

Deep learning and Neuroscience

Emlyon school of business

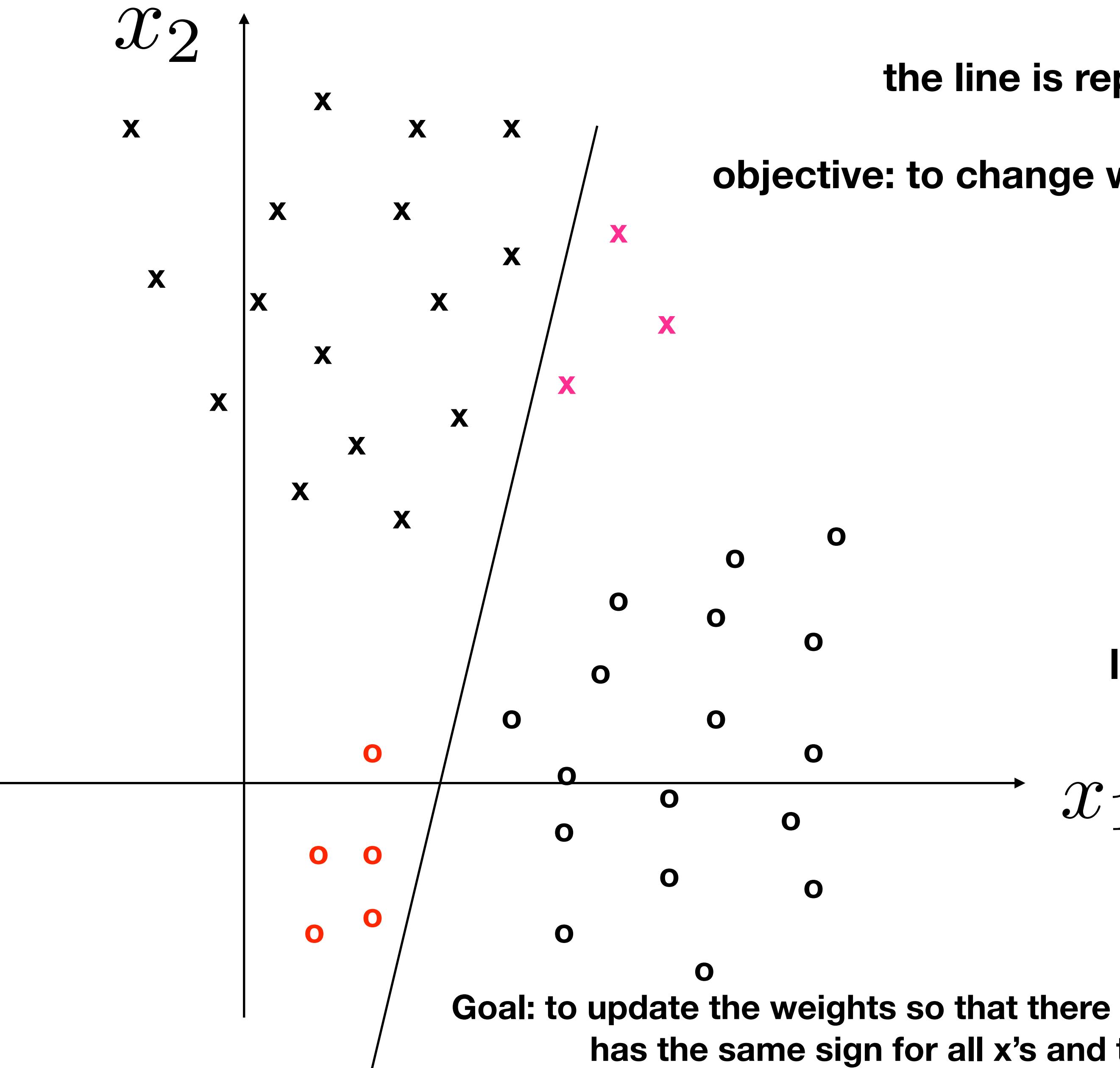
May 2020

Farzad Fathizadeh

*Swansea University and MPI for
Biological Cybernetics*

Course contents

Linear perceptrons, linear regression, logistic regression, deep neural nets and back propagation, convolutional neural nets, recurrent neural nets, large scale statistical inference and signal detection, fast Fourier transform, discussions on geometric dimension reduction, crash course on python and tensorflow.

x_2 

the line is represented by a vector w

objective: to change w in a way that there is less error.

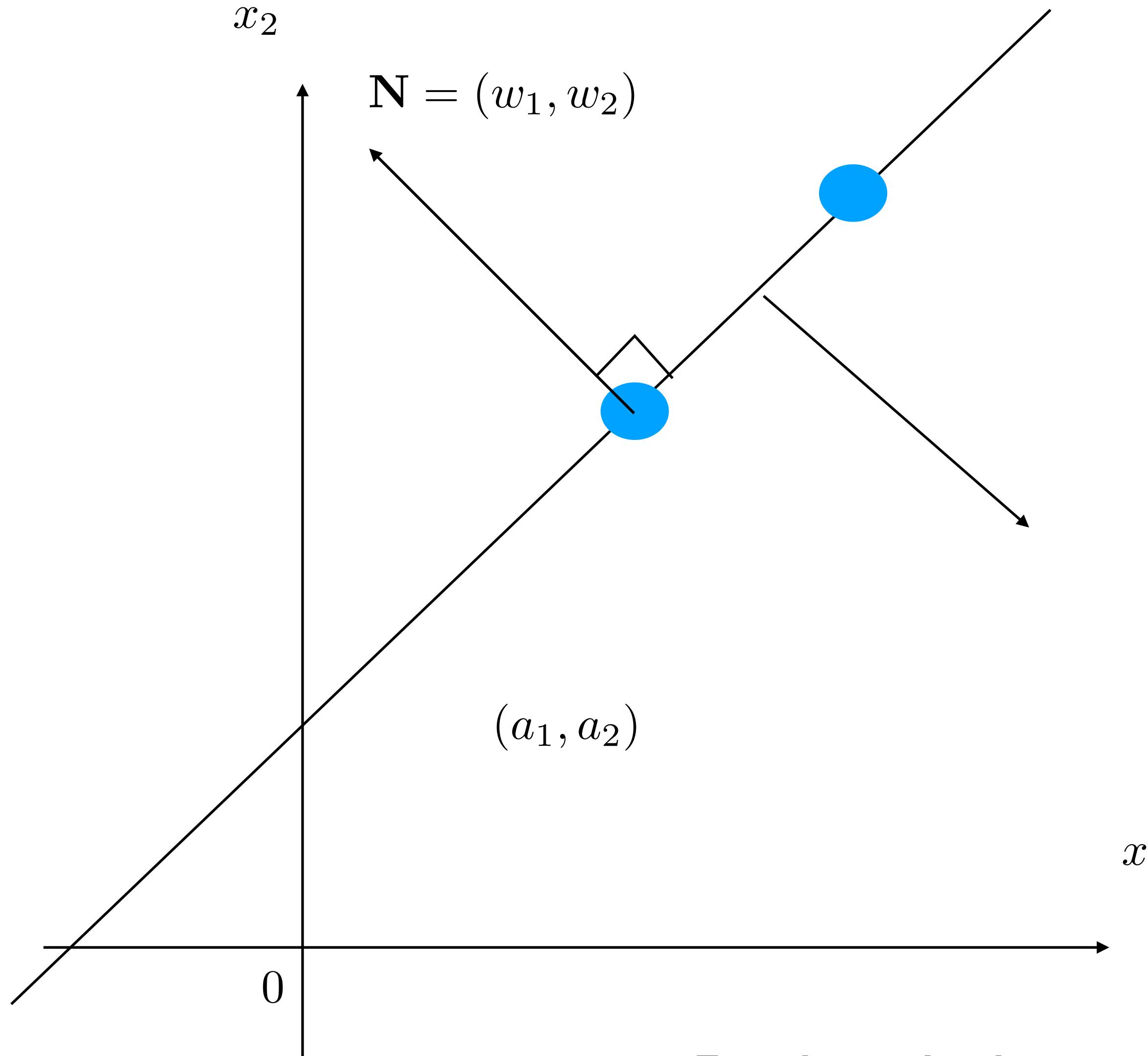
Linear perceptron

red: misclassified points

This is an example of
linearly separable data.

$$w = (w_0, w_1, w_2)$$

Goal: to update the weights so that there is no misclassification: meaning $w \cdot z$ has the same sign for all x's and the same sign for all o's



**The equation of a line determined by:
two points on the line
or
the slope and one point on the line
or
by one point on the line and the normal vector**

$$y = m x + b$$

$$x_2 = m x_1 + b$$

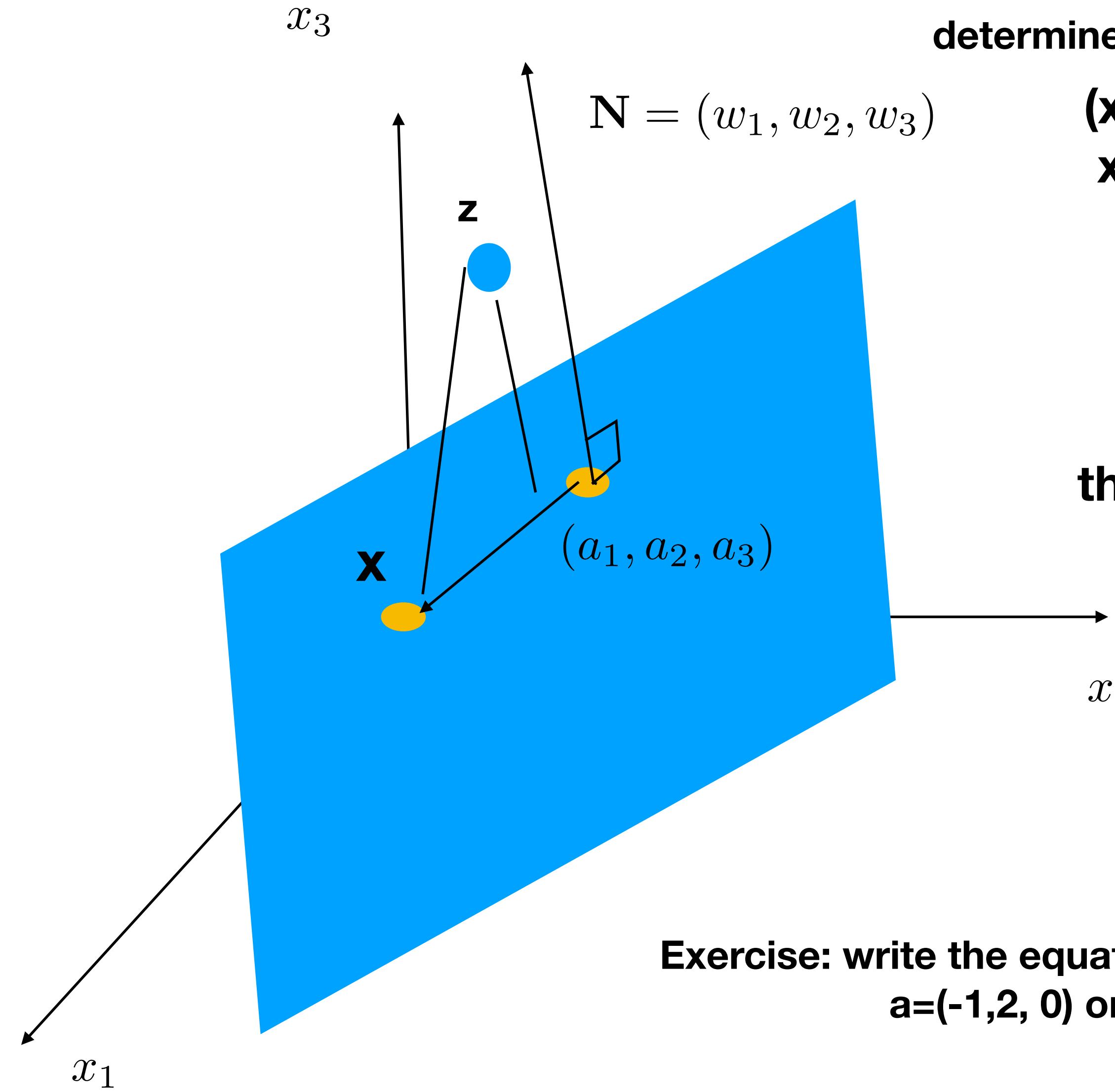
$$\mathbf{w} = (w_0, w_1, w_2)$$

Exercise: write the equation of the line with normal vector $\mathbf{N}=(3, 4)$ with the point $(-1,2)$ on it. Then identify 3 different points on the line.

$$25 = 16 + 9$$

$$5^2 = 4^2 + 3^2$$

The equation of a plane in 3 dimensional space:



determined by one point on the plane and its normal vector

$$(x-z) \cdot N = + \text{ or } - \text{ the distance} = |N|=1$$

$$x \cdot N - z \cdot N = -w_0 - z \cdot N = w \cdot z \quad z=(z_0=1, z_1, \dots, z_n)$$

$$x = (x_1, x_2, x_3)$$

$$a = (a_1, a_2, a_3)$$

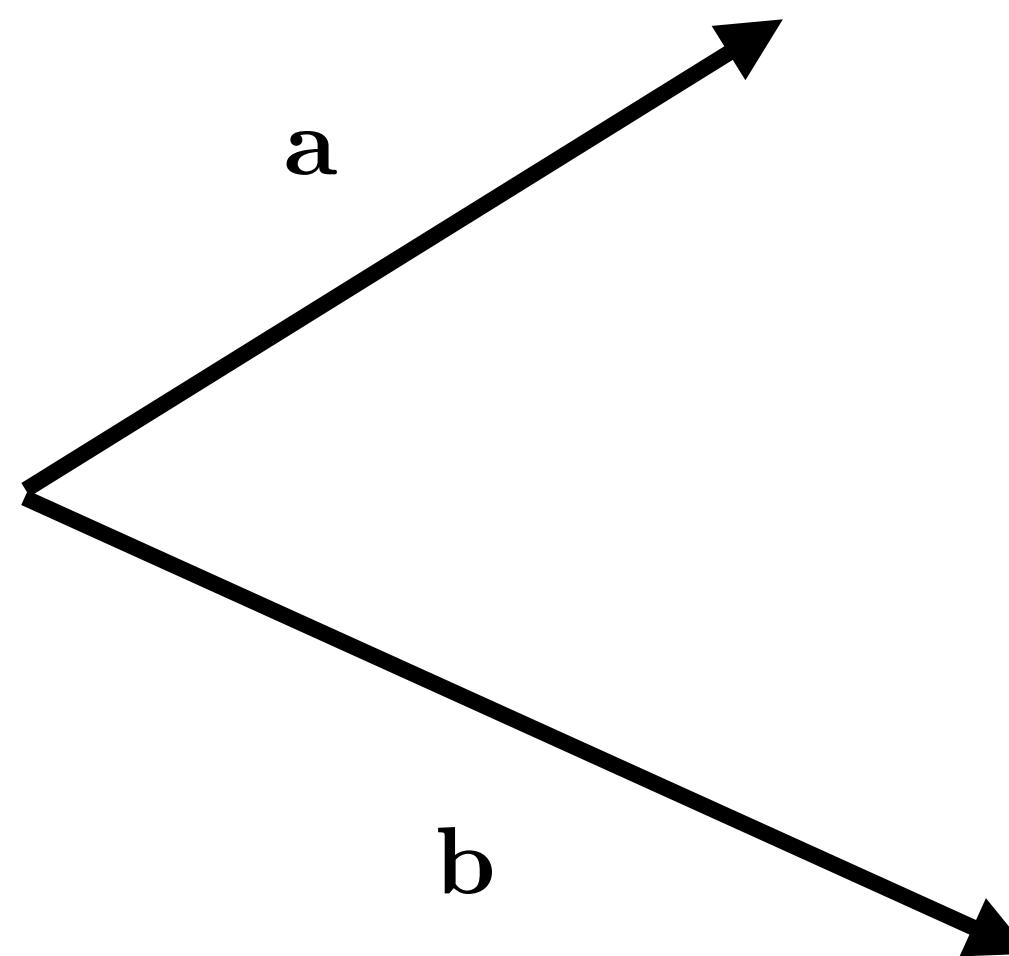
the equation for all x on the surface

$$(x - a) \cdot N = 0$$

orthogonal

Exercise: write the equation of the plane with normal vector $N=(3, 4, 1)$ with the point $a=(-1, 2, 0)$ on it. Then identify 3 different points on the plane.

Calculus of vectors

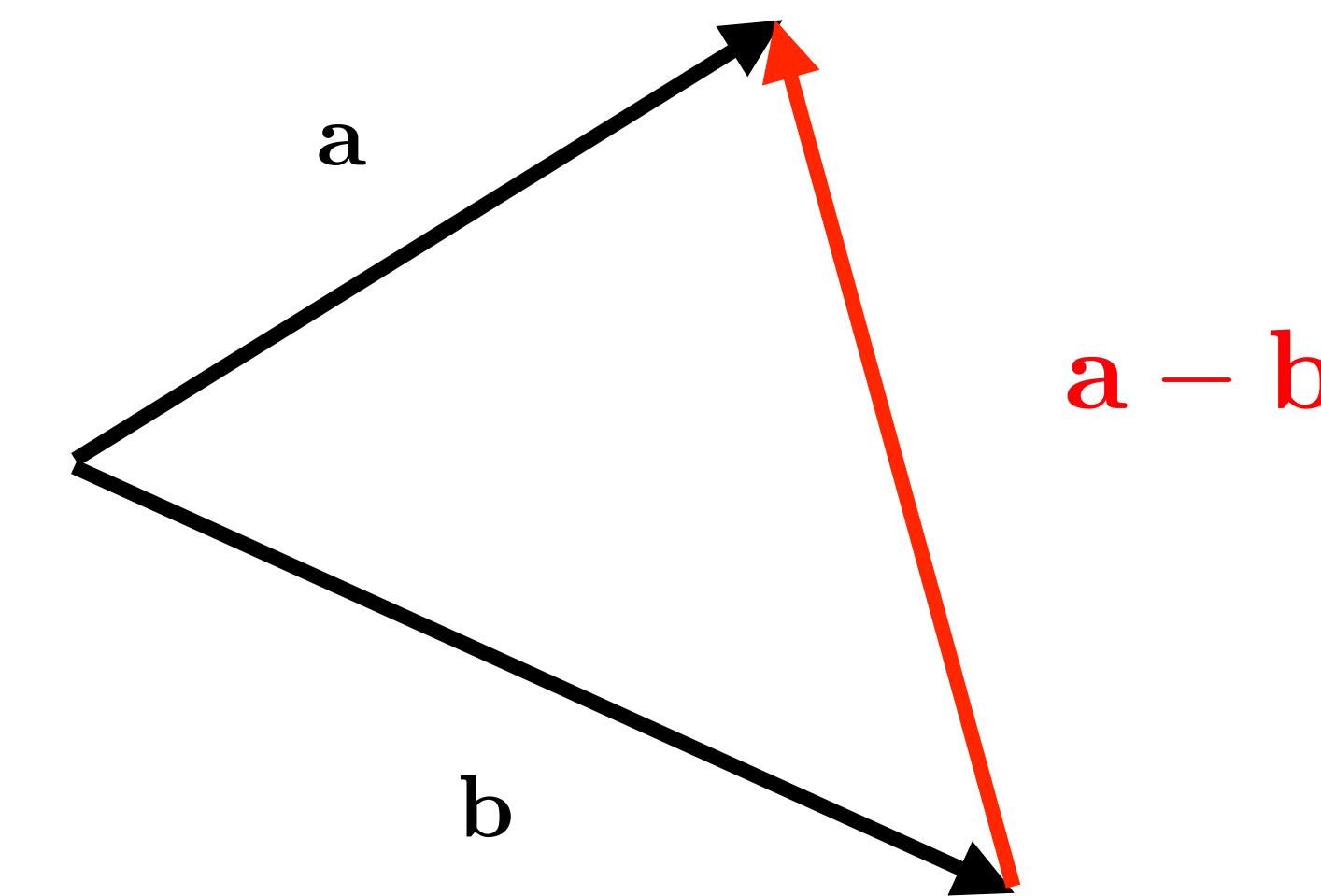
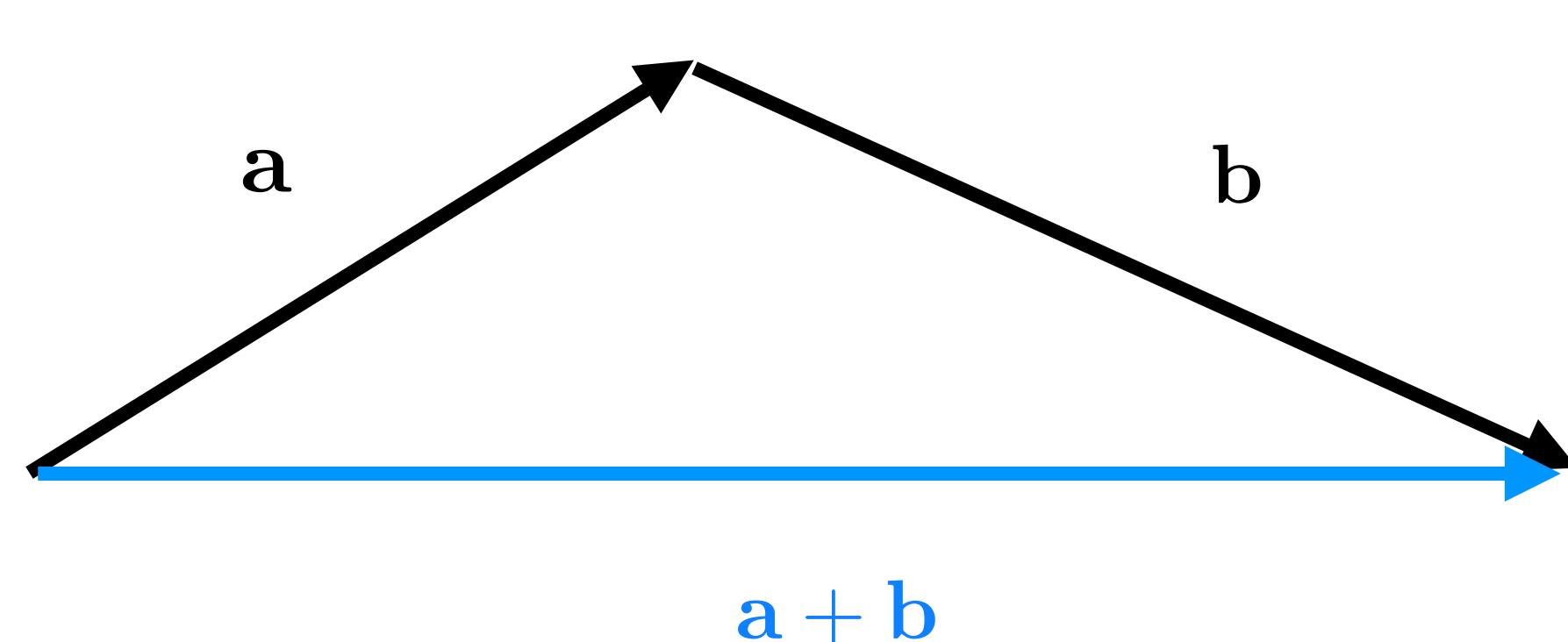


