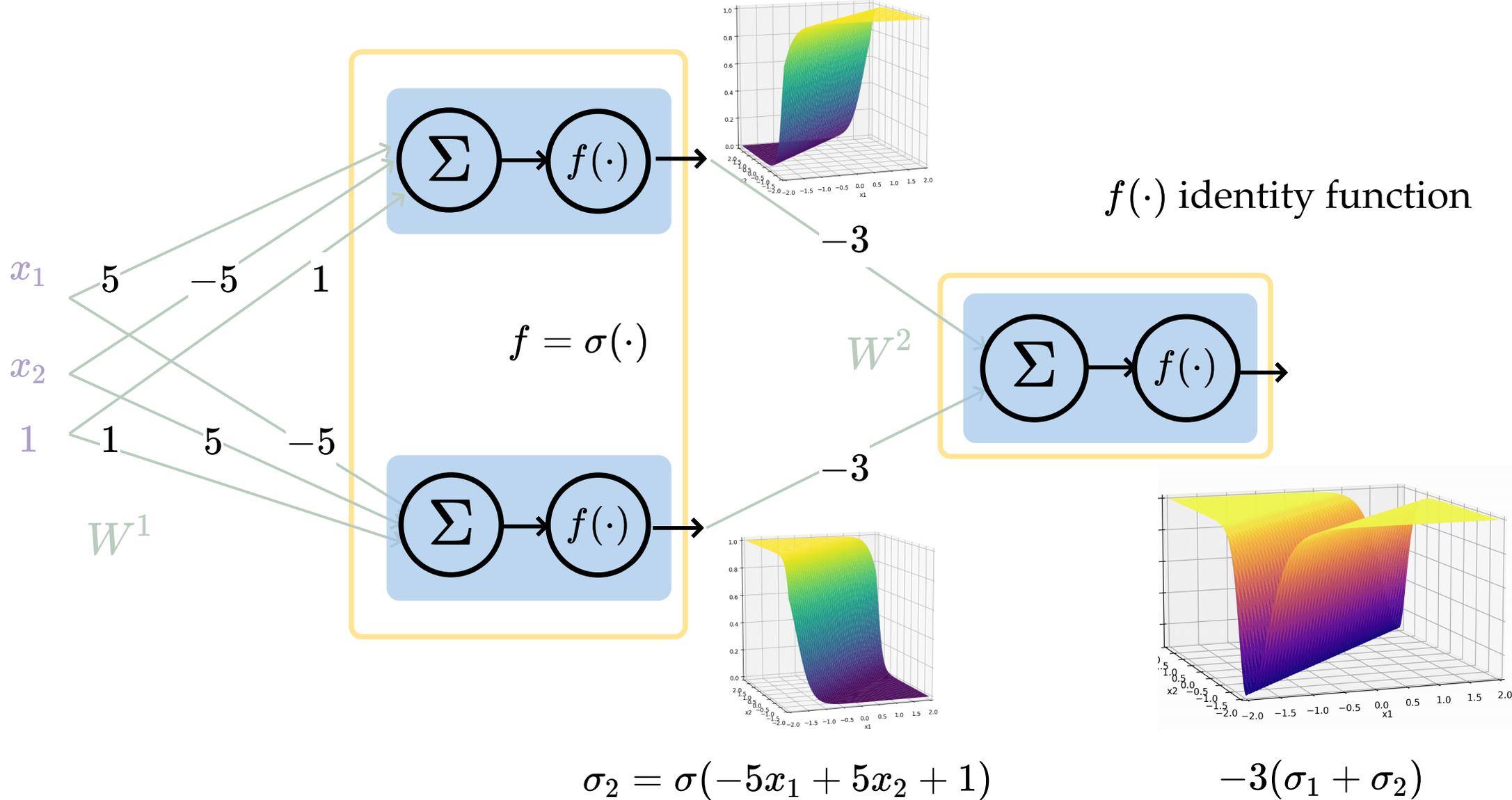


=