$$\mathbf{a} = (a_1, a_2, \dots, a_n)$$

$$\mathbf{b} = (b_1, b_2, \dots, b_n)$$

$$\mathbf{a} + \mathbf{b} = (a_1 + b_1, a_2 + b_2, \dots, a_n + b_n)$$

$$\mathbf{a} - \mathbf{b} = (a_1 - b_1, a_2 - b_2, \dots, a_n - b_n)$$



$$\mathbf{a} = (a_1, a_2, \dots, a_n)$$

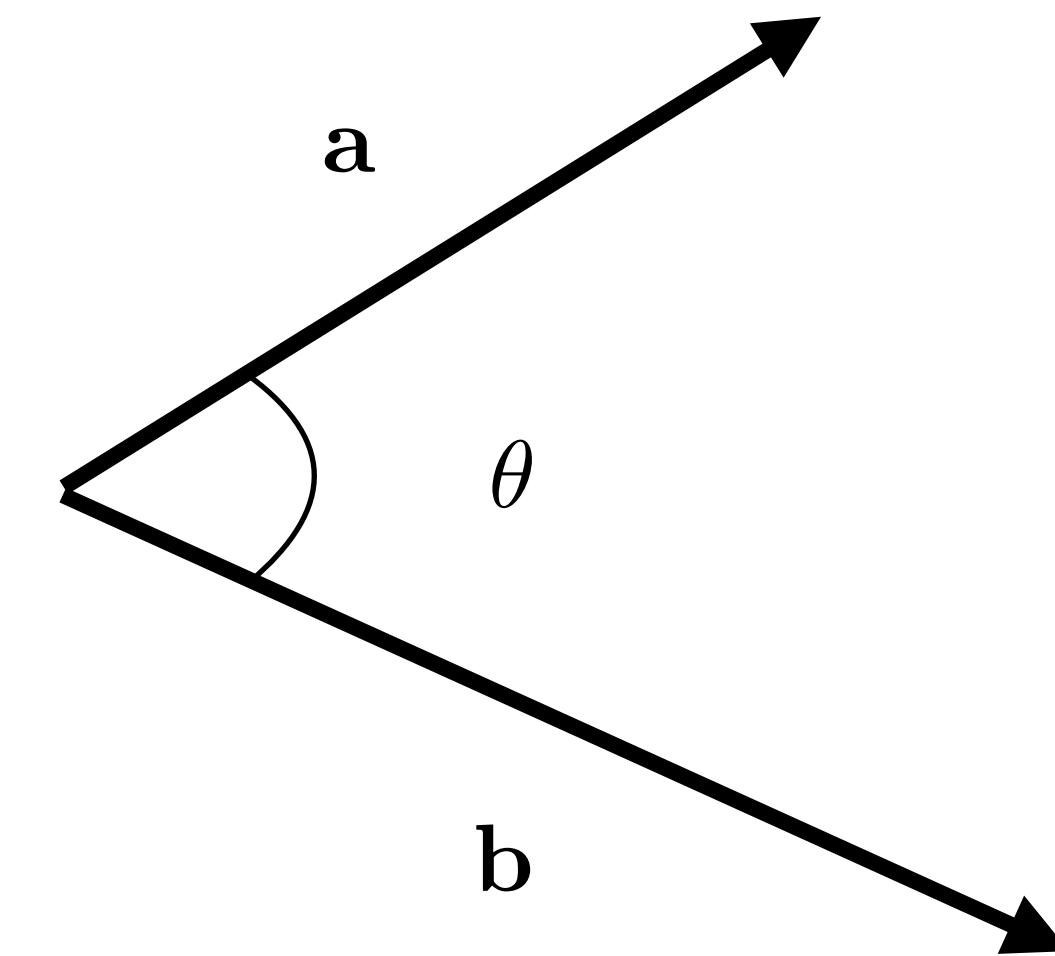
np.array

$$|\mathbf{b}/|\mathbf{b}||=1$$

$$\mathbf{b} = (b_1, b_2, \dots, b_n)$$

$$|\mathbf{a}| = (a_1^2 + \dots + a_n^2)^{1/2}$$

Length



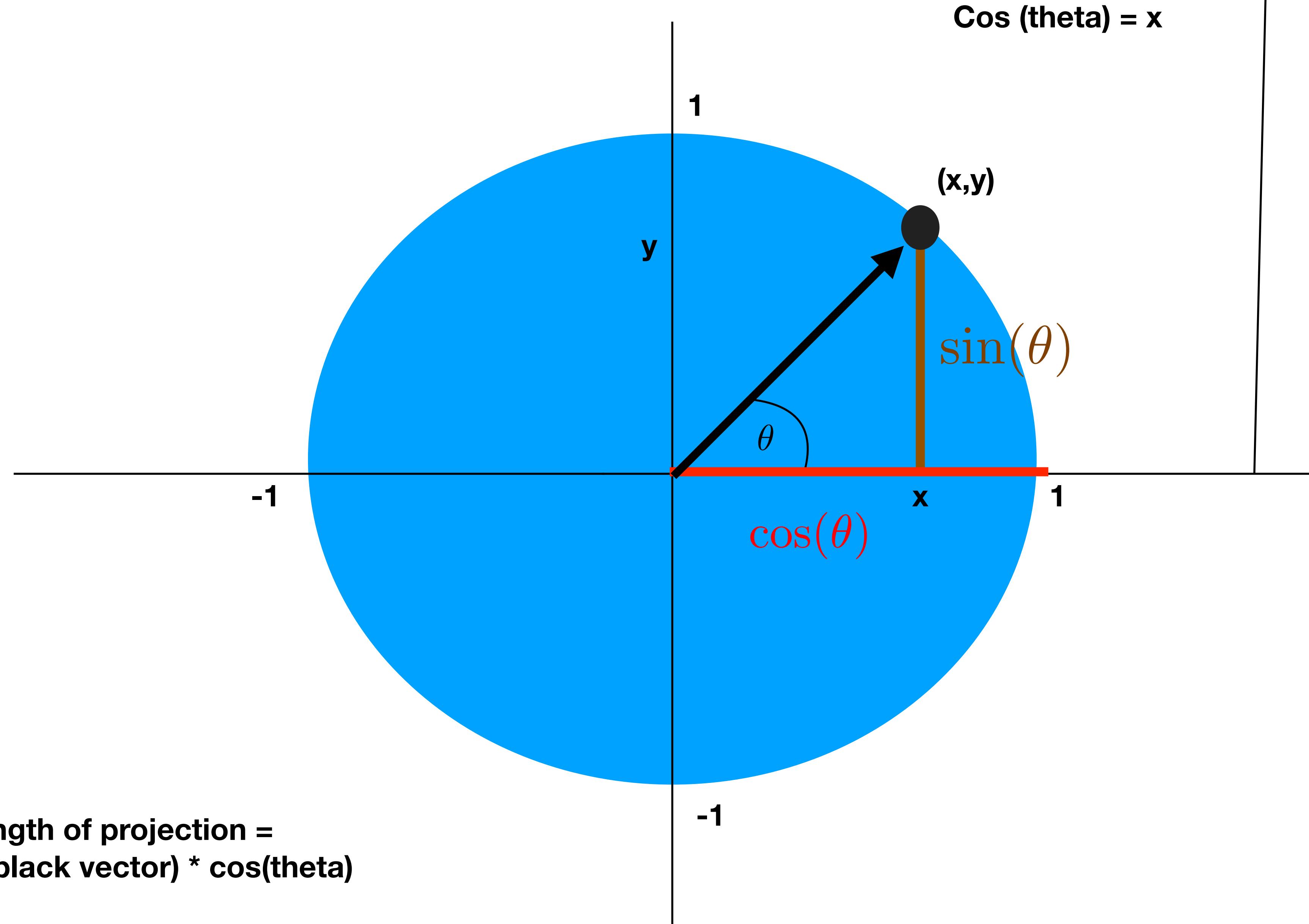
$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos(\theta) \quad (\text{Theorem})$$

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n = \sum_{i=1}^n a_i b_i \quad (1, 1), (1, -1)$$

$$\mathbf{a} \cdot \mathbf{b} = 0 \quad \text{means orthogonality because} \quad \cos(90^\circ) = 0$$

Exercise: write a python code that calculates the angle between two general vectors. Hint: use arccos, \cos^{-1} .

Sin and Cos functions



$$\mathbf{N} = (w_1, w_2, \dots, w_n)$$

We always normalize: take
the length of the normal
vector equal to 1!

$$|N| = 1$$

$$\mathbf{a} = (a_1, a_2, \dots, a_n)$$

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

$$(\mathbf{x} - \mathbf{a}) \cdot \mathbf{N} = 0$$

$$(x_1 - a_1)w_1 + \dots + (x_n - a_n)w_n = 0$$

$$w_1x_1 + \dots + w_nx_n + \textcolor{red}{w_0x_0} = 0$$

$$\textcolor{red}{x_0 = 1}$$

The equation of a hyper-plane in dimension n
with normal vector \mathbf{N} , with the vector (point \mathbf{a}) on the hyper-plane

$$\textcolor{red}{w_0x_0 = w_0 = -a_1w_1 - \dots - a_nw_n}$$

The equation of an (n-1)-dimensional plane (hyper-plane) in n-dimensional space is determined by a vector of weights:

$$\mathbf{w} = (w_0, w_1, \dots, w_n)$$

The points on the hyper-plane are all

$$\mathbf{x} = (x_1, \dots, x_n)$$

such that

$$w_1x_1 + \dots + w_nx_n + \textcolor{red}{w_0x_0} = 0$$

$$\textcolor{red}{x_0 = 1}$$

$$\mathbf{w} = (w_0, w_1, \dots, w_n)$$

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

$$\textcolor{red}{x_0 = 1}$$

$$\mathbf{w} \cdot \mathbf{x} = 0$$

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ \vdots \\ x_n \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ \vdots \\ w_n \end{bmatrix}$$

$$\mathbf{w}^T \mathbf{x} = [w_0 \quad w_1 \quad \dots \quad w_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} = 0$$

Homework 1: implement the following algorithm for linear perceptron with python and try it on data sets that are linearly separable. Also check the performance the non linearly separable data.

Training a linear perceptron: use the misclassified points from the training data to change the weights of the classifying plane to have less and less error.

Simple way (particular for linear perceptron):
initialize w randomly, label points

with $y=+1$ (correct) and $y=-1$ (misclassified)

pick a misclassified point and update w :

(z, y)

$$w \xrightarrow{\text{replace}} w + y z$$

meaning of misclassification

$$\text{sign}(w \cdot z) \neq y$$

Hint:

$$(x-z) \cdot N = + \text{ or } - \text{ the distance} = x \cdot N - z \cdot N = -w_0 - z \cdot N = w \cdot z$$

Exercise:

$$\text{error} = \text{distance of the misclassified point from the hyper-plane of } w = -y w \cdot z$$

$$\text{error} = \text{distance of the misclassified point from the hyper-plane of updated } w =$$

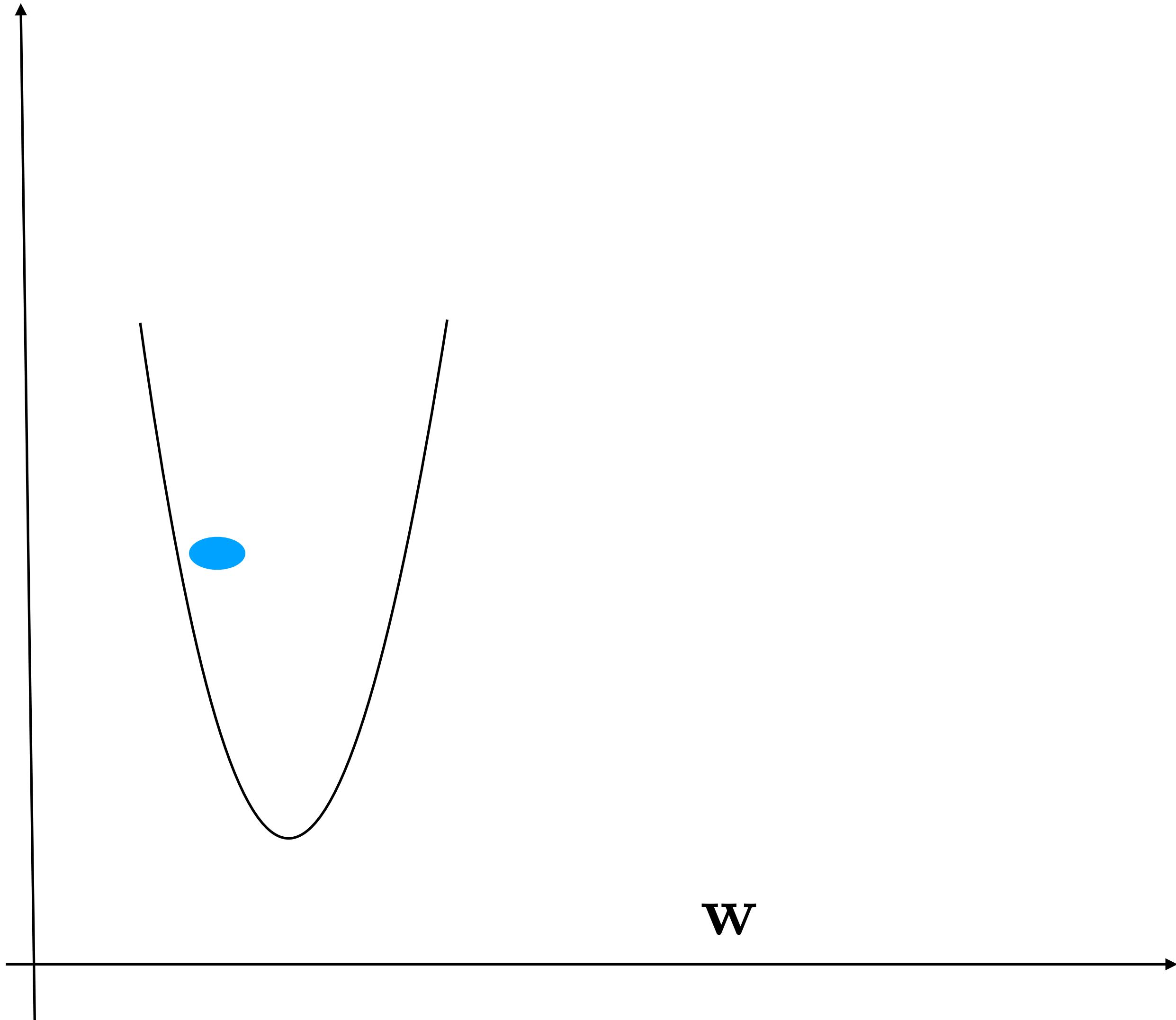
$$-y(w + yz) \cdot z = -y w \cdot z - y^2 z \cdot z$$

General way: using gradient

$$sign(s) = \begin{cases} +1, & s > 0 \\ -1, & s < 0 \end{cases}$$

Gradient descent (1-dimensional case)

Error = $f(\mathbf{w})$



Replace (update)

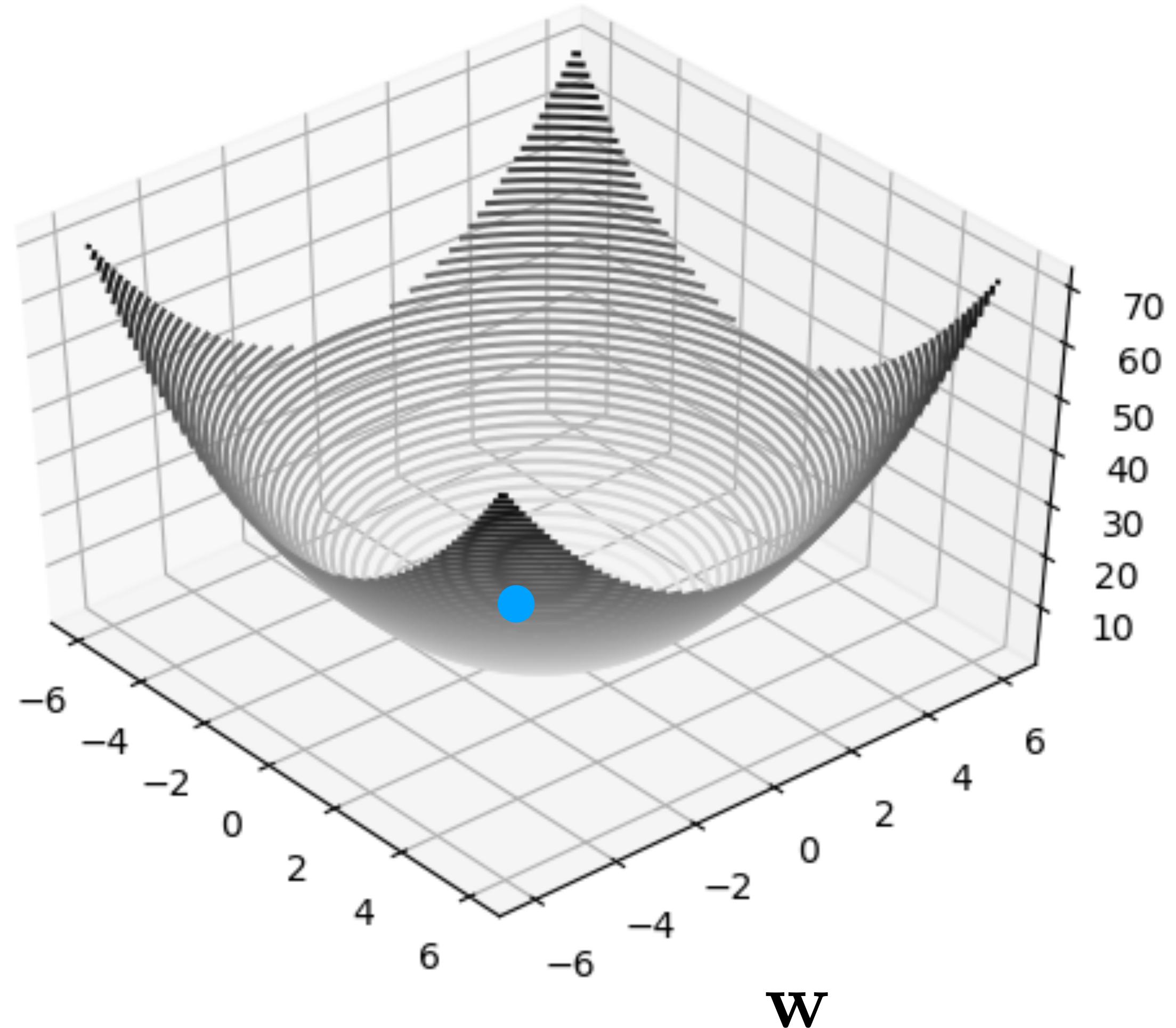
$$\mathbf{w} \longrightarrow \mathbf{w} - \eta f'(\mathbf{w})$$

η

learning rate
(positive)

**In which direction we have the maximum decrease
in the error? in the direction of minus gradient of error**

Error = $f(\mathbf{w}) = f(w_0, w_1, \dots, w_n)$



$$-\nabla f(\mathbf{w})$$

Gradient of f :

$$\nabla f(\mathbf{w}) = \left(\frac{\partial f}{\partial w_0}(\mathbf{w}), \frac{\partial f}{\partial w_1}(\mathbf{w}), \dots, \frac{\partial f}{\partial w_n}(\mathbf{w}) \right)$$

$$\mathbf{w} \longrightarrow \mathbf{w} - \eta \nabla f(\mathbf{w})$$

Example of calculating partial derivatives:

$$f(w) = w^k$$

$$f'(w) = \frac{df}{dw}(w) = kw^{k-1}$$

$$f(\mathbf{w}) = f(w_0, w_1, w_2) = 5w_0^2 w_1^3 w_2$$

$$\nabla f(\mathbf{w}) = \left(\frac{\partial f}{\partial w_0}, \frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2} \right)$$

$$= (10w_0 w_1^3 w_2, 15w_0^2 w_1^2 w_2, 5w_0^2 w_1^3)$$

The reason for the gradient descent working

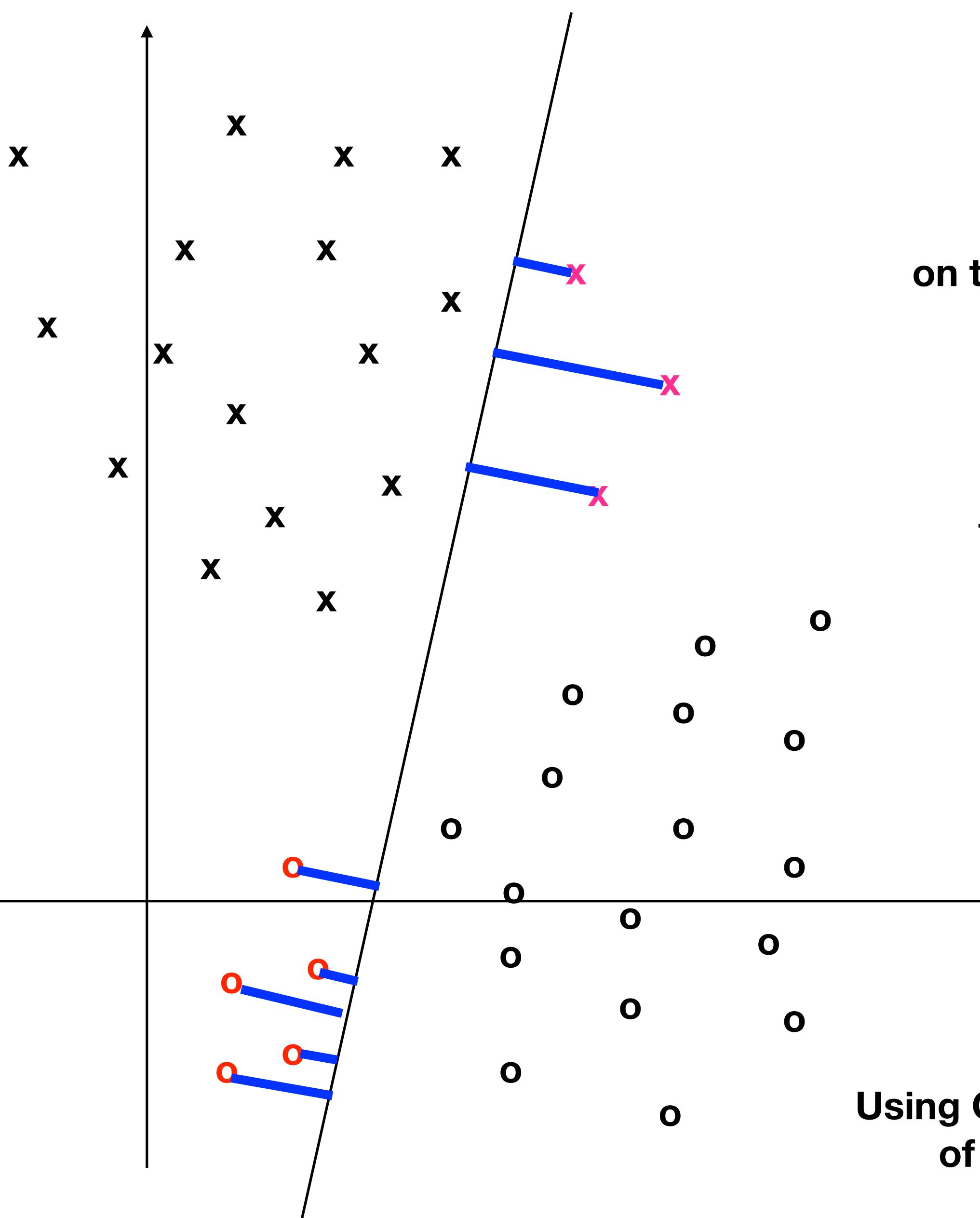
$$f(w + h) \sim f(w) + f'(w)h$$

$$f(\mathbf{w} + \mathbf{h}) \sim f(\mathbf{w}) + \nabla f(\mathbf{w}) \cdot \mathbf{h}$$

$$f(\mathbf{w} - \eta \nabla f(\mathbf{w})) \sim f(\mathbf{w}) - \eta \nabla f(\mathbf{w}) \cdot \nabla f(\mathbf{w})$$

$$f(w + h) = \sum_{k=0}^n \frac{f^k(w)}{k!} h^k$$

Suggested reading:
Taylor series
in one and several variables



Linear perceptron

on the hyper-plane

$$\mathbf{W} \cdot \mathbf{X} = 0$$

training data: $(\mathbf{z}_1, y_1), (\mathbf{z}_2, y_2), \dots, (\mathbf{z}_N, y_N)$

$$y_i = 1$$

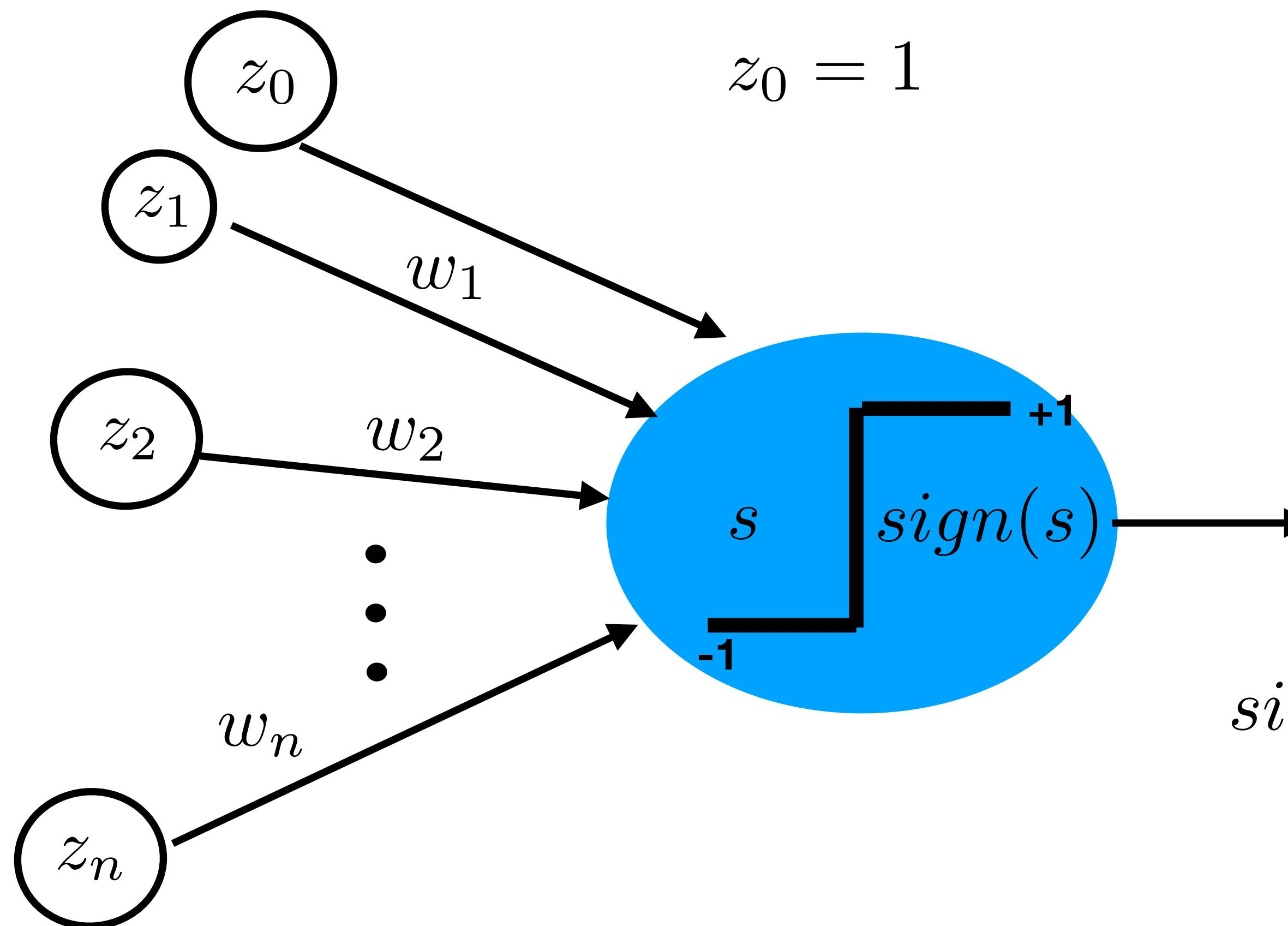
$$\text{or } y_i = -1$$

error = sum of the distances of the misclassified points
from the line.

Using Gradient descent we decrease this error (by changing the parameters
of the line and thereby the line) until there is no misclassified point.

$$\mathbf{z} = (z_1, z_2, \dots, z_n)$$

Perceptron (linear classifier)

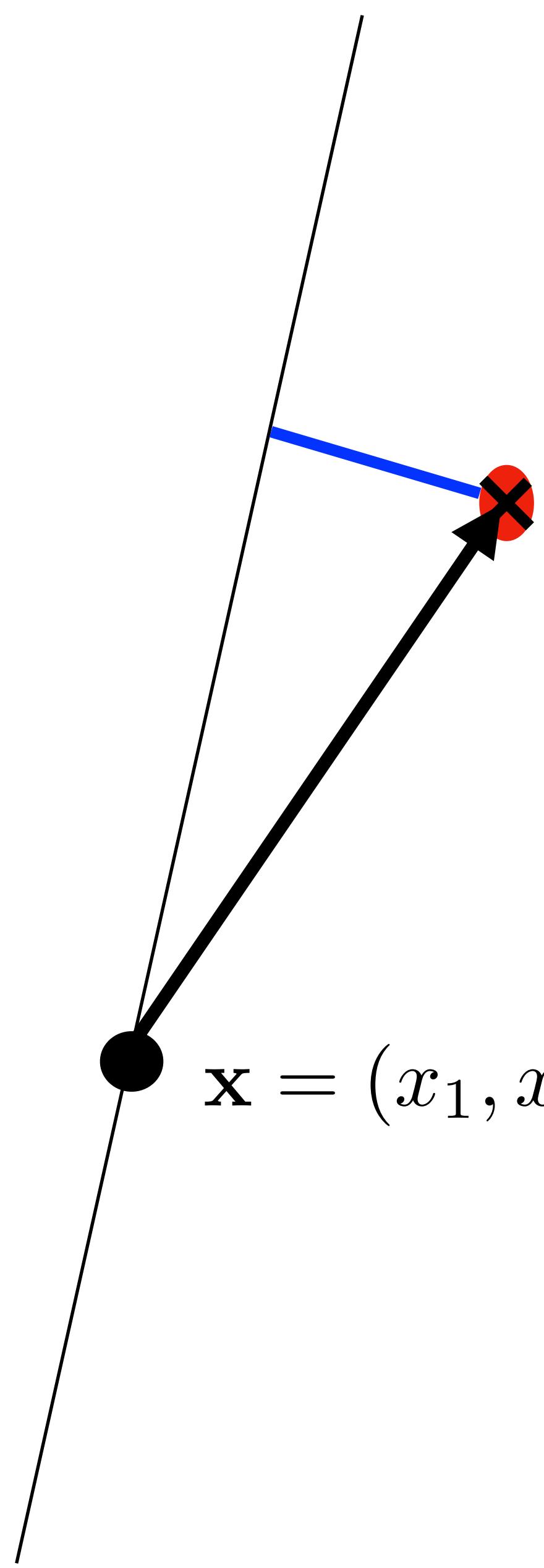


$$h(z) = sign(\mathbf{w} \cdot \mathbf{z})$$

$$s = \sum_{i=1}^n w_i z_i + w_0$$

$$sign(s) = \begin{cases} +1, & s > 0 \\ -1, & s < 0 \end{cases}$$

Eventually, this classifier should give the same sign for all x's and the same sign for all o's.



On the hyper-plane:

$$\sum_{i=1}^n w_i x_i + w_0 = 0$$

$$\mathbf{N} = (w_1, w_2, \dots, w_n)$$

the distance = $|(\mathbf{z} - \mathbf{x}) \cdot \mathbf{N}|$ = $-y \left(\sum_{i=1}^n z_i w_i + w_0 \right)$

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

Stochastic gradient descent

Pick a misclassified point randomly

Calculate the error

Update the weights using gradient descent to decrease the error (do not forget to normalize the weights)

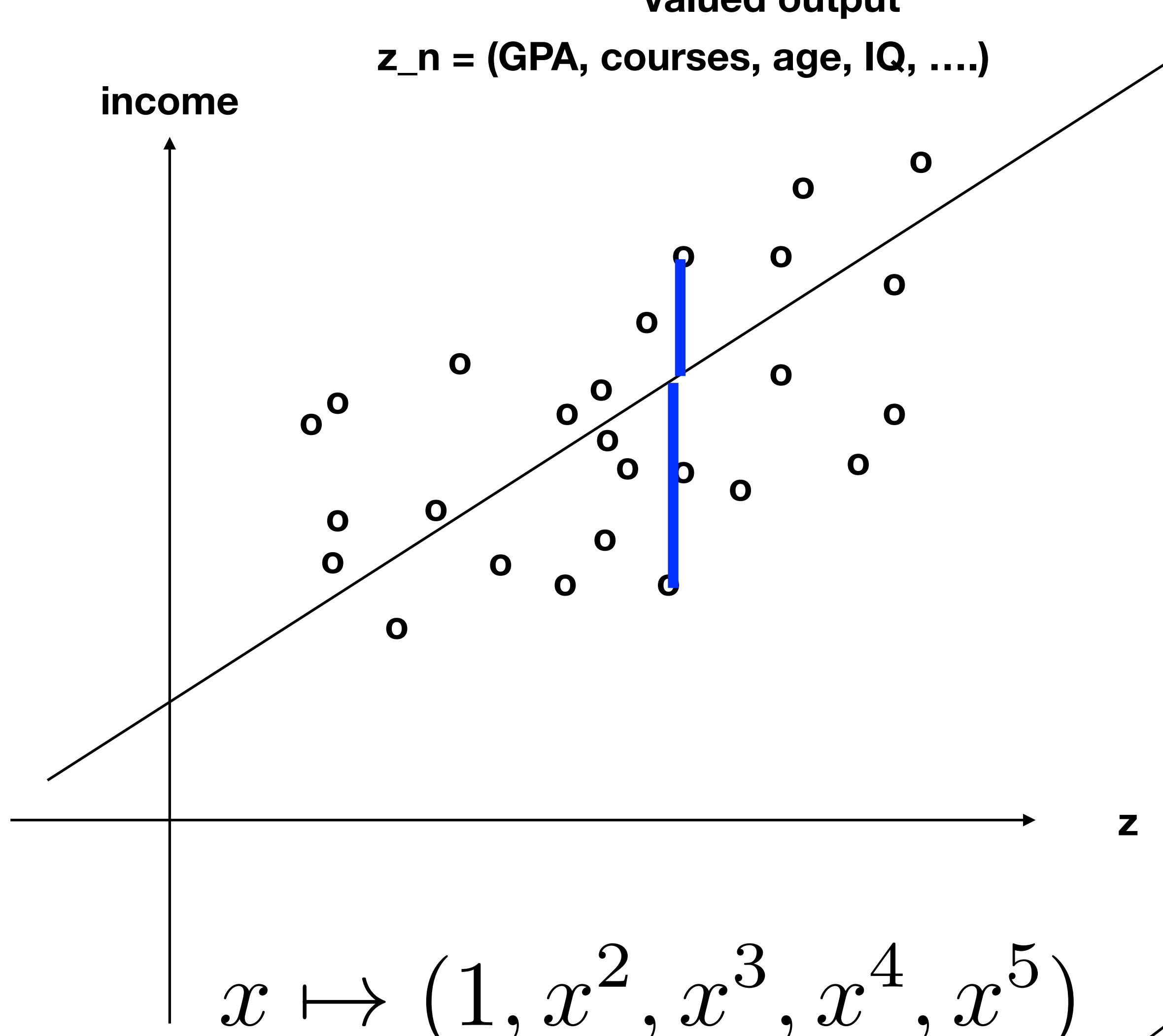
Continue this process until there is no misclassified point in the training data set

Homework: Write a python code to implement the above algorithm for the linear perceptron and show the performance on a linearly separable data set

**Linear regression: rather than classification,
the goal is to fit the data and produce real-
valued output**

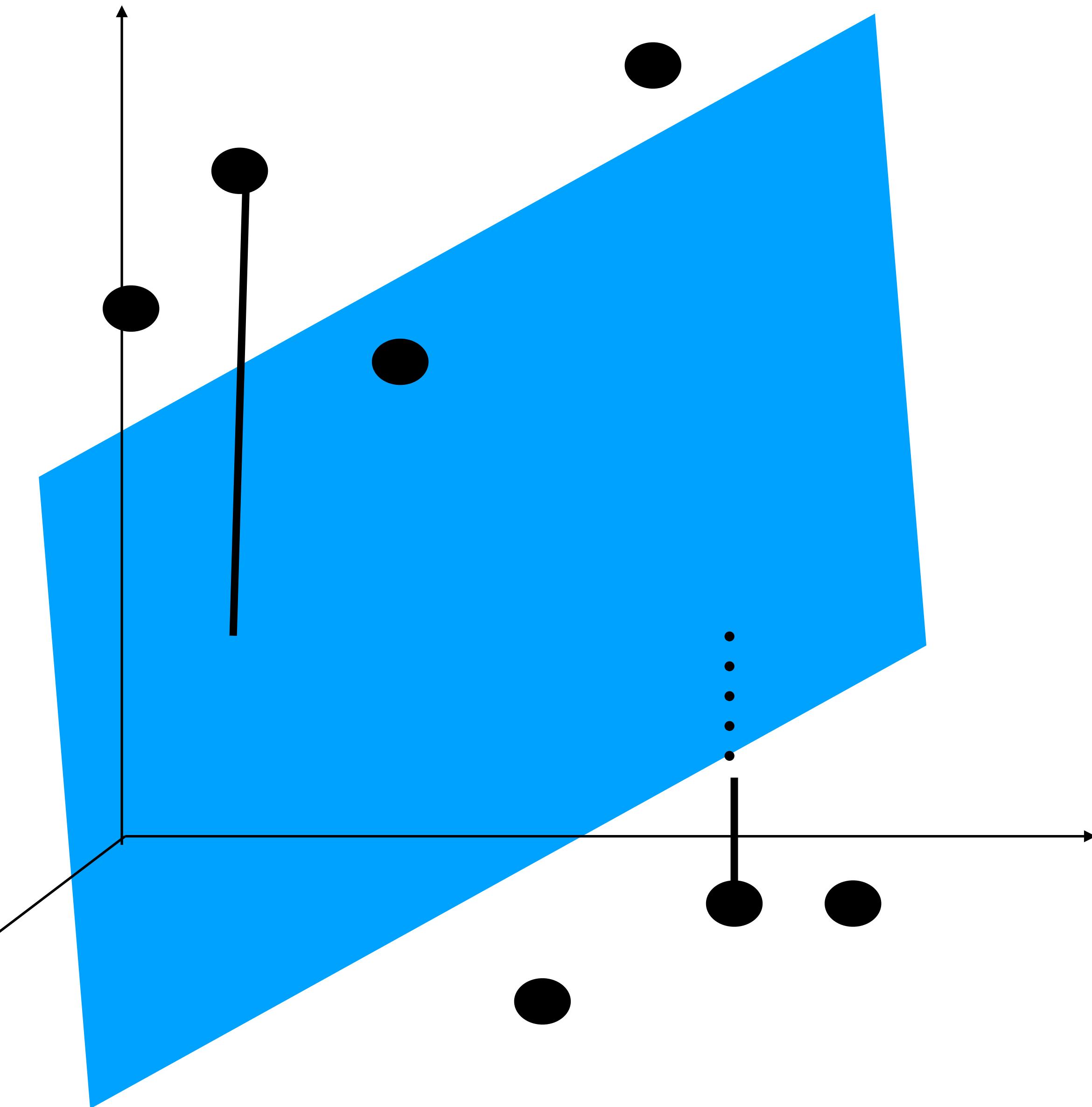
$z_n = (\text{GPA, courses, age, IQ, ...})$

income



$$x \mapsto (1, x^2, x^3, x^4, x^5)$$

$$(x_1, x_2) \mapsto (1, x_1, x_1^2, x_2, x_2^2, x_1 x_2)$$



Application of linear regression: prediction

For example to predict the future income of a student based on the grades and courses

Historical data: $(\mathbf{z}_1, y_1), (\mathbf{z}_2, y_2), \dots, (\mathbf{z}_N, y_N)$

$\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ **vectors in** $\{1\} \times \mathbb{R}^n$

y_1, y_2, \dots, y_N **real numbers**

Goal: Find the weights $\mathbf{w} = (w_0, w_1, \dots, w_n)$

such that the hyper-plane defined by

$$\mathbf{w} \cdot \mathbf{x} = 0 \iff w_1x_1 + \cdots + w_nx_n + \color{red}{w_0x_0} = 0$$
$$\mathbf{x} = (x_0, x_1, \dots, x_n)$$
$$\color{red}{x_0 = 1}$$

has the minimum sum of squares of distances for the historical data

$$(\mathbf{z}_1, y_1), (\mathbf{z}_2, y_2), \dots, (\mathbf{z}_N, y_N)$$

Sum of squares of distances

$$d = \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{z}_i - y_i)^2$$

Note: remember that 1 is the first coordinate of each of $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$

The vector of weights can be found by solving:

$$\nabla d(\mathbf{w}) = 0$$

This is a particular situation!

Goal: Find \mathbf{w} in such a way that d is the minimum possible!

our predictor will look like this:

$$h(z) = \mathbf{w} \cdot \mathbf{z} = \sum_{i=0}^n w_i z_i$$

Homework: write the linear regression algorithm in python and use it for prediction

$$\mathbf{w} = \mathbf{Z}^\dagger \mathbf{y}$$

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

(n+1)*N

The minimizing vector of weights

shape of the pseudo inverse = (n+1)*N

shape of w = (n+1)*1

Pseudo-inverse of Z

$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \vdots \\ \mathbf{z}_N \end{bmatrix}_{N \times (n+1)}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

shape of the $\mathbf{Z}^T \mathbf{Z} = (n+1)*(n+1)$

Matrix multiplication

$$C = AB$$

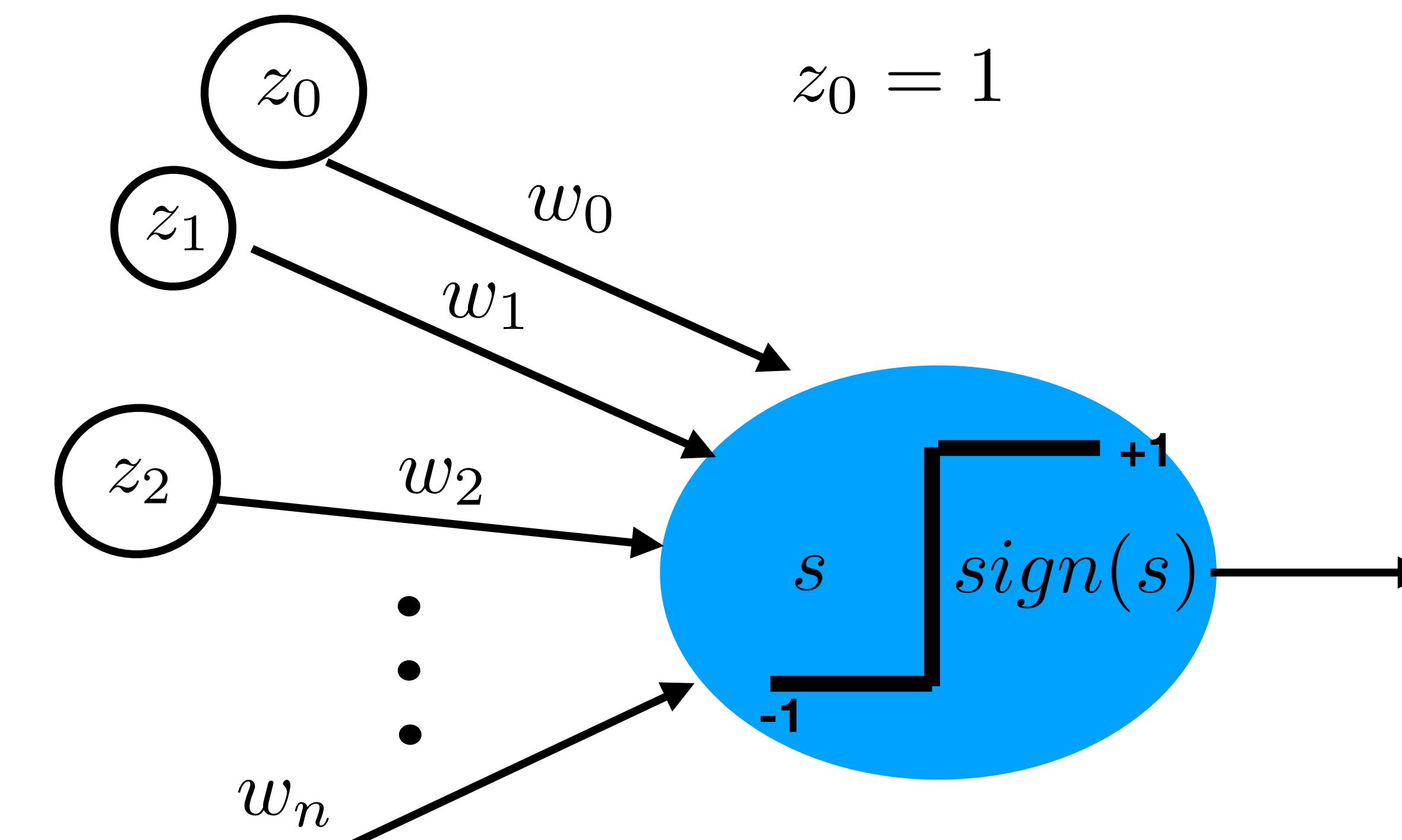
Suggested reading:
matrices from quantum mechanics
(Heisenberg)

$$C_{m \times n} = A_{m \times k} B_{k \times n}$$

C_{ij} = the element in row i and column j of the matrix C
the dot product of of i-th row A with the j-th
column of B

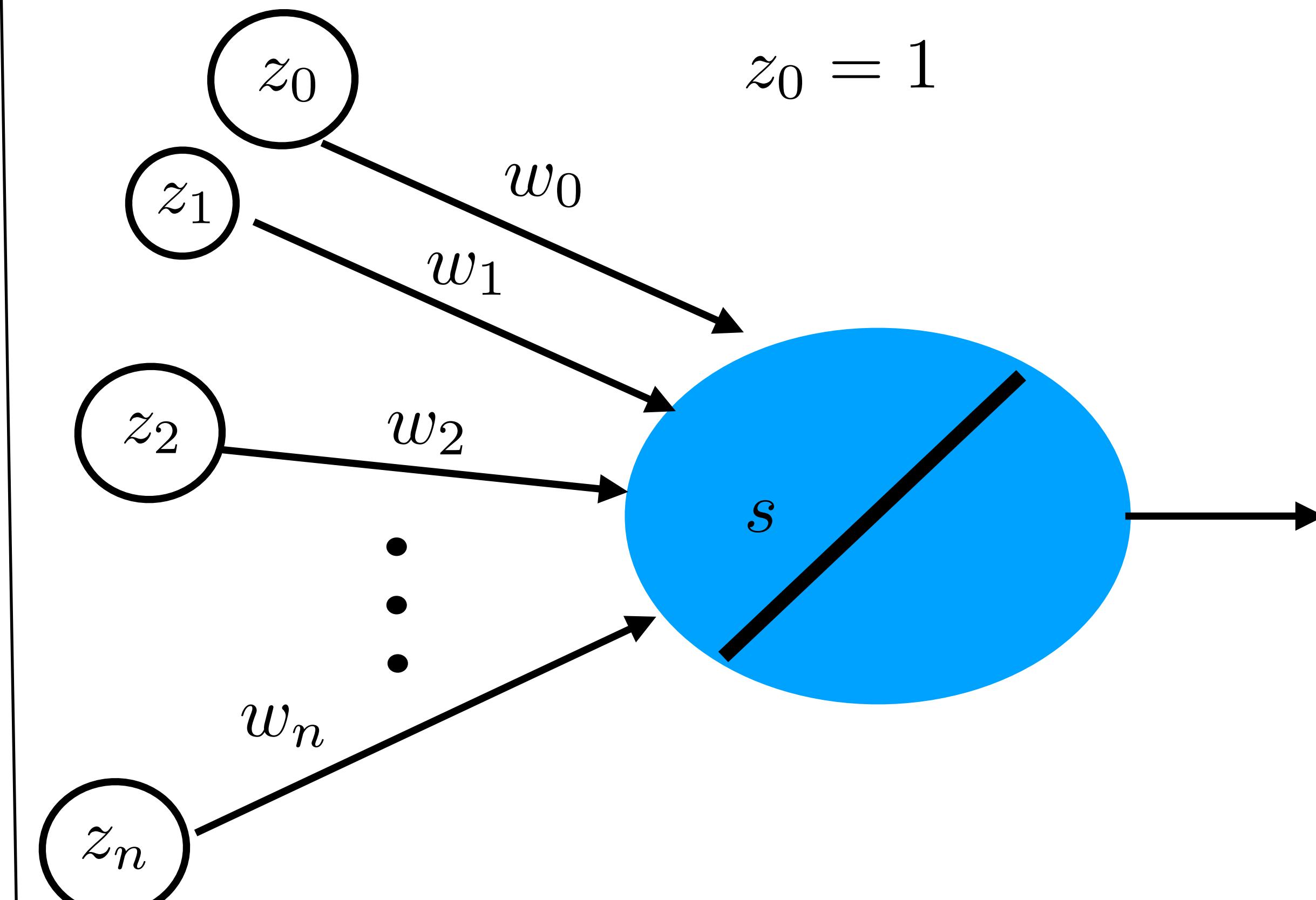
$$\begin{pmatrix} 1 & -1 \\ 4 & 2 \\ 1 & 7 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 & -1 \\ 5 & 6 & 1 & 2 \end{pmatrix} = \begin{pmatrix} -4 & -4 & 2 & -3 \\ 14 & 20 & 14 & 0 \\ 36 & 44 & 10 & 13 \end{pmatrix}$$

linear perception: classifier



$$h(z) = \text{sign}(\mathbf{w} \cdot \mathbf{z})$$

linear regression: prediction: output
is a real number



$$h(z) = \mathbf{w} \cdot \mathbf{z}$$

$$p = \text{Probability}(X | H) < 0.005$$

Statistical inference

X = observation

H = an assumption about X

infer: the assumption is wrong!!!!!!!