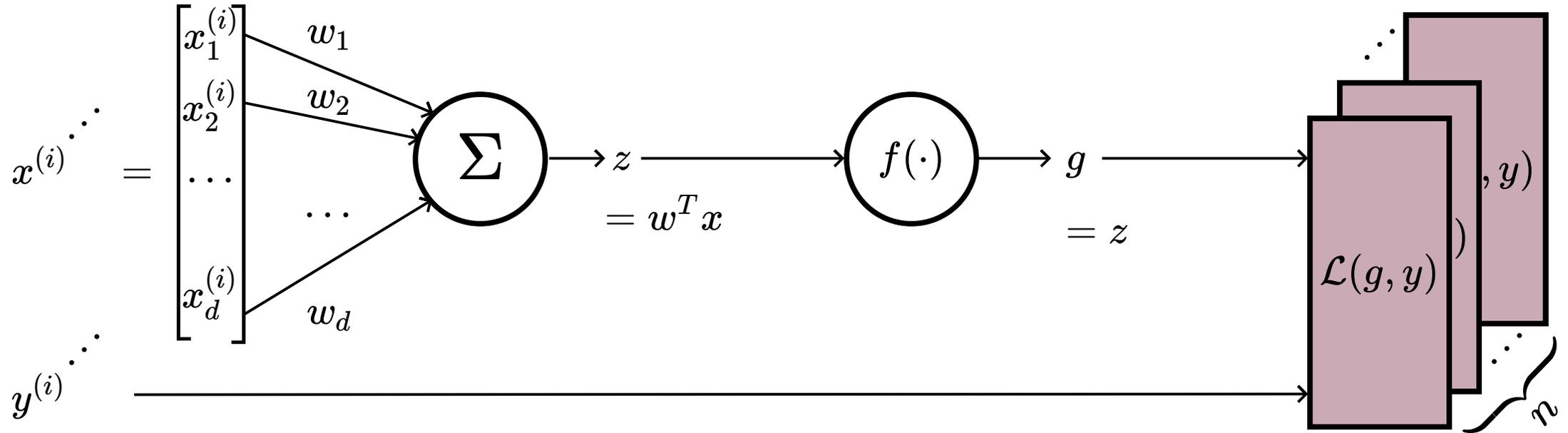


can be represented as:

$$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$$

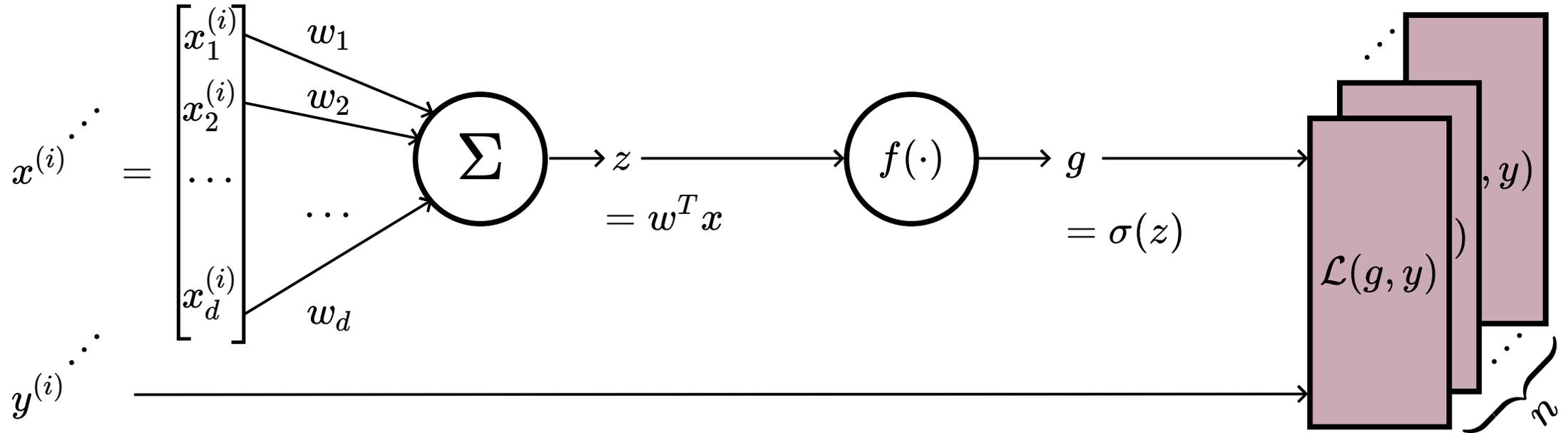


e.g. forward-pass of a linear regressor



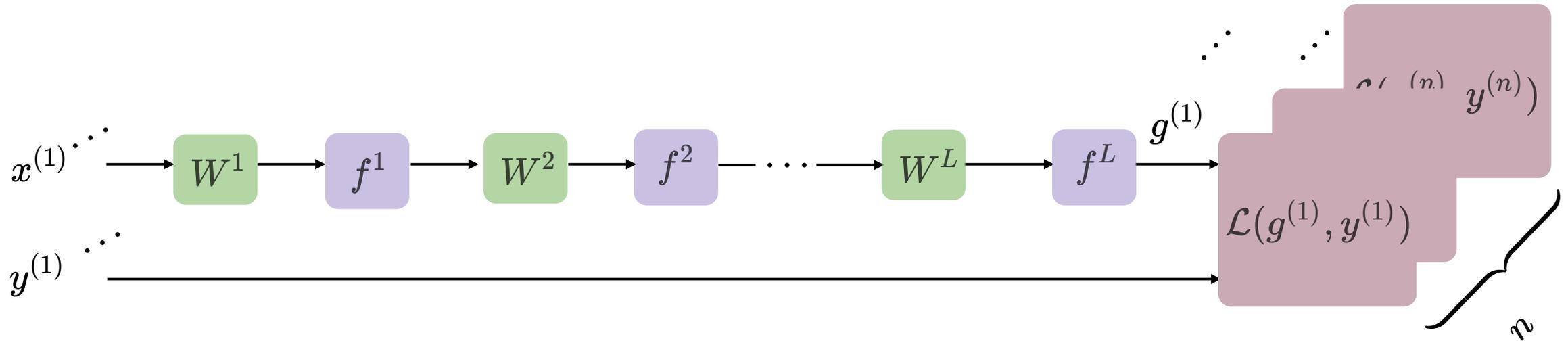
- Activation  $f$  is chosen as the identity function
- Evaluate the loss  $\mathcal{L} = (g^{(i)} - y^{(i)})^2$
- Repeat for each data point, average the sum of  $n$  individual losses

e.g. forward-pass of a linear logistic classifier



- Activation  $f$  is chosen as the sigmoid function
- Evaluate the loss  $\mathcal{L} = -[y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log (1 - g^{(i)})]$
- Repeat for each data point, average the sum of  $n$  individual losses

Forward pass: evaluate, *given* the current parameters



- the model outputs  $g^{(i)} = f^L (\dots f^2 ( f^1 (\mathbf{x}^{(i)}; \mathbf{W}^1); \mathbf{W}^2) ; \dots \mathbf{W}^L)$
- the loss incurred on the current data  $\mathcal{L}(g^{(i)}, y^{(i)})$
- the training error  $J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(g^{(i)}, y^{(i)})$