Maximum Likelihood

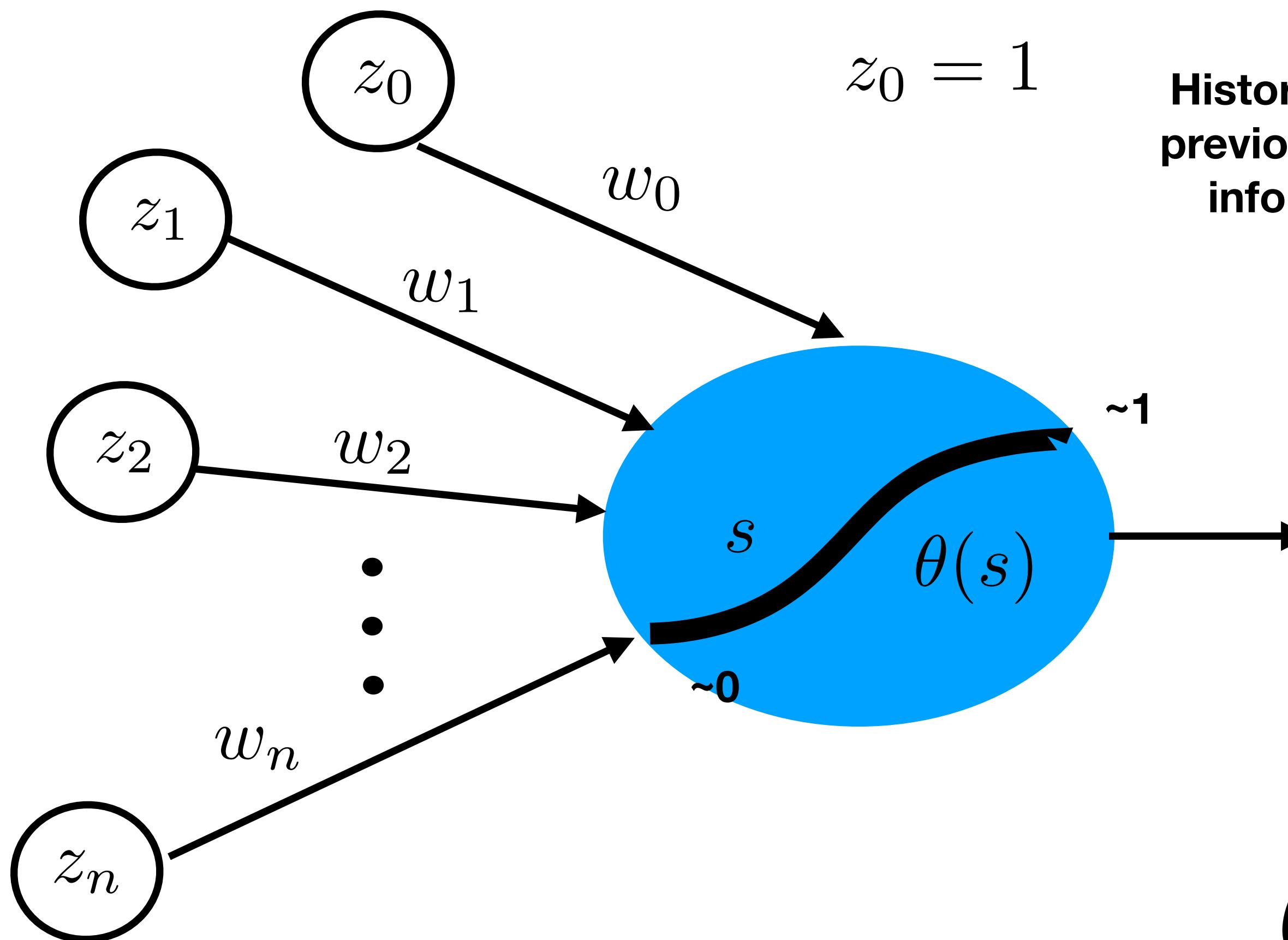
$$p = \text{prob}(X | f_{\Theta})$$

Θ

is unknown, and
we wish to identify it

This is done by maximizing p: finding the Theta that maximizes the probability p of observing X

Logistic regression: the goal is to give a probability



Historical data:
previous patient
information

$$\mathbf{z} = (z_1, z_2, \dots, z_n)$$

$$y = +1 \text{ or } -1$$

Goal: $P(y|\tilde{\mathbf{z}}) = \theta(\mathbf{w} \cdot \tilde{\mathbf{z}})$

\mathbf{w} will be identified
using the historical or training data and
and maximum likelihood

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

The logistic function or the sigmoid function

Historical data: $(\mathbf{z}_1, y_1), (\mathbf{z}_1, y_1), \dots, (\mathbf{z}_N, y_N)$

$\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ vectors in $\{1\} \times \mathbb{R}^n$

y_1, y_2, \dots, y_N are +1 or -1

Exercise:

$$\theta(-s) = 1 - \theta(s)$$

$$P(y|\mathbf{z}) = \begin{cases} \theta(\mathbf{w} \cdot \mathbf{z}), & y = +1 \\ 1 - \theta(\mathbf{w} \cdot \mathbf{z}), & y = -1 \end{cases}$$
$$= \theta(y \mathbf{w} \cdot \mathbf{z})$$

Likelihood of the historical data:

$$\prod_{i=1}^N P(y_i|\mathbf{z}_i) = P(y_1|\mathbf{z}_1)P(y_2|\mathbf{z}_2) \cdots P(y_N|\mathbf{z}_N)$$

$$= \prod_{i=1}^N \theta(y_i \mathbf{w} \cdot \mathbf{z}_i)$$

Goal: chose the weights (w) such that this likelihood is maximized.

$$\ell = \prod_{i=1}^N \theta(y_i \mathbf{w} \cdot \mathbf{z}_i)$$

Maximize is equivalent to

$\ln(\ell)$ **maximizing this, which is equivalent to**

$-\ln(\ell)$ **minimizing this**

$$-\frac{1}{N} \ln(\ell) = \frac{1}{N} \sum_{i=1}^N \ln \left(\frac{1}{\theta(y_i \mathbf{w} \cdot \mathbf{z}_i)} \right) = \frac{1}{N} \sum_{i=1}^N \ln (1 + e^{-y_i \mathbf{w} \cdot \mathbf{z}_i})$$

“Error of logistic regression” $= \frac{1}{N} \sum_{i=1}^N \ln (1 + e^{-y_i \mathbf{w} \cdot \mathbf{z}_i})$

risk score;
average of cross entropy errors

Now we can try to minimize the error!

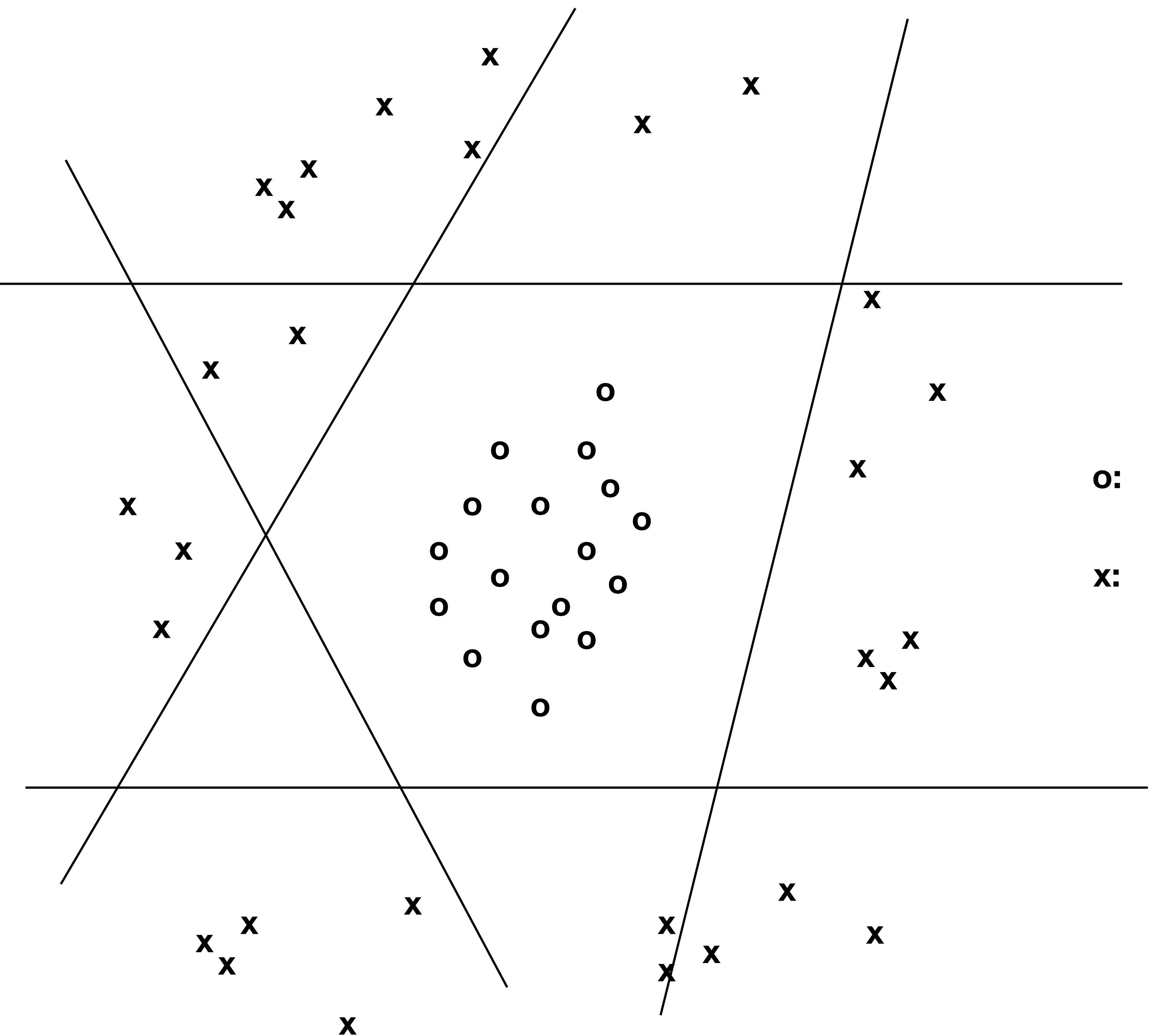
Exercise:

$$\nabla \text{“Error of logistic regression”} = \nabla E_L$$

$$= -\frac{1}{N} \sum_{i=1}^N \frac{y_i \mathbf{z}_i}{1+e^{y_i \mathbf{w} \cdot \mathbf{z}_i}}$$

Linear regression algorithm

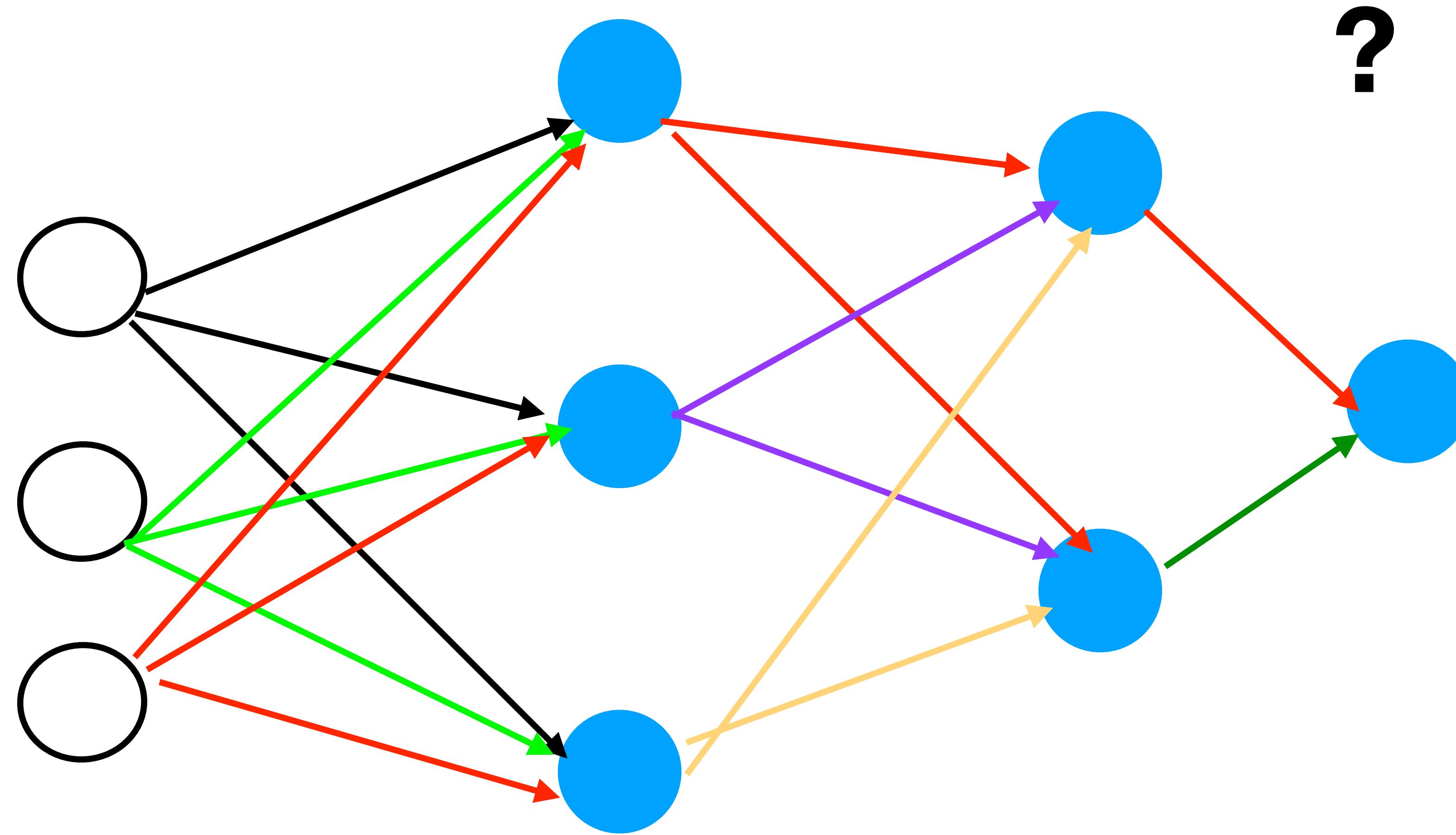
- 1. Choose the weights of the vector \mathbf{w} randomly**
- 2. compute the “error”** $= \frac{1}{N} \sum_{i=1}^N \ln (1 + e^{-y_i \mathbf{w} \cdot \mathbf{z}_i})$
- 3. Update the weights using the formula** $\mathbf{w} \rightarrow \mathbf{w} - \eta \nabla \text{“error”}$
- 4. Iterate the same method until it is time to stop.**

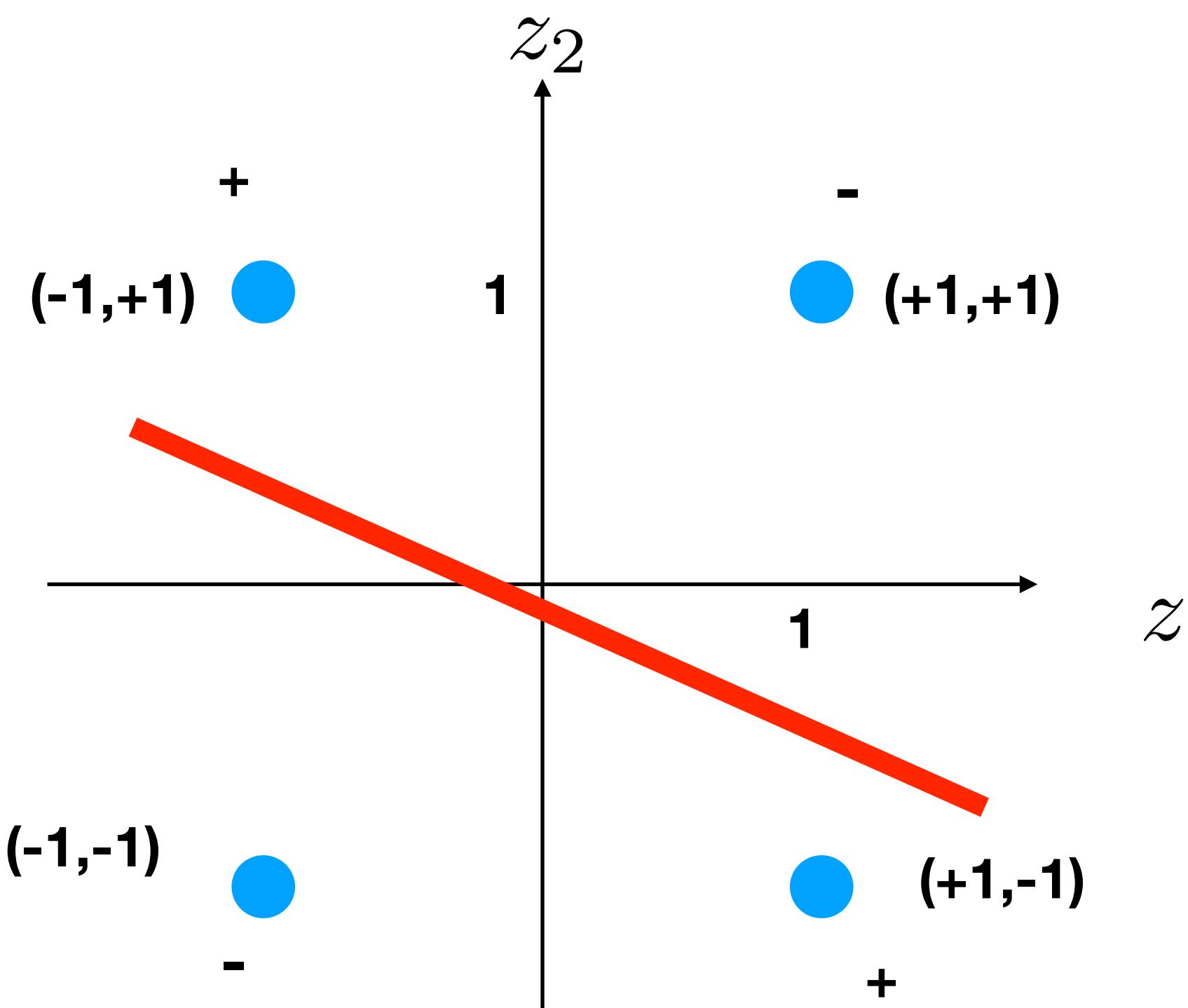


$$x_1^2 + x_2^2 = r^2$$

$$\text{o: } x_1^2 + x_2^2 < r^2$$

$$\text{x: } x_1^2 + x_2^2 > r^2$$





XOR function is defined on 4 points: $(-1, -1), (+1, -1), (+1, +1), (-1, +1)$

$$\mathbf{z} = (z_0 = 1, z_1, z_2)$$

Question: is the XOR function a linear perceptron?

Can we write:

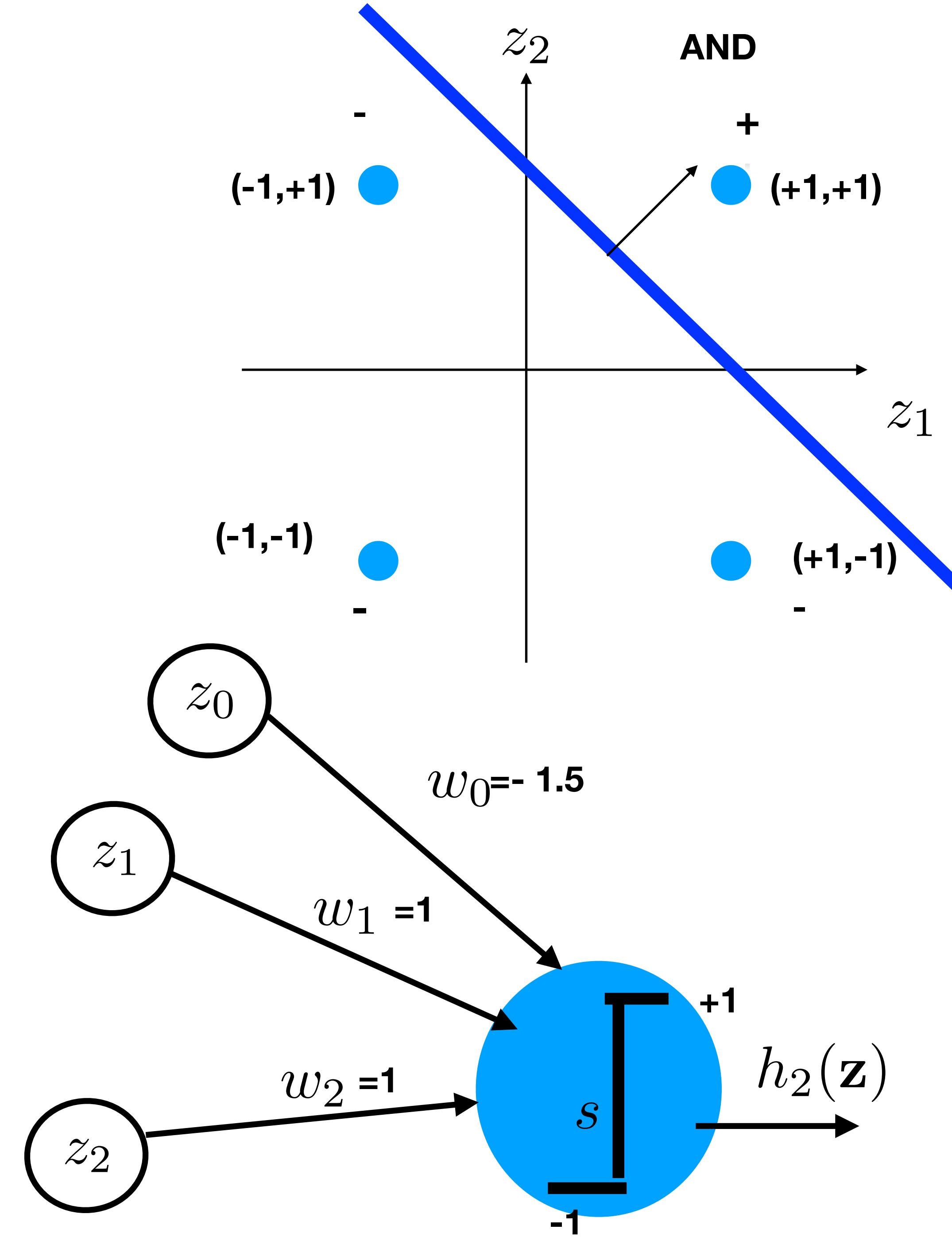
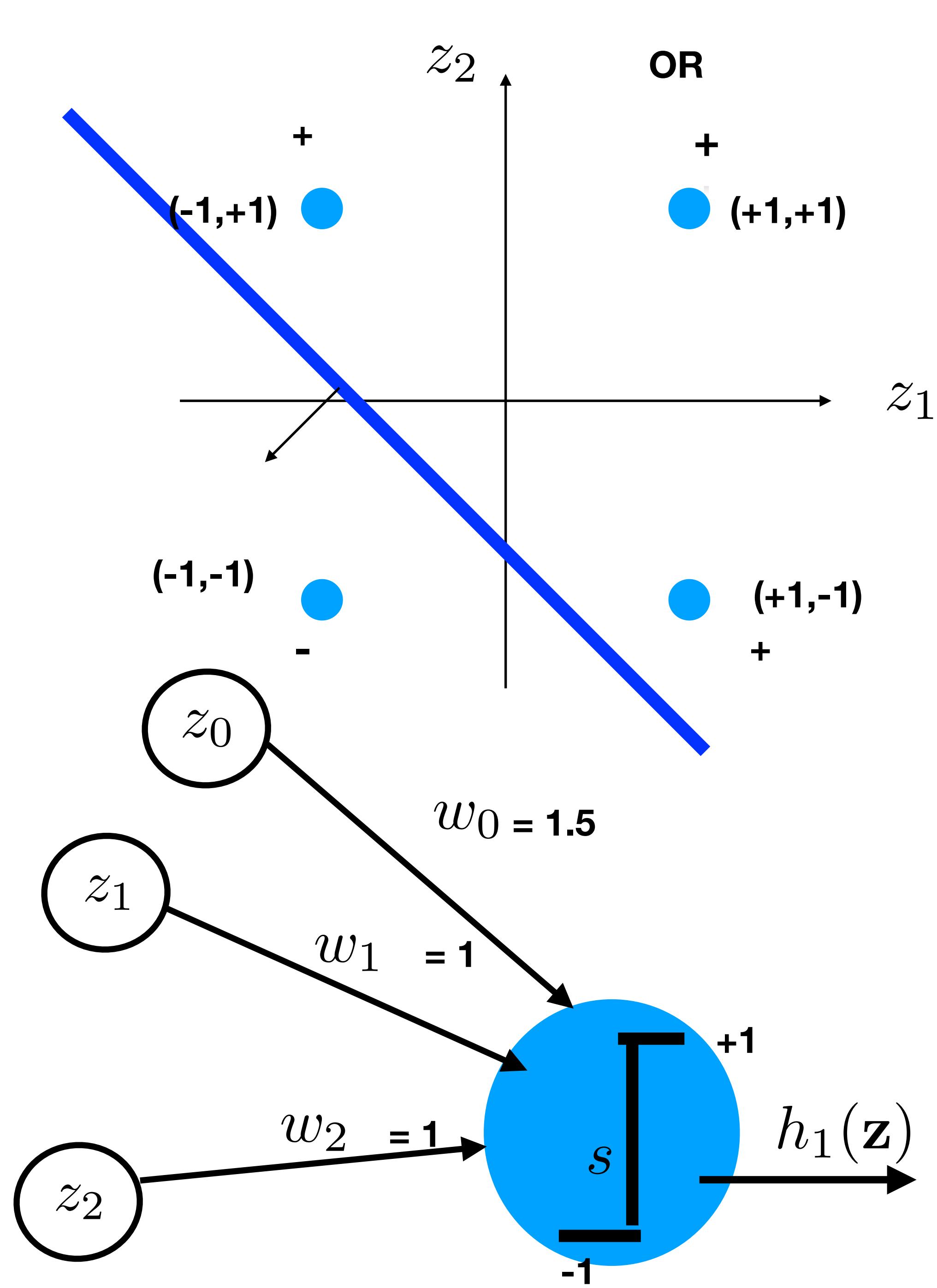
$$XOR(\mathbf{z}) = \text{sing}(\mathbf{w} \cdot \mathbf{z}) ?$$

Answer: NO because there is no straight line classifying the +'s and -'s in our graph.

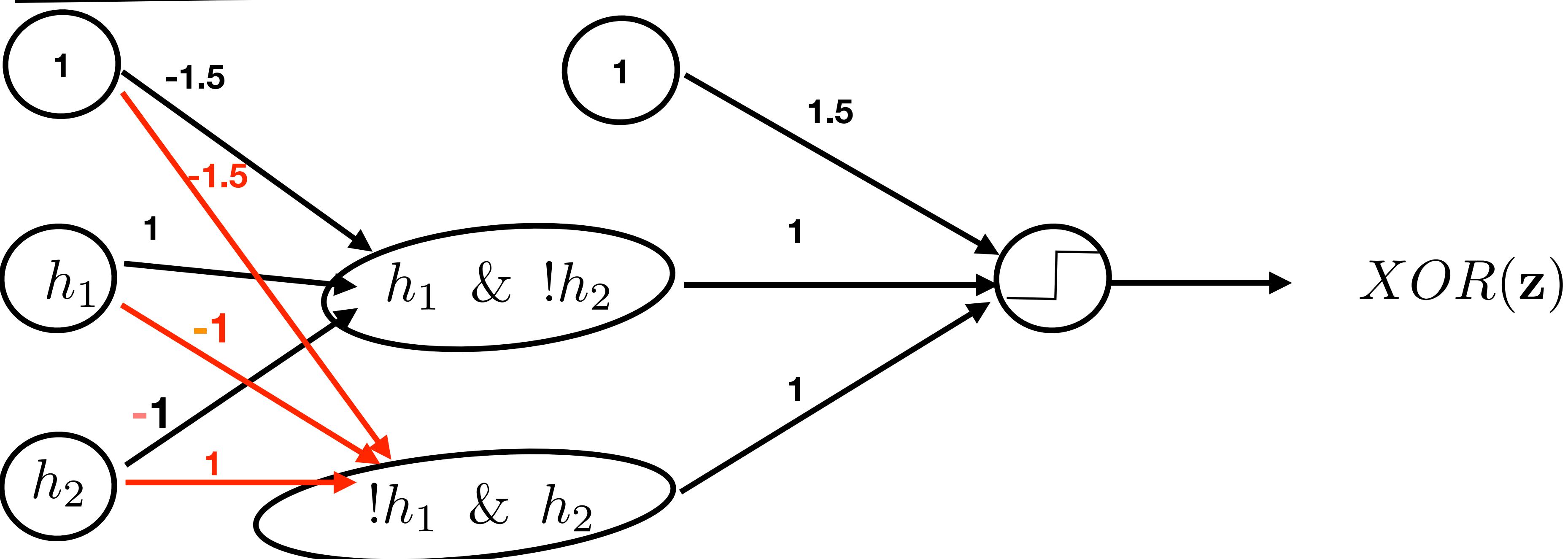
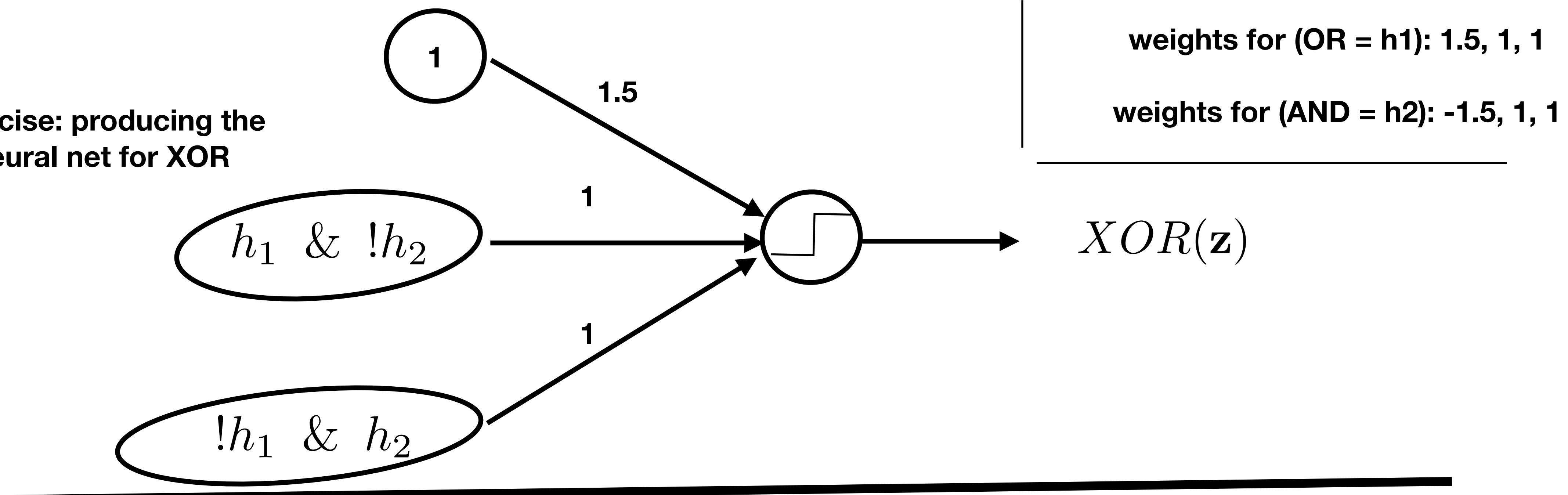
$$XOR(\mathbf{z}) = \begin{cases} +1, & \text{if exactly one coordinate is 1 } (z_1 \text{ or } z_2) \\ -1 & \text{otherwise} \end{cases}$$

Question: Can we use a combination of linear perceptrons to construct the XOR function?

Yes, because the AND and OR functions can be written as linear perceptrons.



Exercise: producing the
neural net for XOR



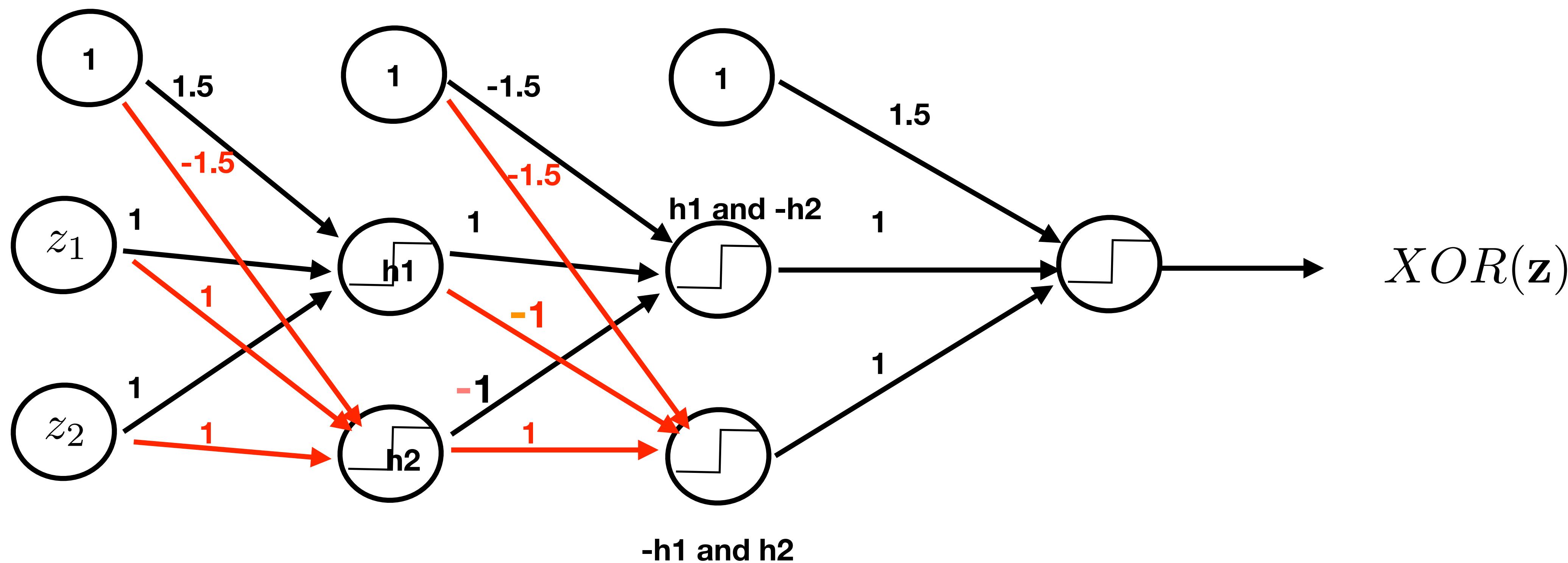
XOR

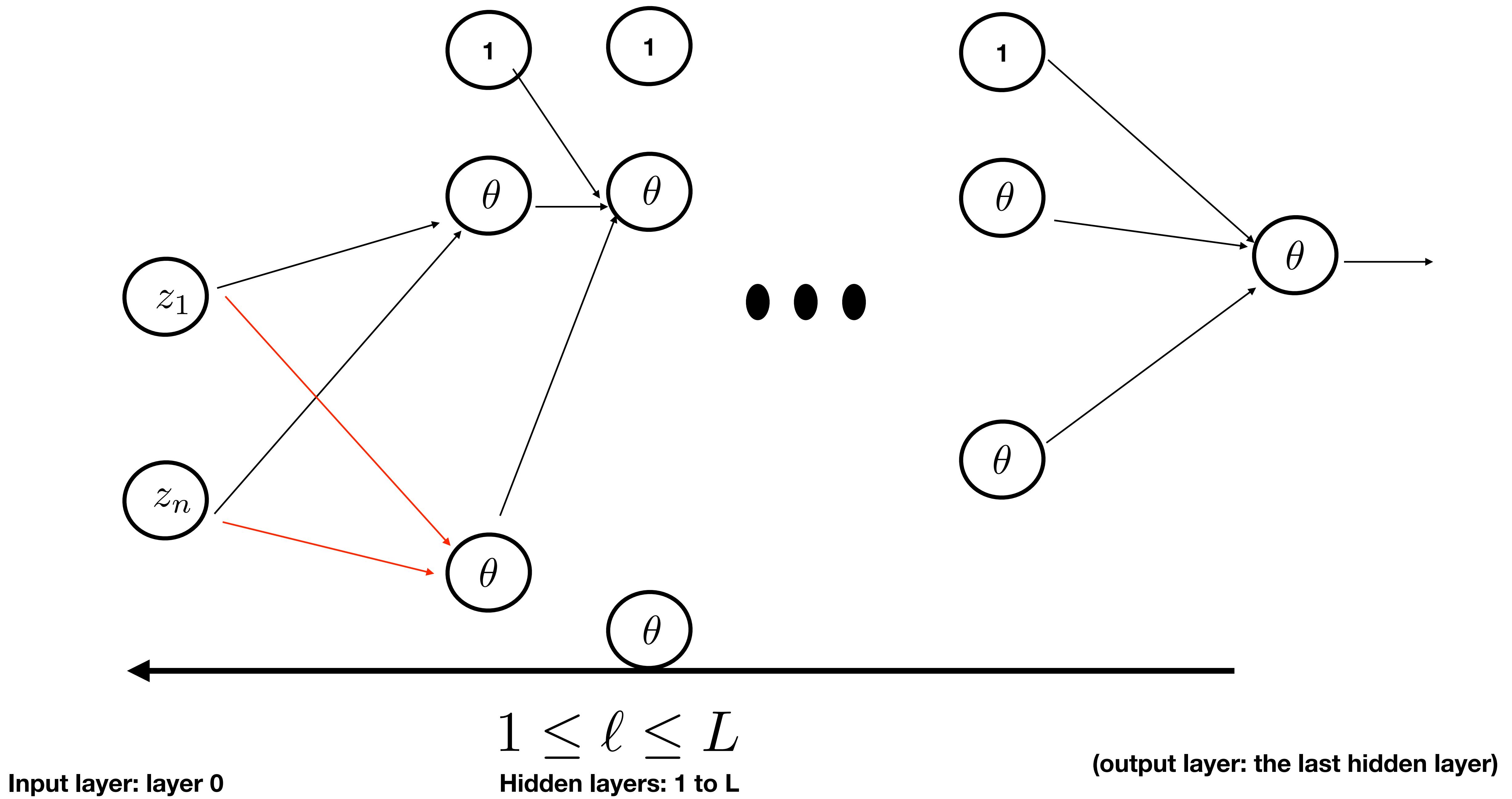
weights for (OR = h1): 1.5, 1, 1

weights for (AND = h2): -1.5, 1, 1

((2=3) or (3=3)) and ((3=4) and 3=3) False = -1

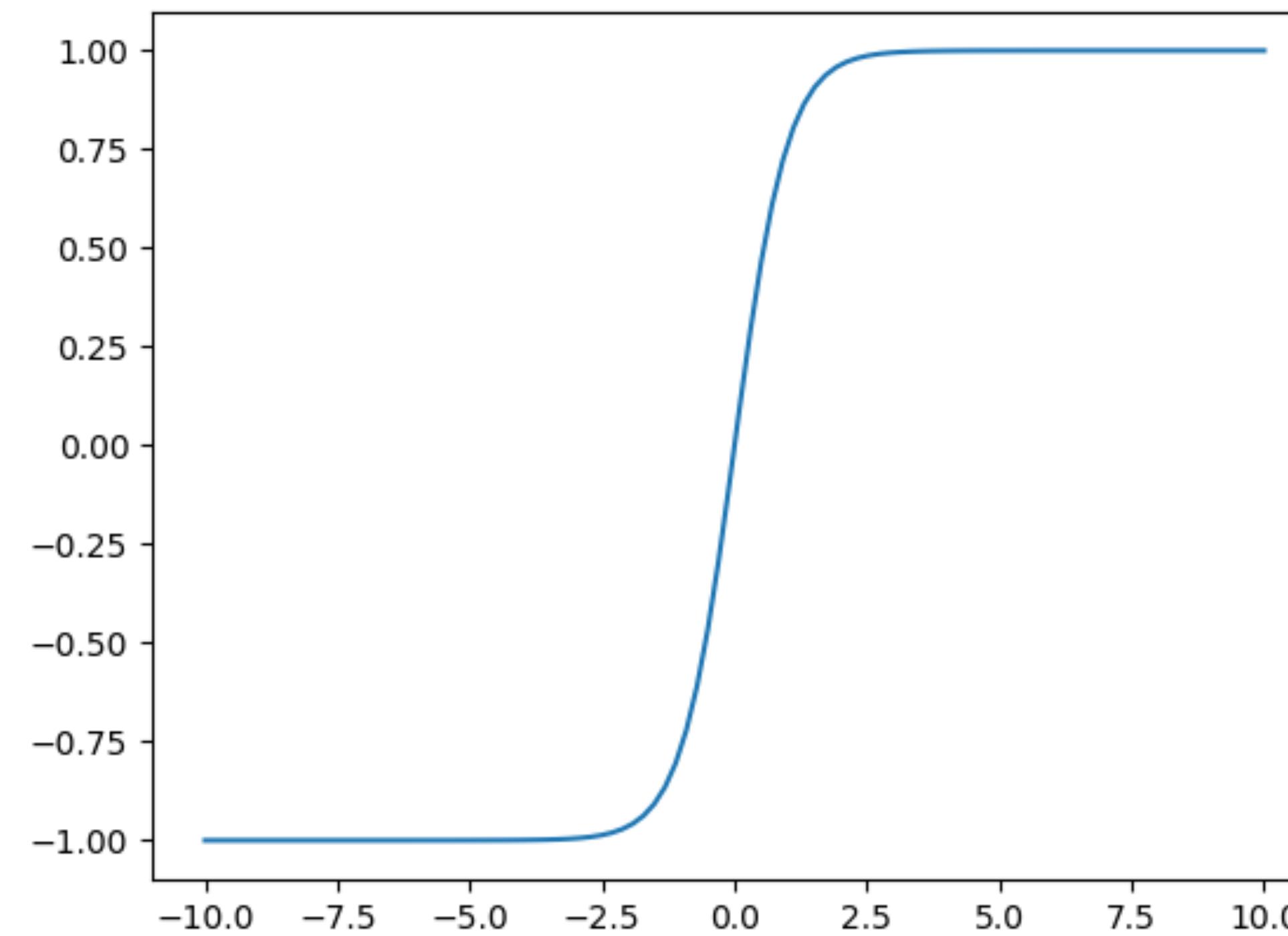
((2=3) and (3=3)) or ((3=4) or 3=3) True= +1





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

$$\theta'(s) = 1 - \theta(s)^2$$



Training a neural network

Training data: $(\mathbf{z}_1, y_1), (\mathbf{z}_1, y_1), \dots, (\mathbf{z}_N, y_N)$

i input of the layer
j output from the layer

weights: $\mathbf{w} = \left(w_{ij}^{(\ell)} \right) \quad 1 \leq \ell \leq L \quad 0 \leq i \leq d^{(\ell-1)} \quad 1 \leq j \leq d^{(\ell)}$

$z_j^{(\ell)} = \theta(s_j^{(\ell)}) = \theta \left(\sum_{i=0}^{d^{(\ell-1)}} w_{ij}^{(\ell)} z_i^{(\ell-1)} \right)$

Apply this to a training data to get an output $h(\mathbf{z}) = z_1^{(L)}$

measure the error $e(\mathbf{w}) = e(h(\mathbf{z}), y)$ **for example** $|h(\mathbf{z}) - y|^2$

use gradient descent to decrease the error by updating the weights

Back propagation: a fast algorithm for implementing the gradient descent to neural networks

$$\nabla e(\mathbf{w})$$

$$\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(\ell)}} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(\ell)}} \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(\ell)}} z_i^{(\ell-1)} = \delta_j^{(\ell)} z_i^{(\ell-1)}$$

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

In the output layer: $\ell = L, j = 1$

For example: $e(\mathbf{w}) = (z_1^{(L)} - y)^2$

$$z_1^{(L)} = \theta(s_1^{(L)})$$

Back propagate

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

$$\delta_i^{(\ell-1)} = (1 - (z_i^{(\ell-1)})^2) \sum_{j=1}^{d^{(\ell)}} w_{ij}^{(\ell)} \delta_j^{(\ell)}$$

Back propagation algorithm for deep forward neural networks

1.