linear combination

(nonlinear) activation

loss function

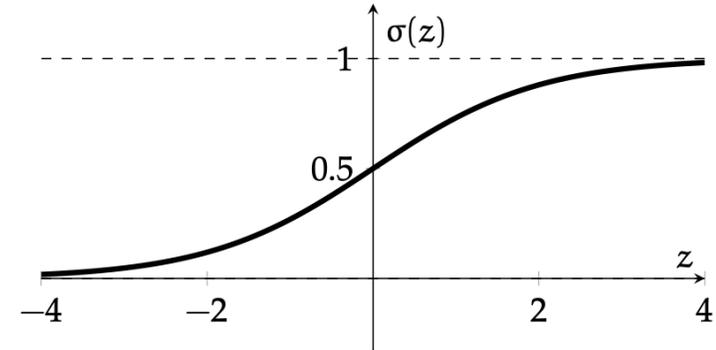
# Outline

- Systematic feature transformations
- Neural networks
  - Terminology
  - Design choices

## Hidden layer activation function $f$ choices

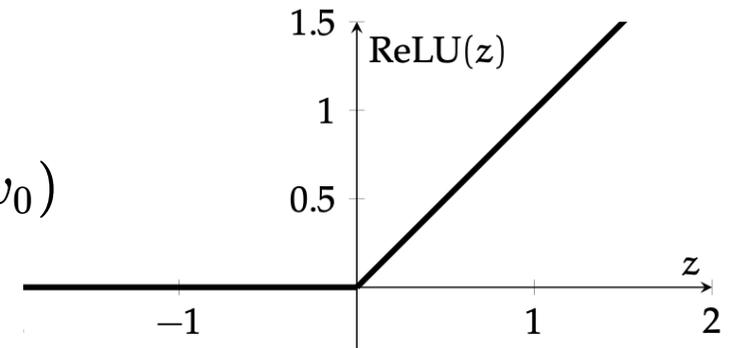
$\sigma$  used to be the most popular

- firing rate of a neuron
- elegant gradient  $\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$



nowadays, the default choice:

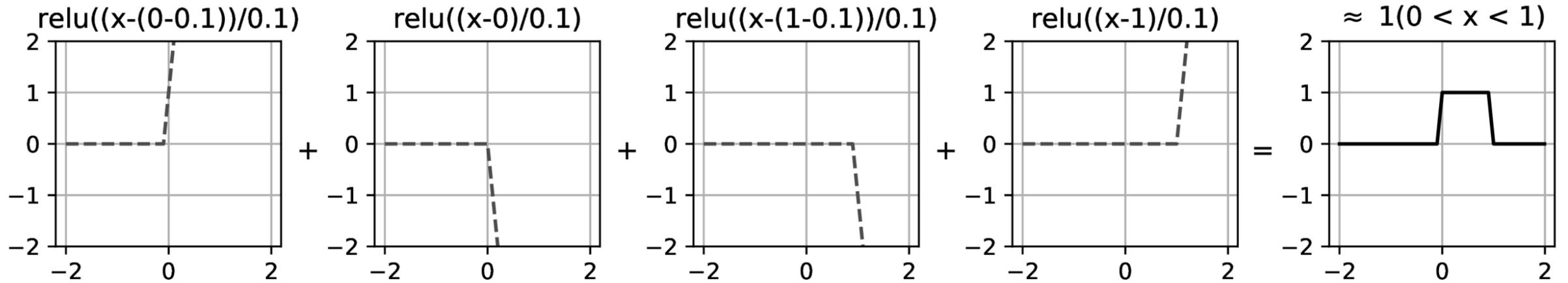
$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise} \end{cases} = \max(0, z) = \max(0, w^T x + w_0)$$