Initialize the weights $w_{ij}^{(\ell)}$ at RANDOM

3. Go forward in the NN by calculating all

$$z_j^{(\ell)}$$

2. Pick a training point (\mathbf{z}, y)

until you get the output for the chosen training point

4. Use back propagation to calculate all

$$\delta_j^{(\ell)}$$

5. Update the weights by this formula:

$$w_{ij}^{(\ell)} \longrightarrow w_{ij}^{(\ell)} - \eta \delta_j^{(\ell)} z_i^{(\ell-1)}$$

6. Iterate this method until it is time to stop.

7. Return the final weights.

Chain rule and partial differentiation

Regularization

Why learning is feasible?

$$y = f(\mathbf{z})$$

hypothesis set

learning algorithm finds an element g in \mathcal{H}

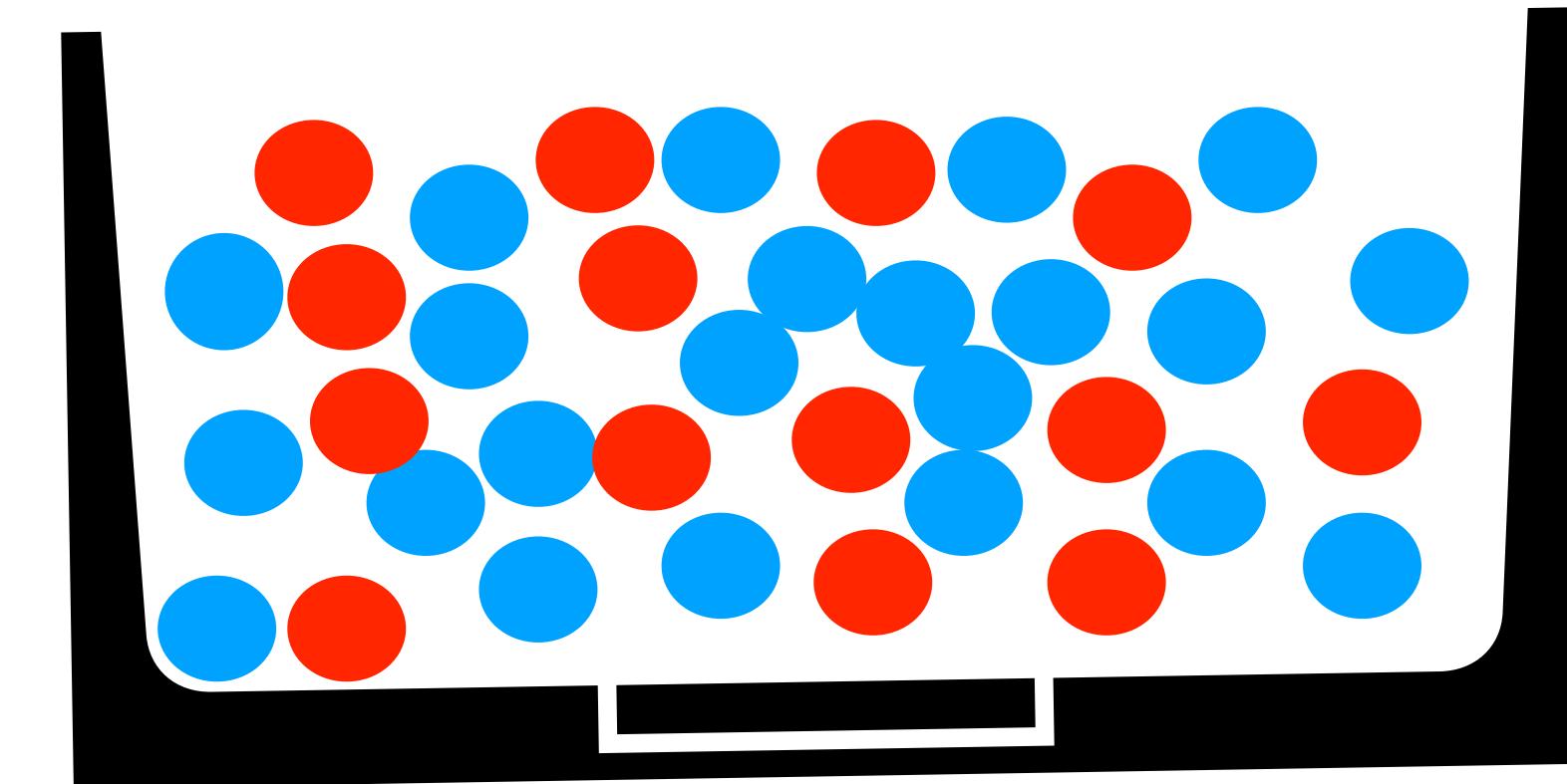
$$g \approx f$$

$$\mu = P(red)$$

$$1 - \mu = P(blue)$$

$$\mu \ ?$$

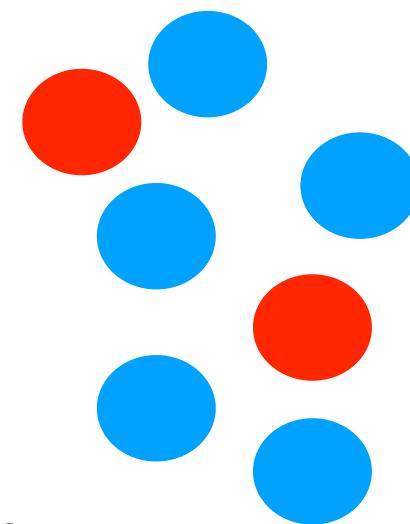
bin with blue and red marbles



sample (training)

$$\nu$$

proportion of red in the sample



$$\nu \longrightarrow \mu \ ? \text{ Probably}$$

$$P(|\nu - \mu| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

N size of sample

Hoeffding's inequality

Relation to learning

$$h(z) = f(z) \quad z$$

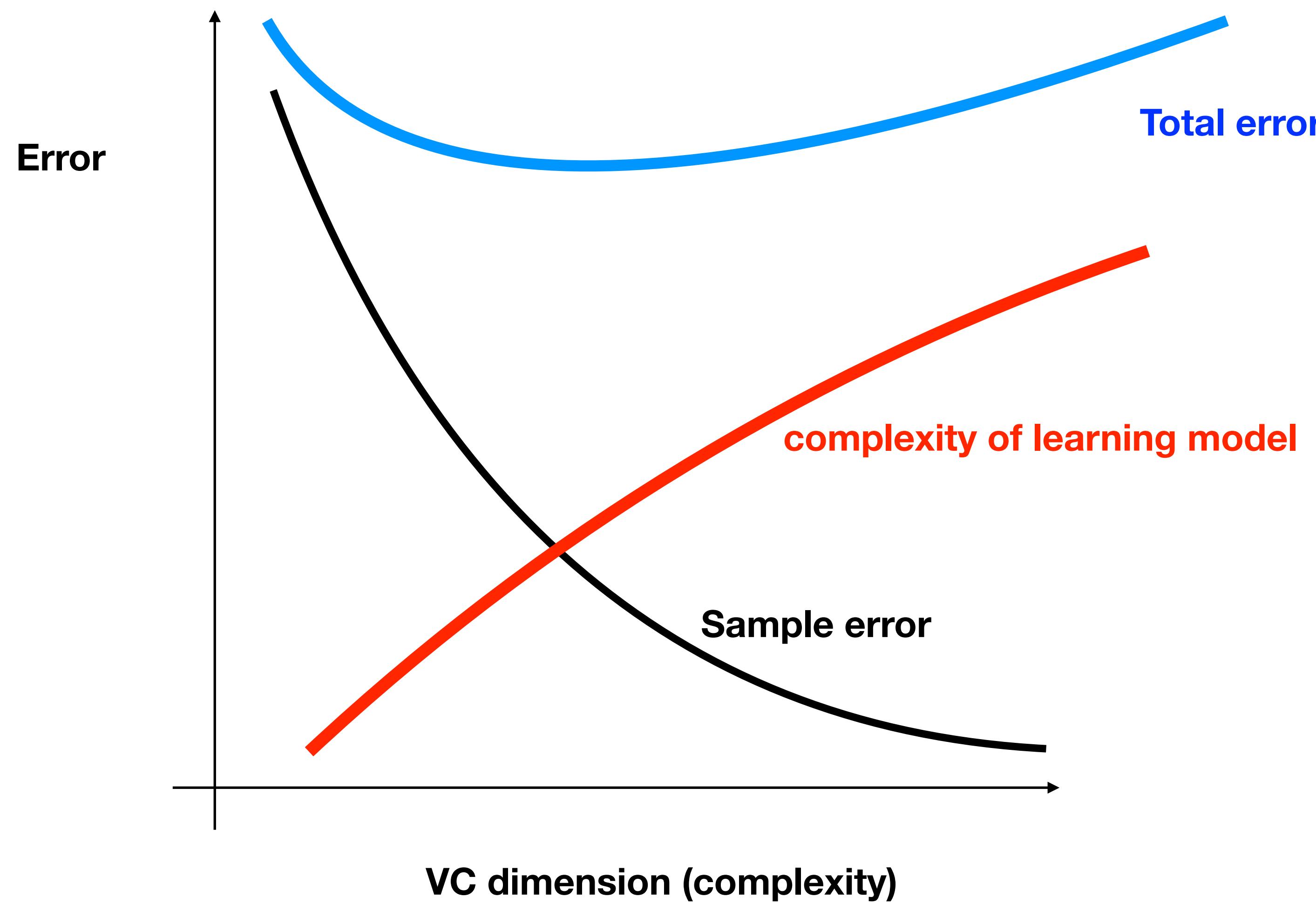
$$h(z) \neq f(z) \quad z \quad \mathcal{H} = \{h_1, h_2, \dots, h_M\}$$

$$P(|E_{sample}(g) - E_{total}(g)| > \epsilon) \leq 2M e^{-2\epsilon^2 N}$$

Learning well-done if:

$$|E_{sample}(g) - E_{total}(g)| \quad \text{small}$$

$$E_{sample}(g) \quad \text{small}$$



$$E_{total} \leq E_{sample} + \Omega$$

Probabilistic inequality for generalization bound

$$\Omega = \Omega(N, \mathcal{H}, \delta)$$

The inequality holds with probability (at least)

$$1 - \delta$$

decreases as N increases

**increases as VC dimension or
complexity of learning model
increases**

Bias-Variance tradeoff

\mathcal{H}

large: more chances of having the target function there but takes a long time to get close to it, large total error

small: better generalization; far from target

$$E_{total}(g^{(D)}) = \mathbb{E} \left((g^{(D)}(z) - f(z))^2 \right)$$

square error summation

$$\bar{g}(z) = \mathbb{E}_D \left(g^{(D)}(z) \right) \approx \frac{1}{K} \sum_{k=1}^K g^{D_k}(z)$$

D_1, D_2, \dots, D_k

k data sets

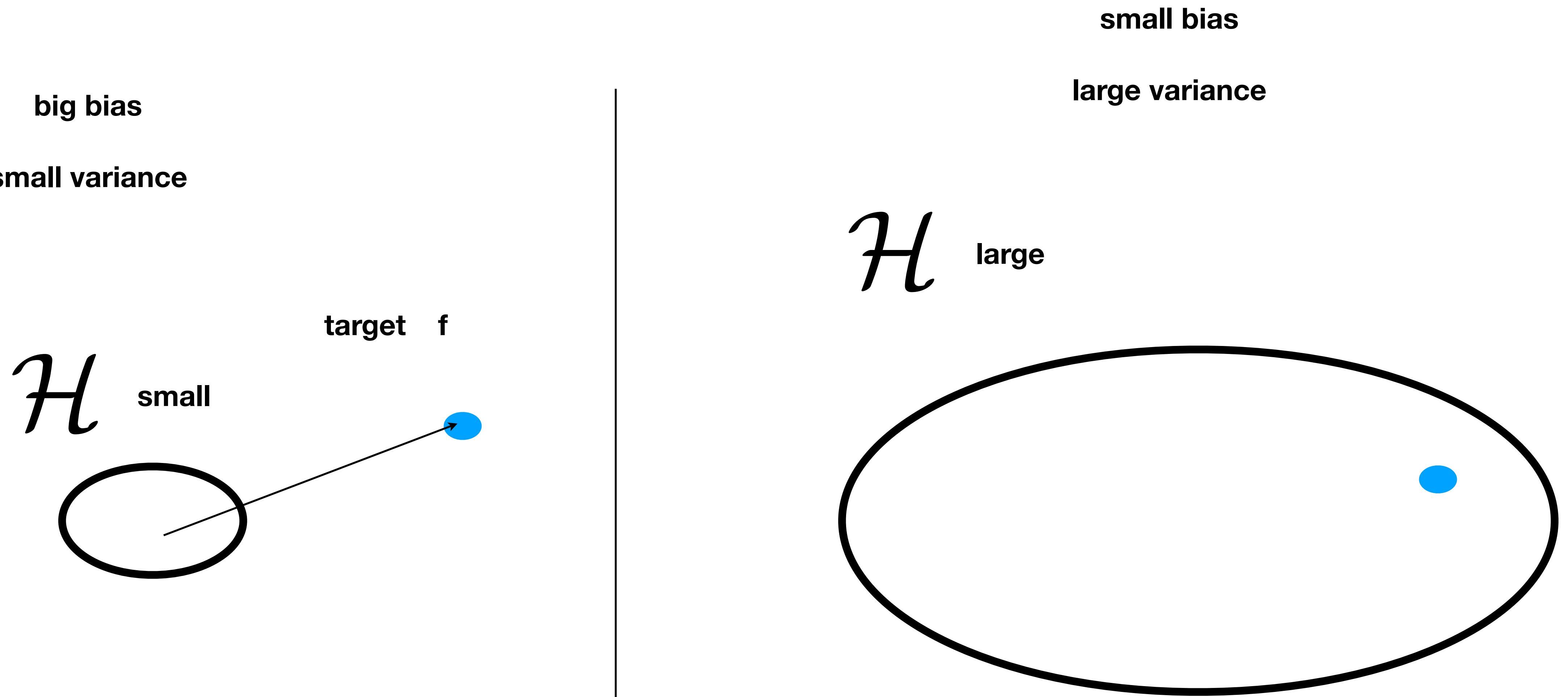
$$\begin{aligned}\mathbb{E}_D \left((g^{(D)}(z) - f(z))^2 \right) &= \mathbb{E}_D \left((g^{(D)}(z) - \bar{g}(z))^2 \right) + (\bar{g}(z) - f(z))^2 \\ &= \text{variance(z)} + \text{bias (z)}\end{aligned}$$

Variance by definition: how much we might vary!

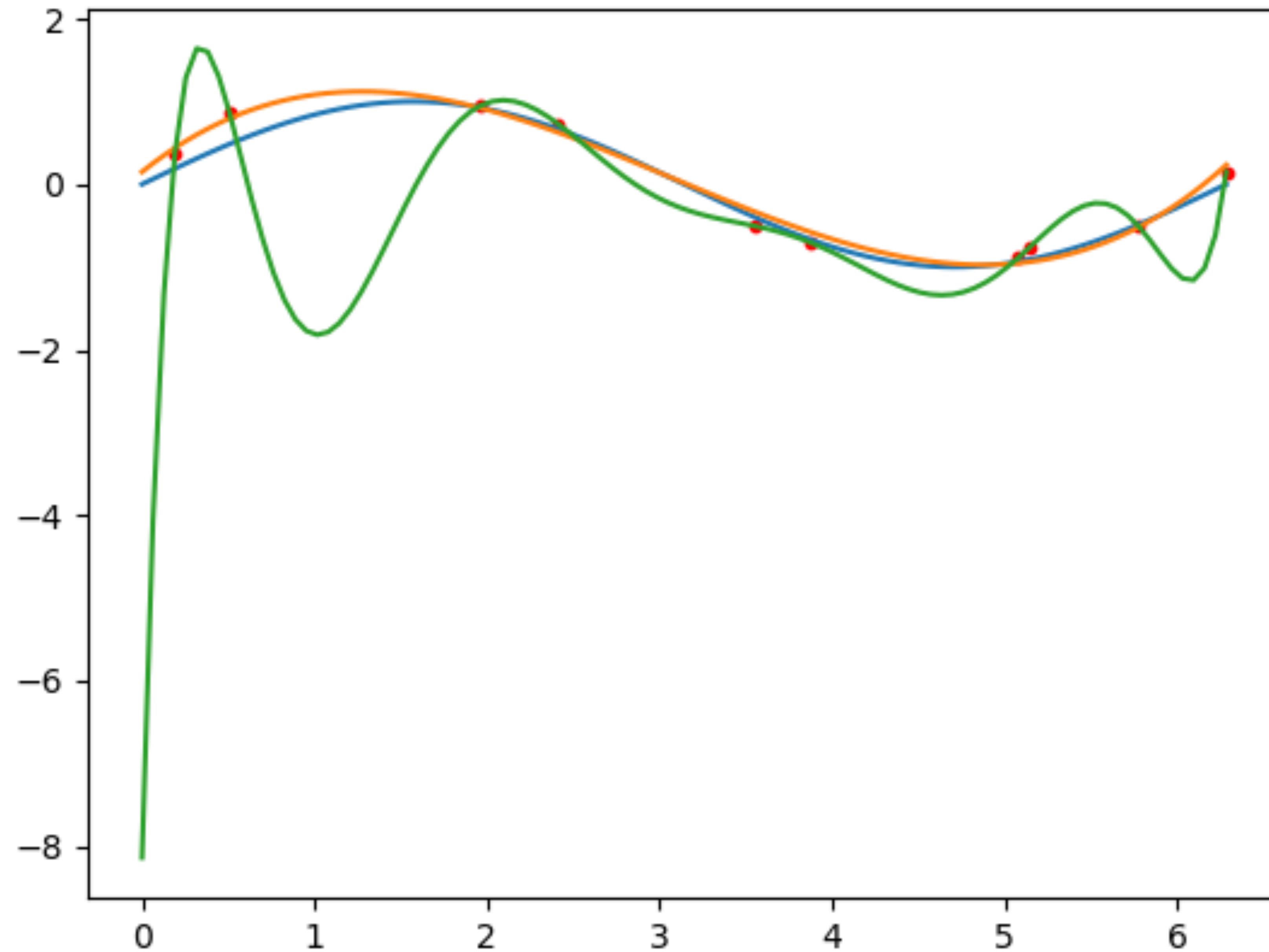
how biased away from the target function we are!

$$\mathbb{E}_D \left(E_{out}(g^{(D)}) \right) = \text{bias} + \text{variance}$$

The choice of complexity should be determined by the available data.



**overfitting
(and
underfitting)**



Remedy for overfitting: regularization

minimizing under constraints on the weights

$$\mathbf{w} \cdot \mathbf{w} = \sum_{i=0}^n w_i^2 \leq C$$

Suggested
reading:

Lagrange
multiplier
method

$$E_{extended} = E_{sample} + \frac{\lambda}{N} \mathbf{w} \cdot \mathbf{w}$$

Gradient Descent, and validation

Convolutional neural networks (CNN)

Convolution, pooling, etc

$$x * w(0) = \sum_{i=-\infty}^{\infty} x(-i)w(i)$$

$$x * w = w * x$$

$$(x * w)(t) = \sum_i x(t - i)w(i)$$

$$(I * K)(m, n) = \sum_{i,j} I(m - i, n - j)K(i, j)$$

I = Image

2 3 4 -1

1 2 -1 1

3 -2 0 1

K = kernel

Convolved feature

*1 *1

*0 *0

6 9 2

6 -1 0

Convolution



average pooling

→ 21/3 10/3

A convolutional neural network is network that has at least one convolutional layer in its hidden layers

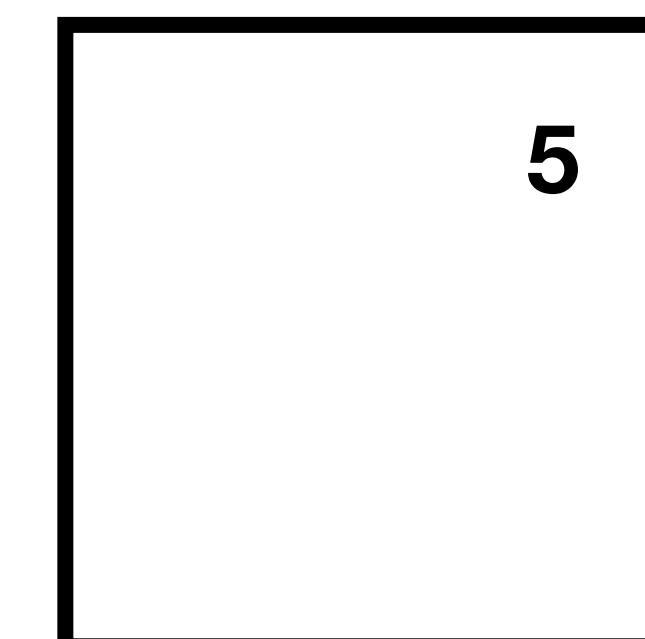
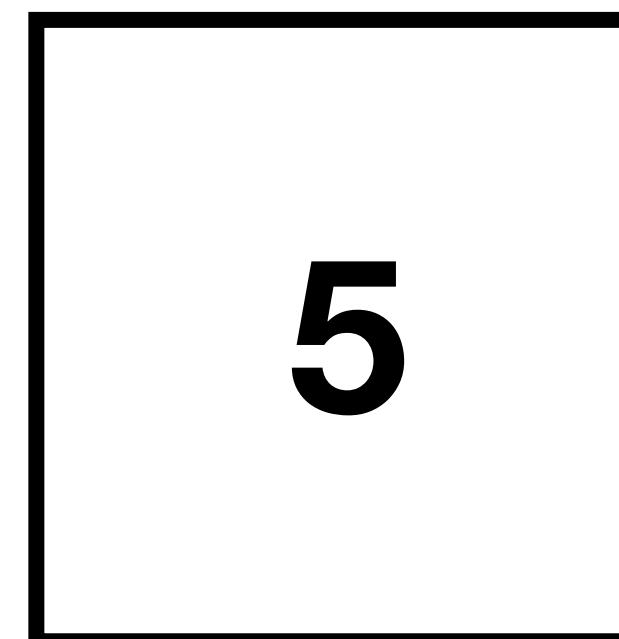
Convolutional layer:

1. **Convolution**
2. **Activation function (nonlinear): detection**
3. **Pooling**

Pooling: Calculating a summary statistics of the output of stage 2 such as mean, max, norm, etc.

e.g. maximum of a rectangular neighborhood

advantage: makes the training less sensitive to translations, therefore digits located in different places can be detected