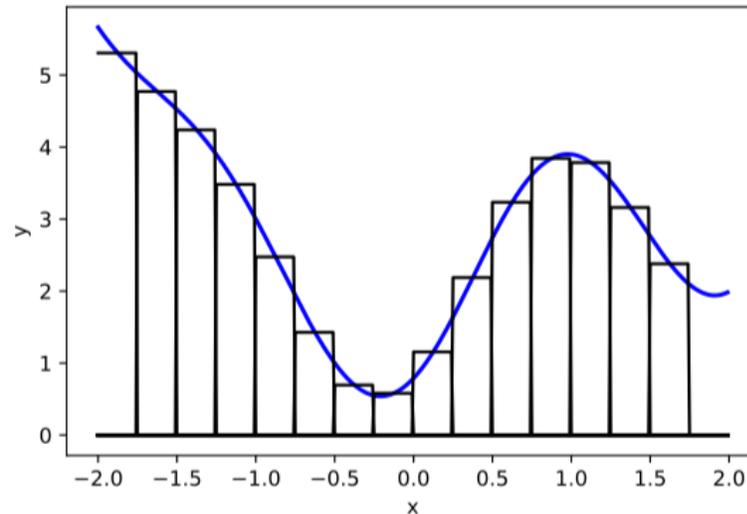
very simple function form (so is the gradient).

$$\frac{\partial \text{ReLU}(z)}{\partial z} = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{otherwise} \end{cases}$$

compositions of ReLU(s) can be quite expressive

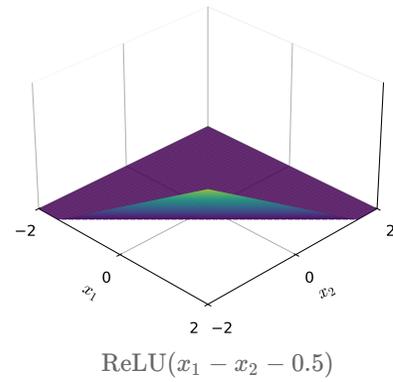
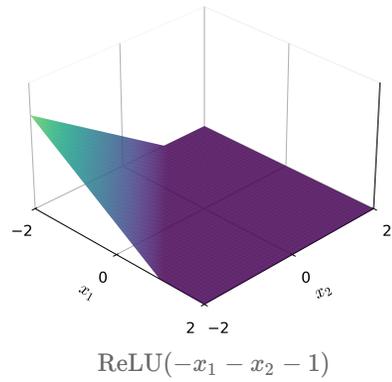


asymptotically, can approximate any continuous function arbitrarily well (for regression)



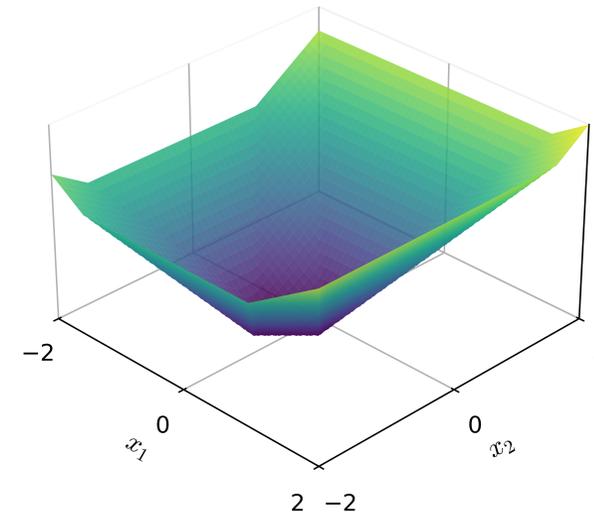
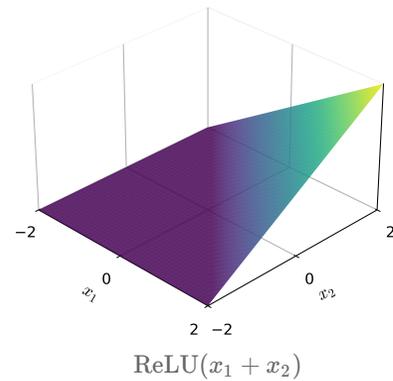
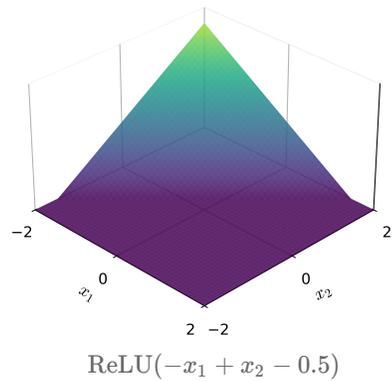
therefore can also approximate arbitrary decision boundaries (for classification)