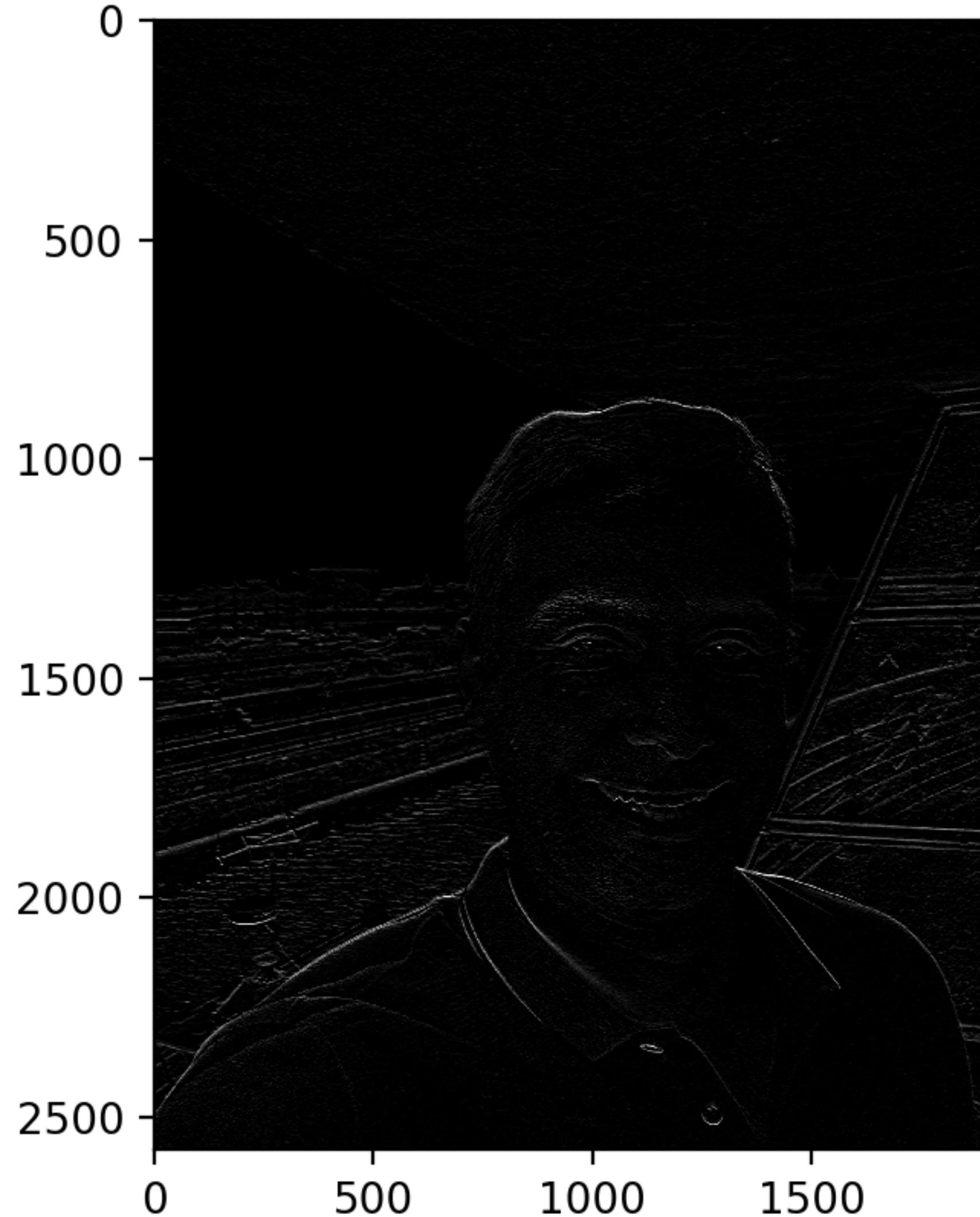


*1

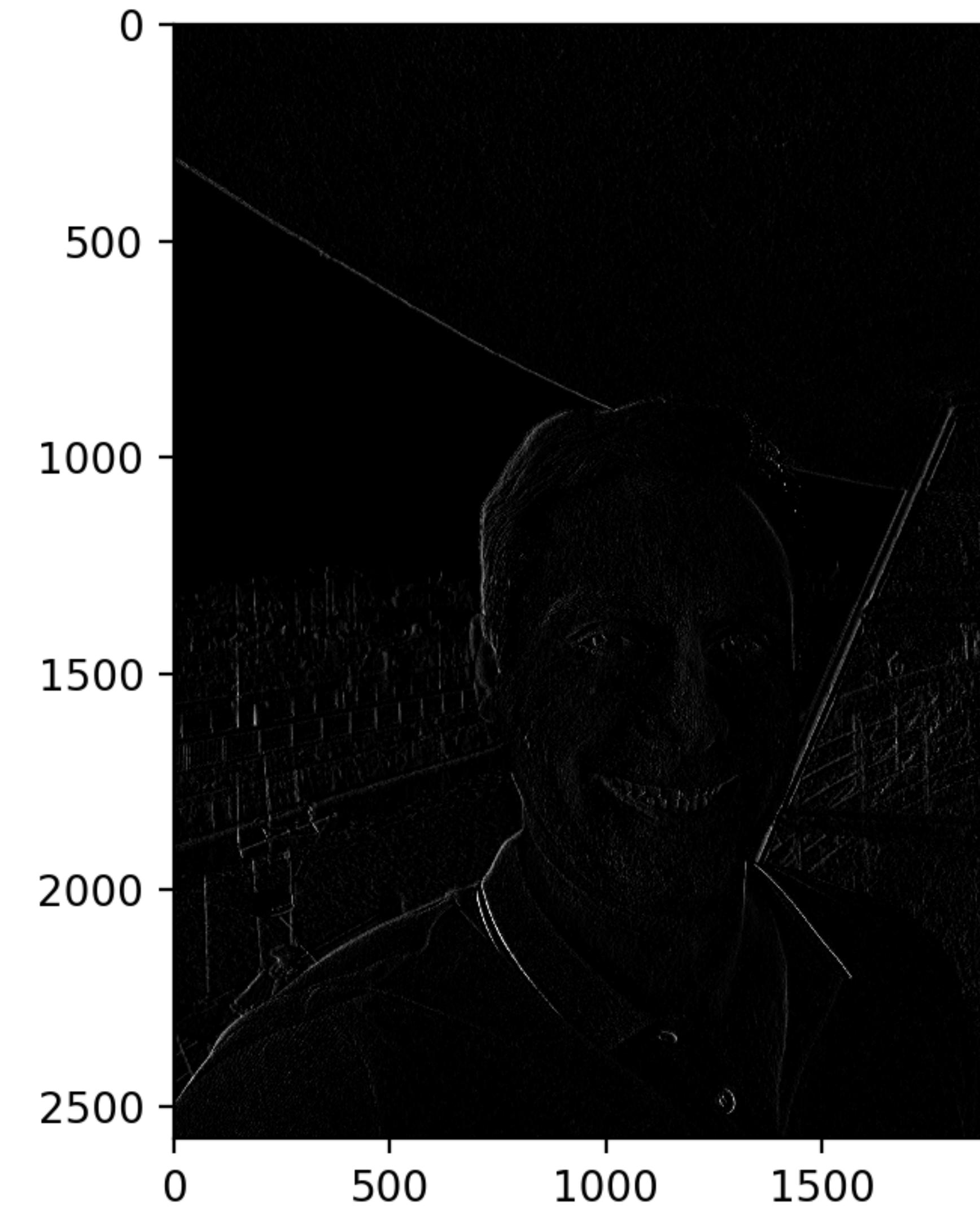
*1

*0

*0



1, 1, 1
0, 0, 0
-1, -1, -1



1, 0, -1
1, 0, -1
1, 0, -1

Recurrent neural networks (RNN)

Sequential data: language processing, sentences, ...

The order of the data is important.

Predicting the forthcoming data.

Time series analysis: prediction of stock prices

Dealing with issues such as: length of input data changes

**Need a number of neural networks connected to each other:
The connection is made by connecting their hidden layers**

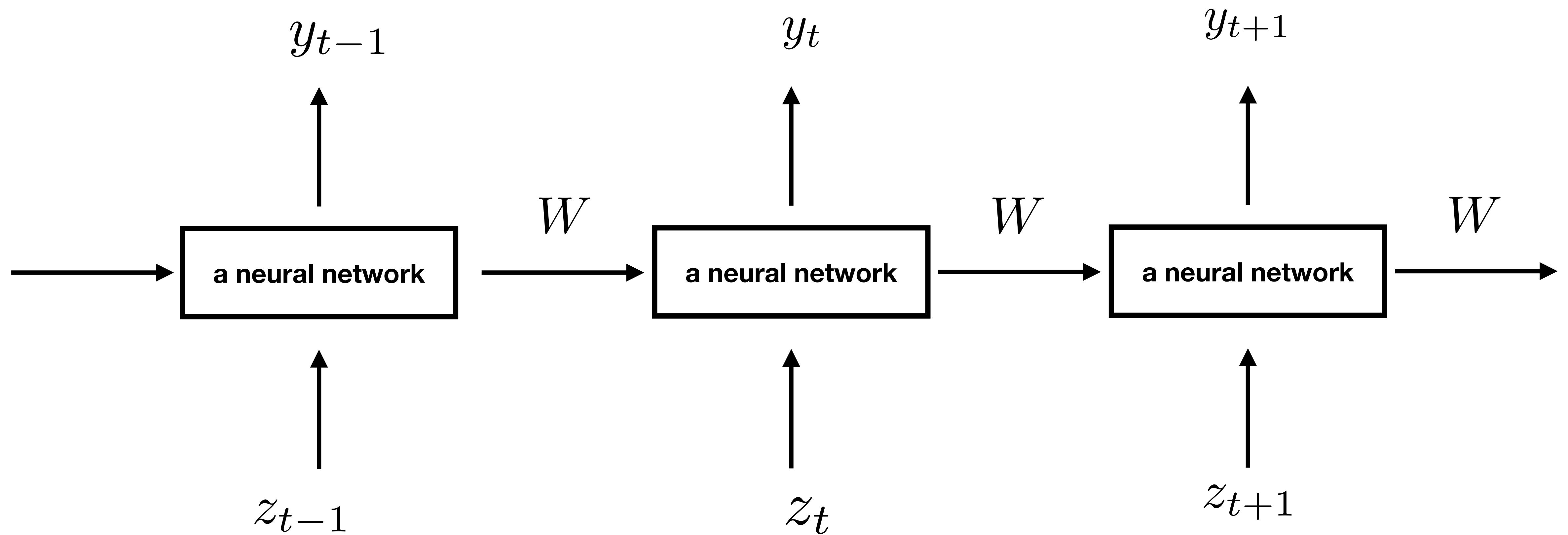
Parameter sharing:
Same W between the hidden layers

$$a_t = b + W s_{t-1} + U x_t$$

$$s_t = \tanh(a_t)$$

$$y_t = c + V s_t$$

$$p_t = \text{softmax}(y_t)$$



$$a_t = b + W s_{t-1} + U x_t$$

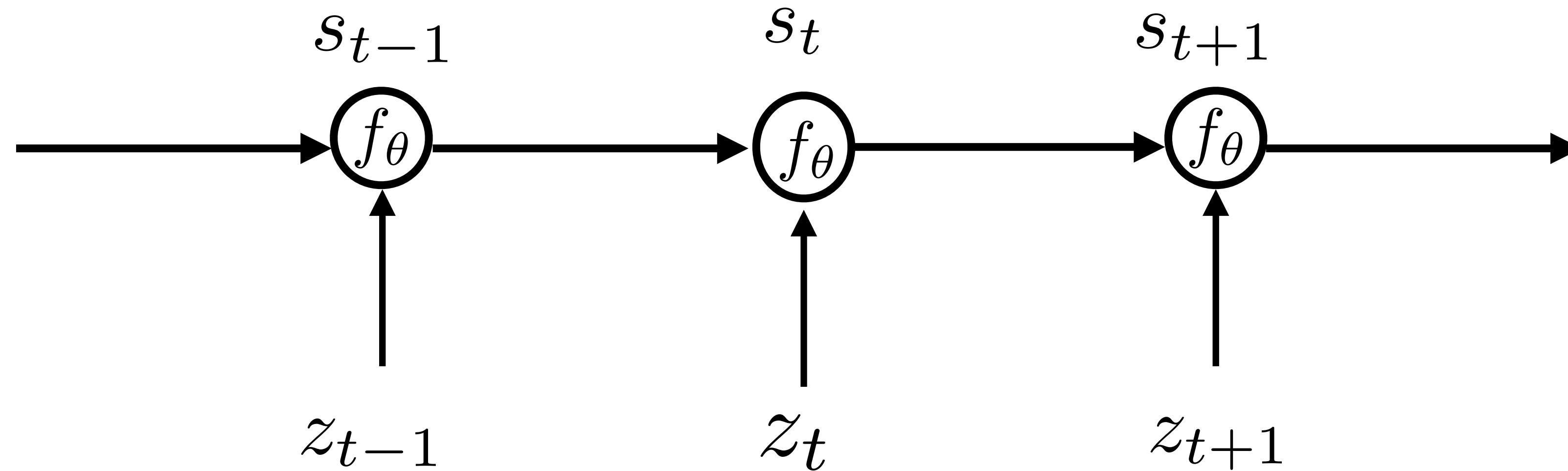
$$s_t = \tanh(a_t)$$

$$y_t = c + V s_t$$

$$p_t = softmax(y_t)$$

Dynamical system

$$s_t = f_\theta(s_{t-1}, z_t)$$



$$(3,2,1) = 3(1,0,0) + 2(0,1,0) + 1(0,0,1)$$

$$e_1=(1,0,0)$$

$$e_2=(0,1,0)$$

$$e_3=(0,0,1)$$

$$\mathbf{x}=(x_1,x_2,x_3)$$

$$= (\mathbf{x}\cdot e_1)e_1 + (\mathbf{x}\cdot e_2)e_2 + (\mathbf{x}\cdot e_3)e_3$$

$$e_n(x) = e^{2\pi i n x} = \cos(2\pi n x) + i \sin(2\pi n x) \quad \text{base functions}$$

$$n = 0, \pm 1, \pm 2, \dots$$

$$i^2 = -1$$

$$(x, y) = x + iy$$

$$(x, y) = x - iy$$

$$\langle f, g \rangle = f(x) \cdot g(x) = \int_0^1 f(x) \bar{g}(x) dx$$

f and g are periodic functions of period 1

$$f(x) = \sum_{n=-\infty}^{\infty} \langle f, e_n \rangle e_n(x)$$

n runs over all integers and period of f is 1!!!

$$\hat{f}(n) = \langle f, e_n \rangle$$

**n-th Fourier
coefficient**

$$\hat{f}(\xi) = \int_{\mathbb{R}^m} f(x) e^{-2\pi i \xi \cdot x} dx$$

**Fourier transform of a
function defined on all
real numbers (not periodic)**

$$f(x) = \int_{\mathbb{R}^m} \hat{f}(\xi) e^{2\pi i \xi \cdot x} d\xi$$

**Fourier inversion formula:
reconstruction from Fourier
transform**

$$(f * g)(x) = \int f(x - t)g(t) dt$$

$$(\hat{f} * \hat{g})(x) = \hat{f}(x) \hat{g}(x)$$

Fast Fourier transform

$$f(0), f(1), \dots, f(N-1)$$

$$\hat{f}(k) = \sum_{n=0}^{N-1} f(n) e^{i(-2\pi kn/N)}$$

$$e^{ix} = \cos(x) + i \sin(x)$$

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

$$\sin(x) = \frac{\exp(ix) - \exp(-ix)}{2i}$$

$$f(n) = \sum_{k=0}^{N-1} \hat{f}(k) e^{i(2\pi kn/N)}$$

k: **-(N-1)/2 , ... , -2, -1, 0, 1, 2, ..., (N-1)/2**
assuming N is odd

**Fast Fourier transform
2 dimensional**

$$f(0, 0), f(1, 0), \dots, f(N_1 - 1, 0)$$

$$N_1 \times N_2$$

$$f(0, 1), f(1, 1), \dots, f(N_1 - 1, 1)$$

• • •

$$f(0, N_2 - 1), f(1, N_2 - 1), \dots, f(N_1 - 1, N_2 - 1)$$

$$\hat{f}(k_1, k_2) = \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} f(n_1, n_2) e^{i(-2\pi k_1 n_1 / N_1)} e^{i(-2\pi k_2 n_2 / N_2)}$$

$$f(n_1, n_2) = \sum_{k_1=0}^{N_1-1} \sum_{k_2=0}^{N_2-1} \hat{f}(k_1, k_2) e^{i(2\pi k_1 n_1 / N_1)} e^{i(2\pi k_2 n_2 / N_2)}$$

Signal detection in electrophysiological measurements

Multi-modal time series of neural activities are accompanied by noise and our goal has been to detect the signals and spikes in such series by way of analyzing their intrinsic statistical properties.

Our approach has been able to **sort the spikes** indicating a difference in their formation, and to **detect faint signals**.

Our method has the power to deal **efficiently** with correlated time series and **correlated noise**.

Statistical inference

Consider an **assumption H** about a sample X.

Calculate the probability that X would be observed given that H is true:

$$p = P(X|H)$$

Choose a *significance level* (a small number) such as:

$$\alpha = 0.05$$

If

$$p \leq \alpha$$

infer that the hypothesis **H** should be **rejected**.

Example of statistical inference

X : a poisson distribution such as the number of accidents on a street in one week with mean 5

$$P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!}, \quad i = 0, 1, 2, \dots$$

Hypothesis **H**: $\lambda = 5$

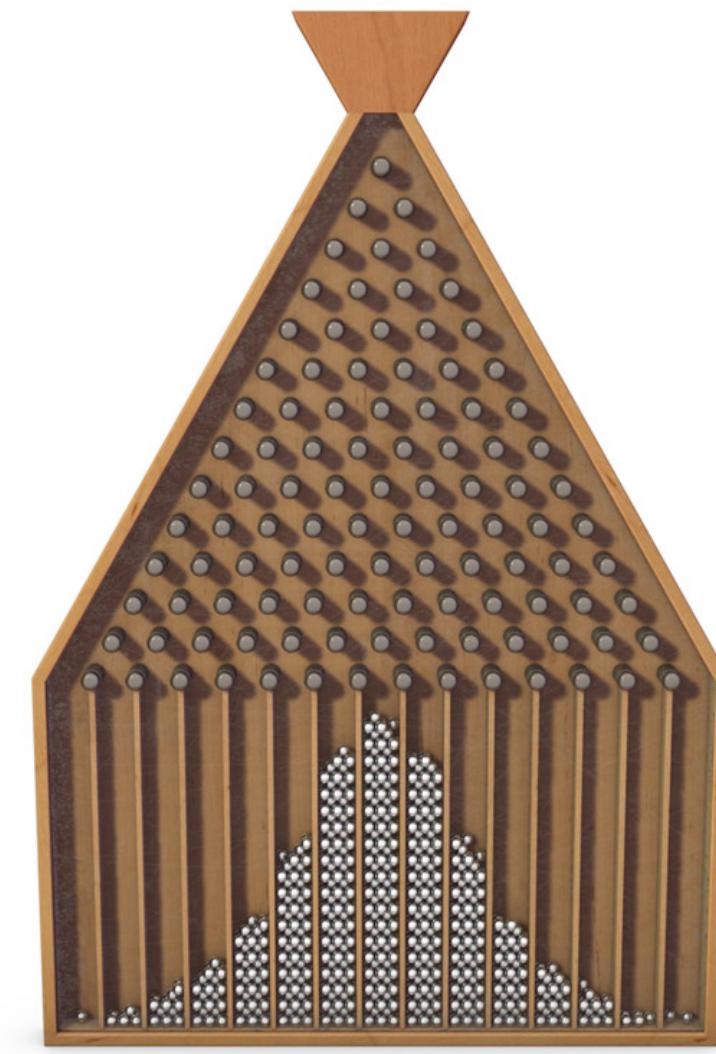
$$P(X = 10|H) = 0.0181328$$

So based on one week that incidentally there were 10 accidents
one will **reject the correct hypothesis**.

Multiple comparisons

How to make statistical inference about a hypothesis from many samples ?

Fundamental idea: To exploit **Bell shaped** distributions that form in nature prevalently when the outcome of an event is the result of the **independent** contribution of many effects with **the same distribution**.



Galton board

De Moivre-Laplace theorem (\mathcal{CLT})

$$f(x) = (1/2\pi) e^{-x^2/2}$$

$$P(\mathcal{N}(0, 1) > 2) = 0.022$$

$$X_1, \dots, X_m$$

$$Y_1, \dots, Y_m$$

defined as:

$$Y_j = 1 \quad \text{if } j\text{-th inference rejects } \mathbf{H} \text{ at significance } \alpha$$

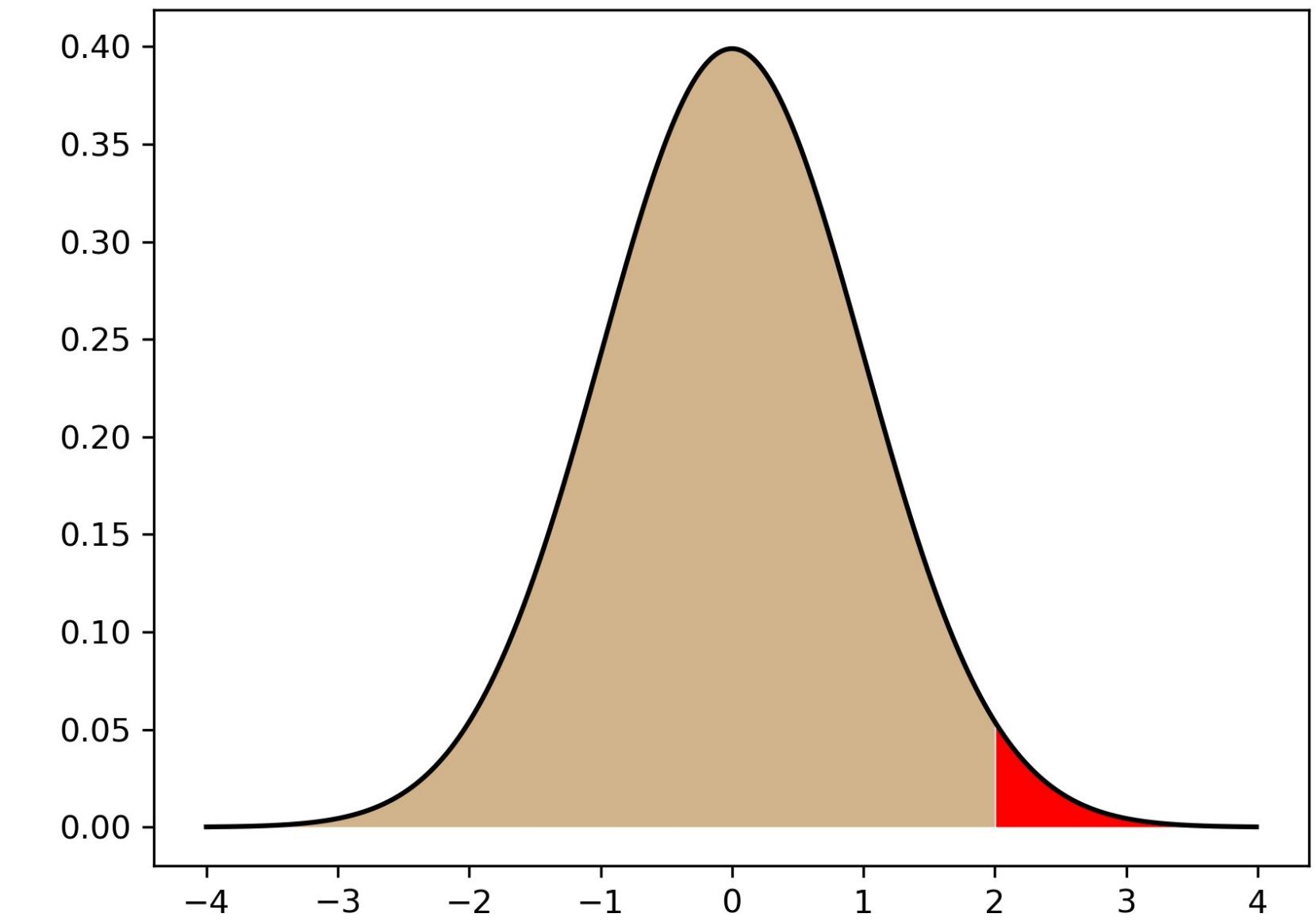
$$Y_j = 0 \quad \text{otherwise}$$

$$HC_{m,\alpha} = \frac{(Y_1 + \dots + Y_m)/m - \alpha}{\sqrt{\alpha(1-\alpha)}/\sqrt{m}}$$