compositions of ReLU(s) can be quite expressive



+

=



*just 4 ReLU neurons already  
carve out a non-trivial surface*

<https://shenshen.mit.edu/demos/2layers.html?embed=relu>



- Width: # of neurons in layers
- Depth: # of layers

- Typically, increasing either the width or depth (with non-linear activation) makes the model more expressive, but it also increases the risk of overfitting.

However, in the realm of neural networks, the precise nature of this relationship remains an active area of research—for example, phenomena like the double-descent curve and scaling laws

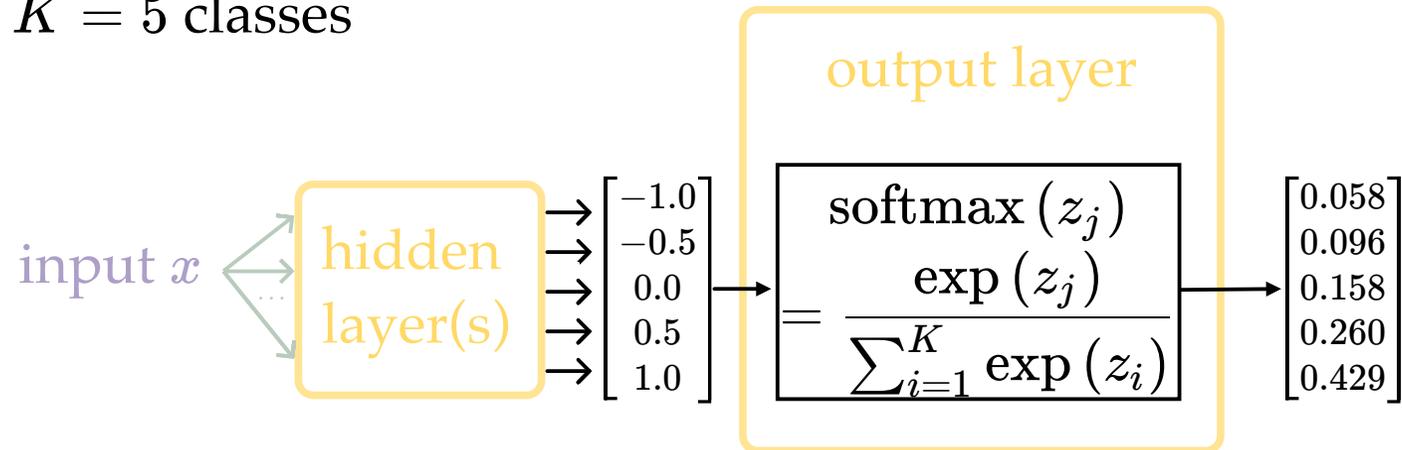
(The demo won't embed in PDF. But the direct link below works.)

<https://shenshen.mit.edu/demos/nn/playground/>

## output layer design choices

- # neurons, activation, and loss depend on the high-level goal.
- typically straightforward.
- Multi-class setup: if predict *one and only one* class out of  $K$  possibilities, then  
output layer:  $K$  neurons, softmax activation, cross-entropy loss

e.g., say  $K = 5$  classes



- For other multi-class settings, see hw/lab.

# Summary

- Linear models are convenient but lack expressiveness for most real-world tasks.
- A **fixed** non-linear feature transformation lets us use linear methods on complex problems, but designing features by hand gets tedious.
- Neural networks automate feature learning: layers alternate parameterized linear maps with non-linear activations (typically ReLU).
- For classification outputs, we apply sigmoid (binary) or softmax (multi-class).

**Next time:** How do we *train* neural networks? (gradient descent + backpropagation)