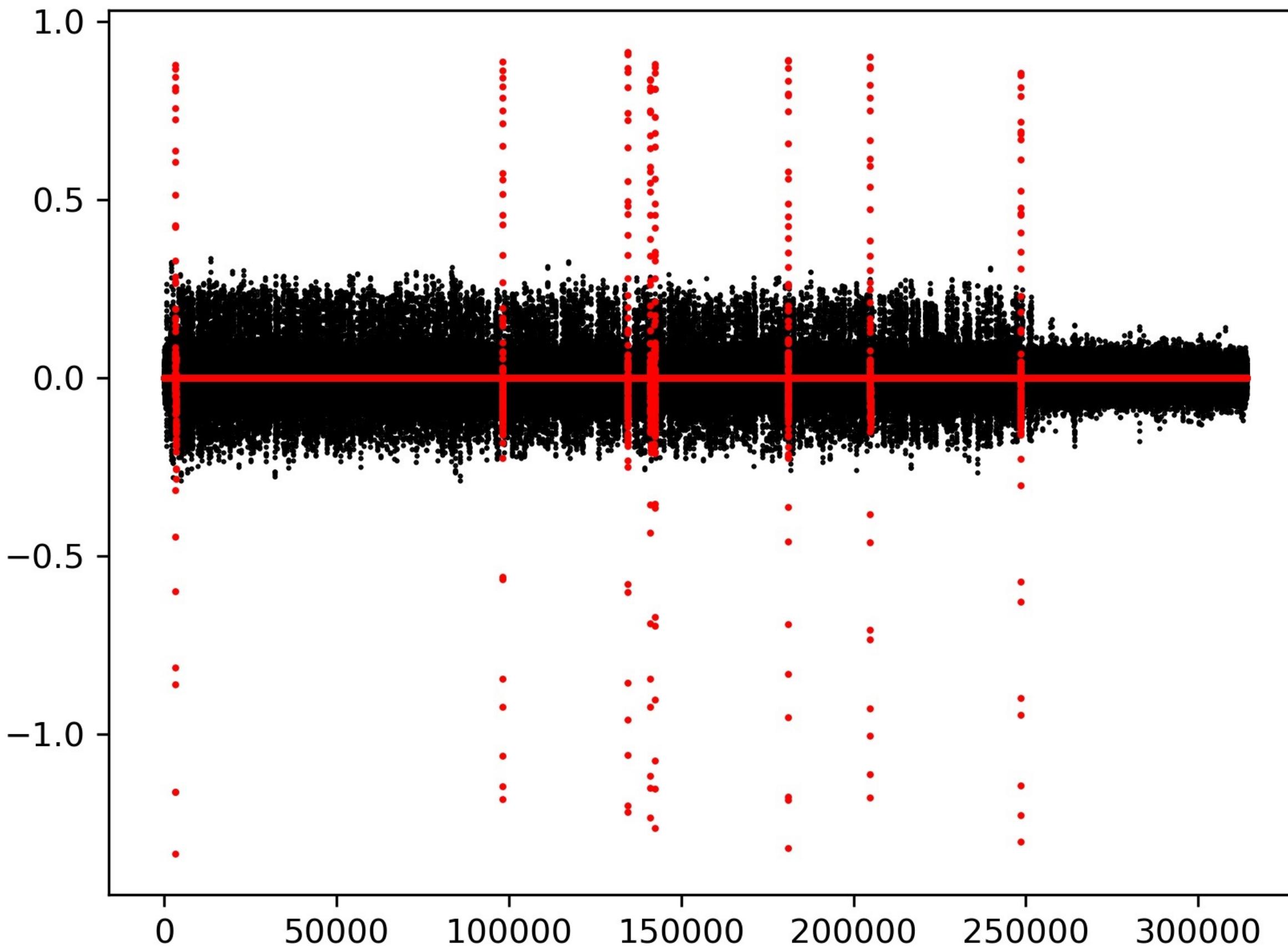
approaches $\mathcal{N}(0, 1)$ as $m \rightarrow \infty$

$$HC_m = \max HC_{m,\alpha}$$

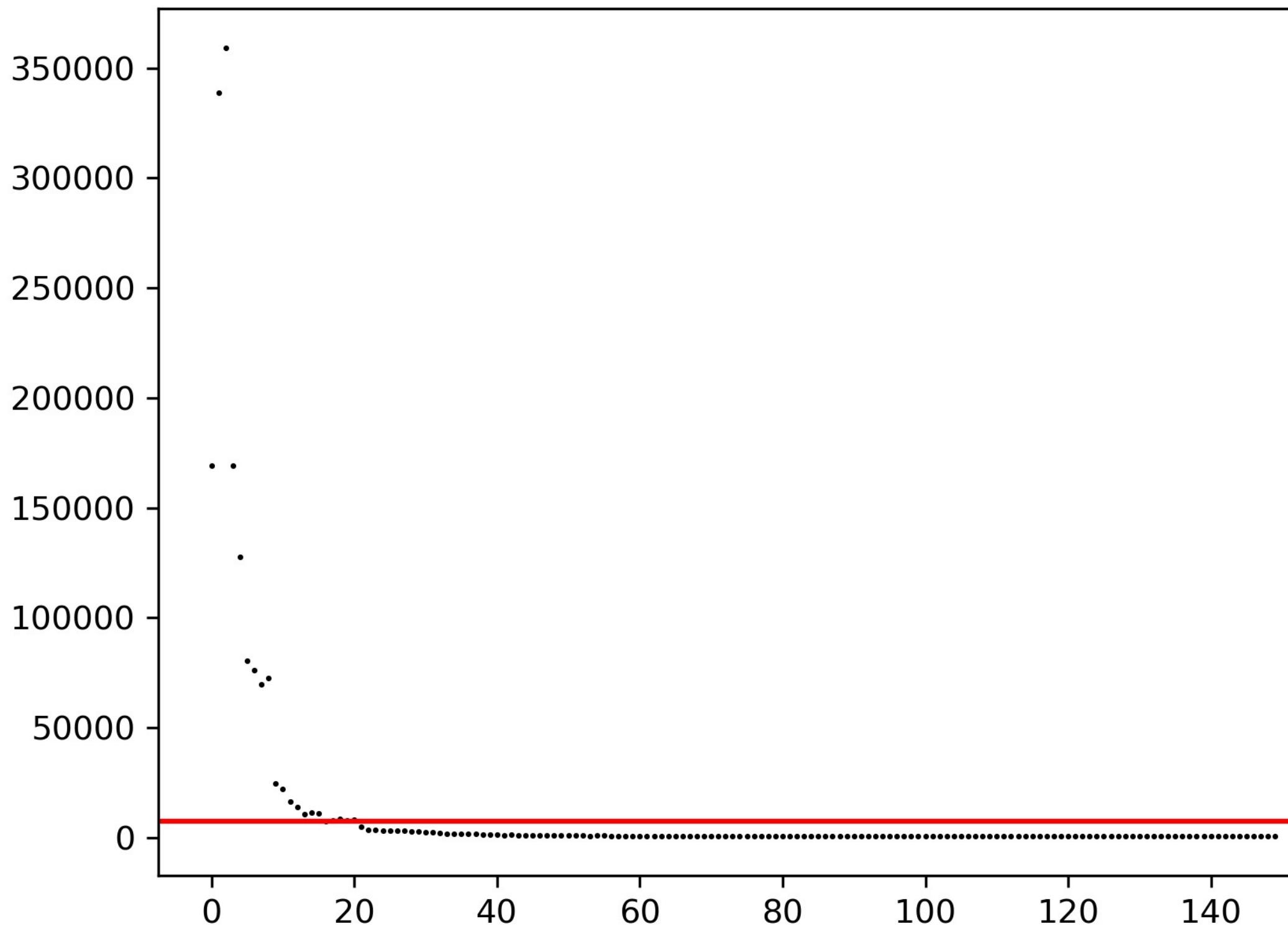
Tukey: reject \mathbf{H} if $HC_m > 2$



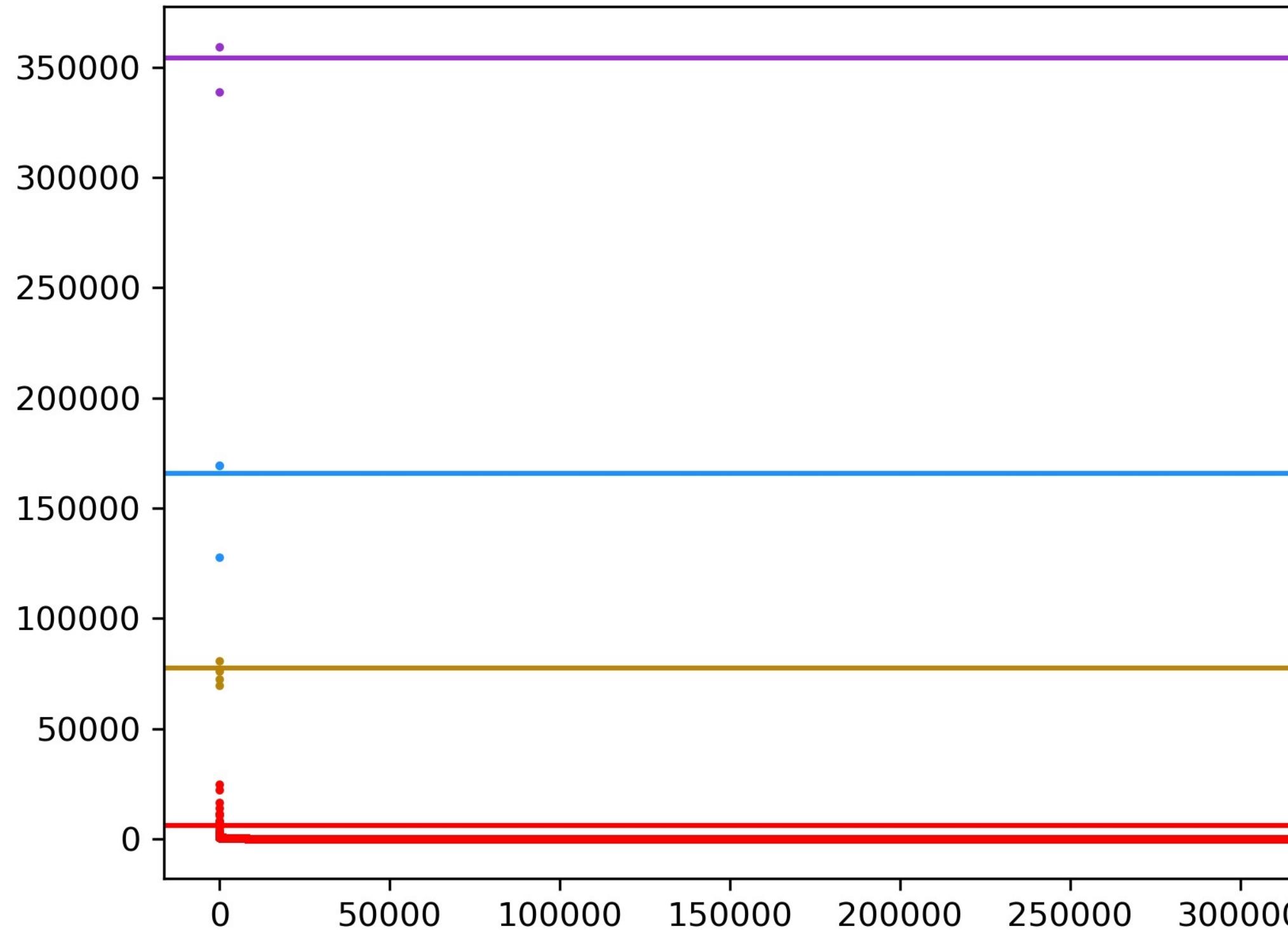
Detection of electrophysiological data spikes with higher criticism

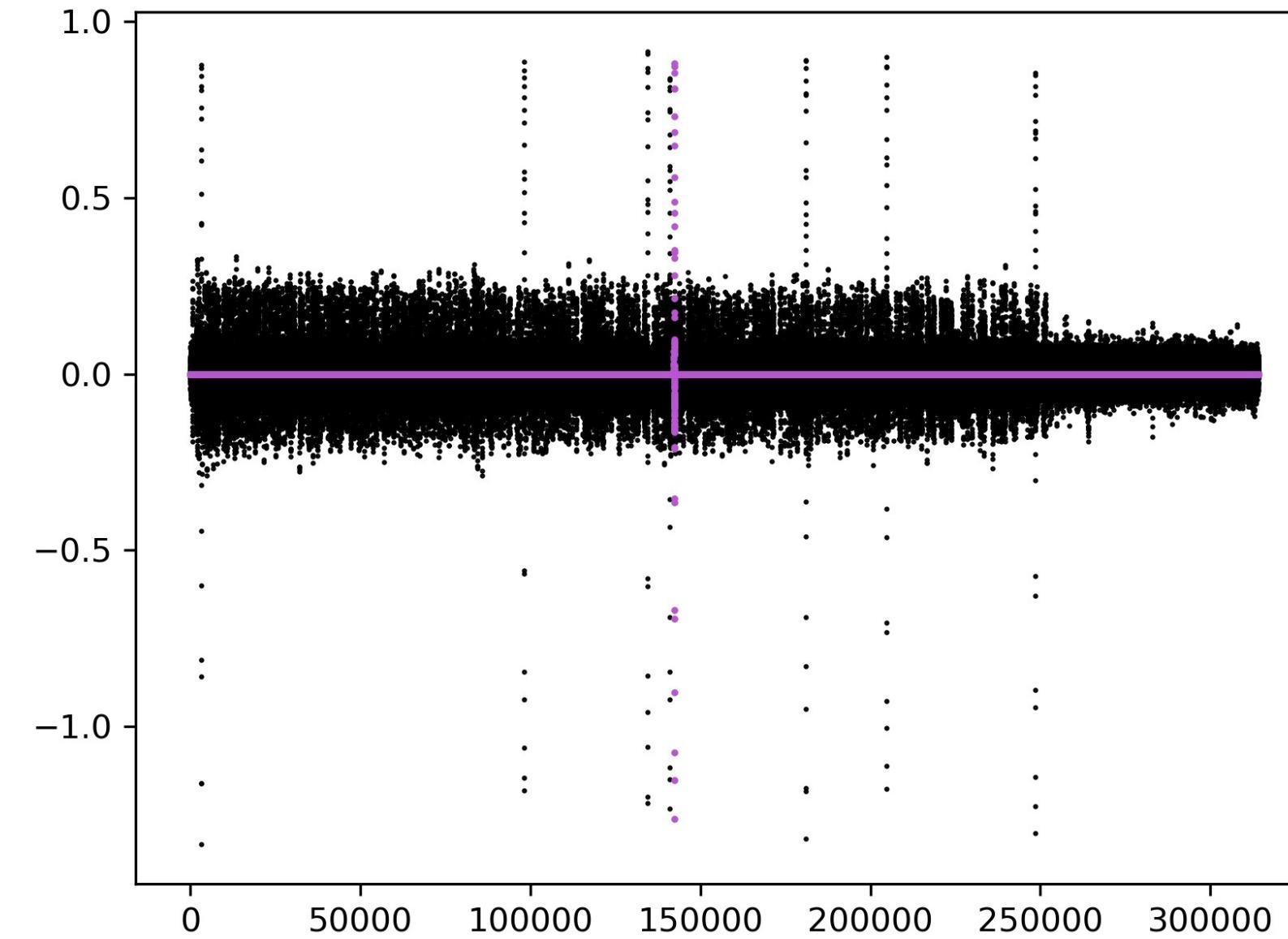
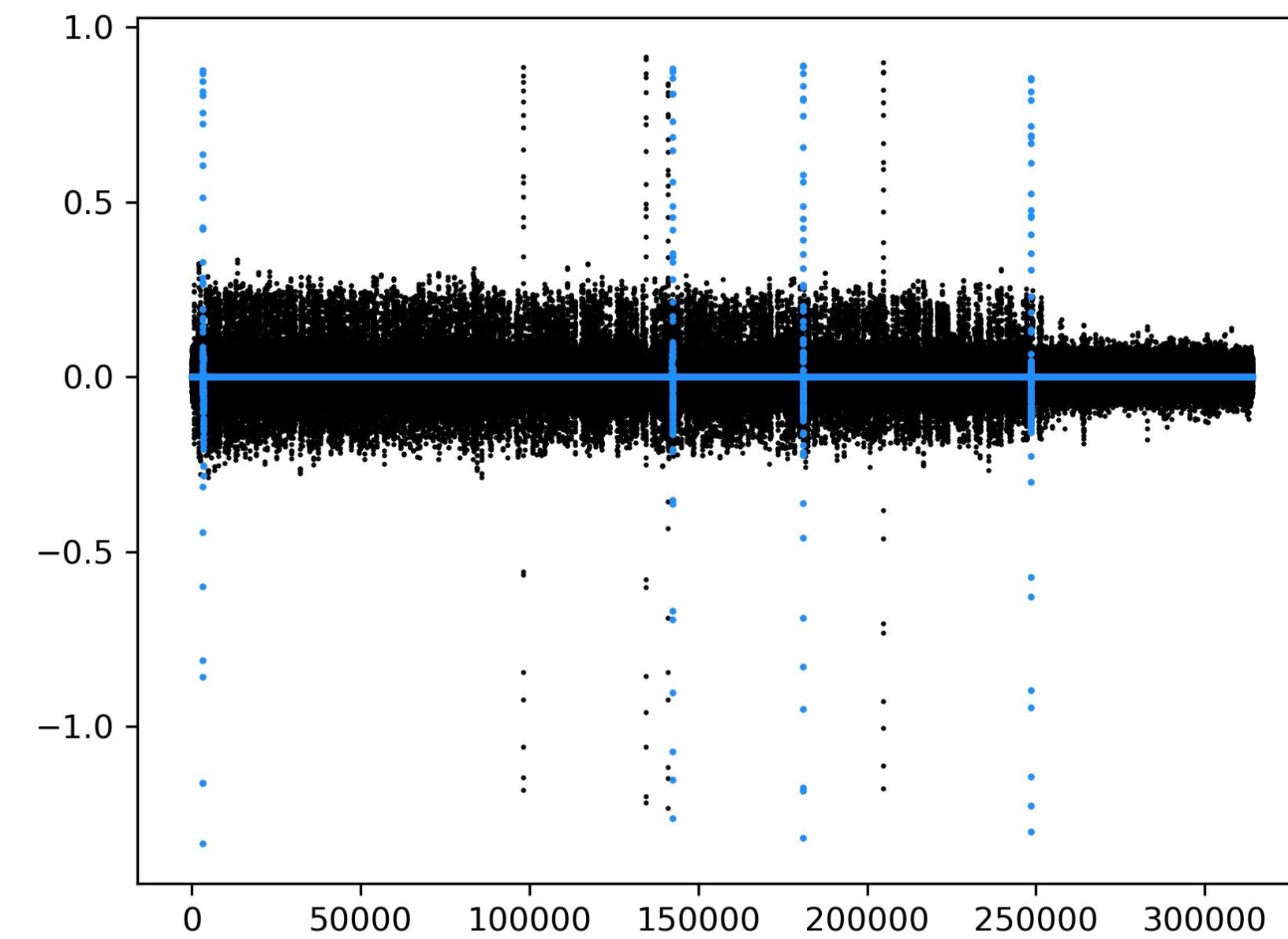
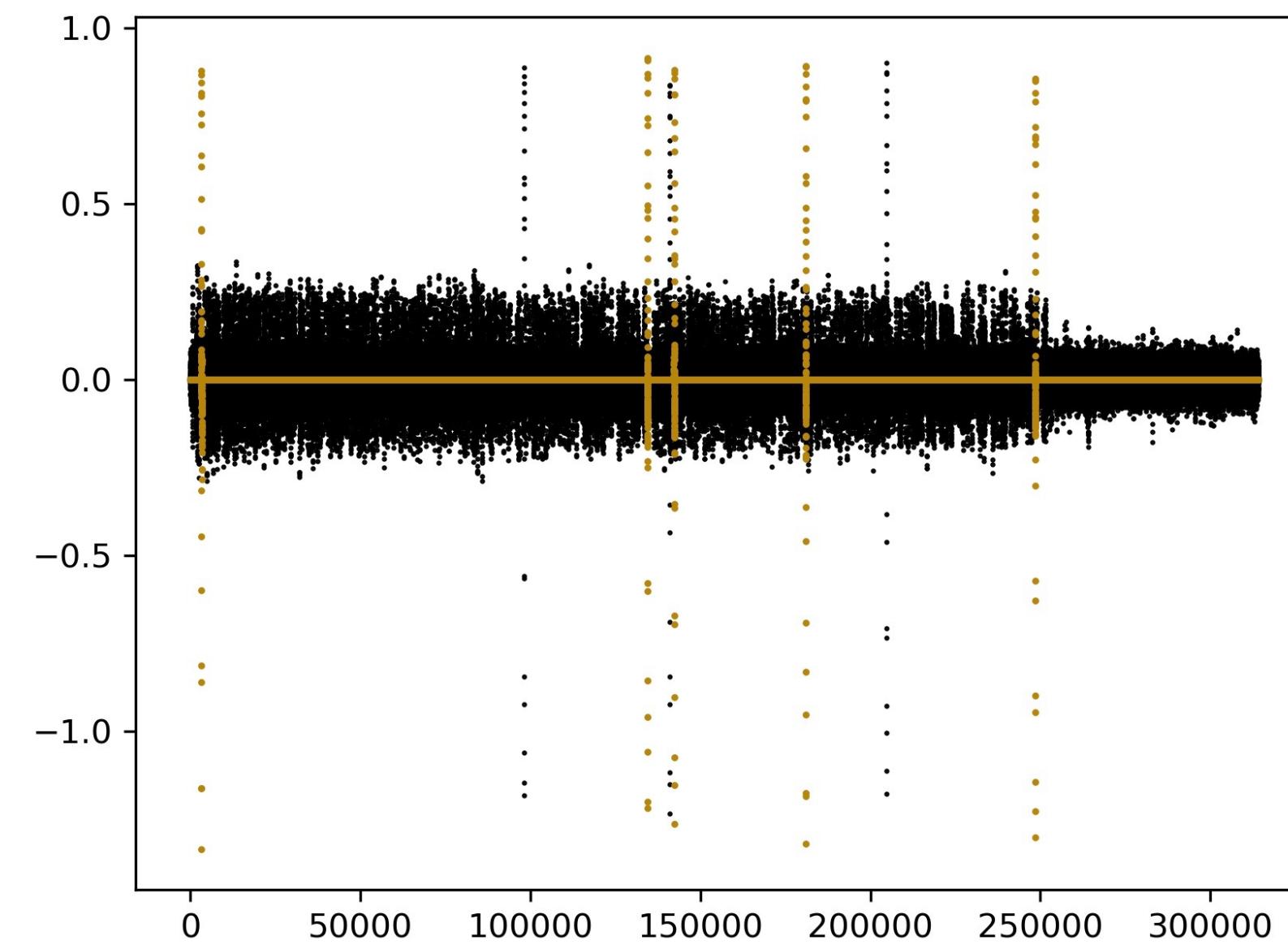
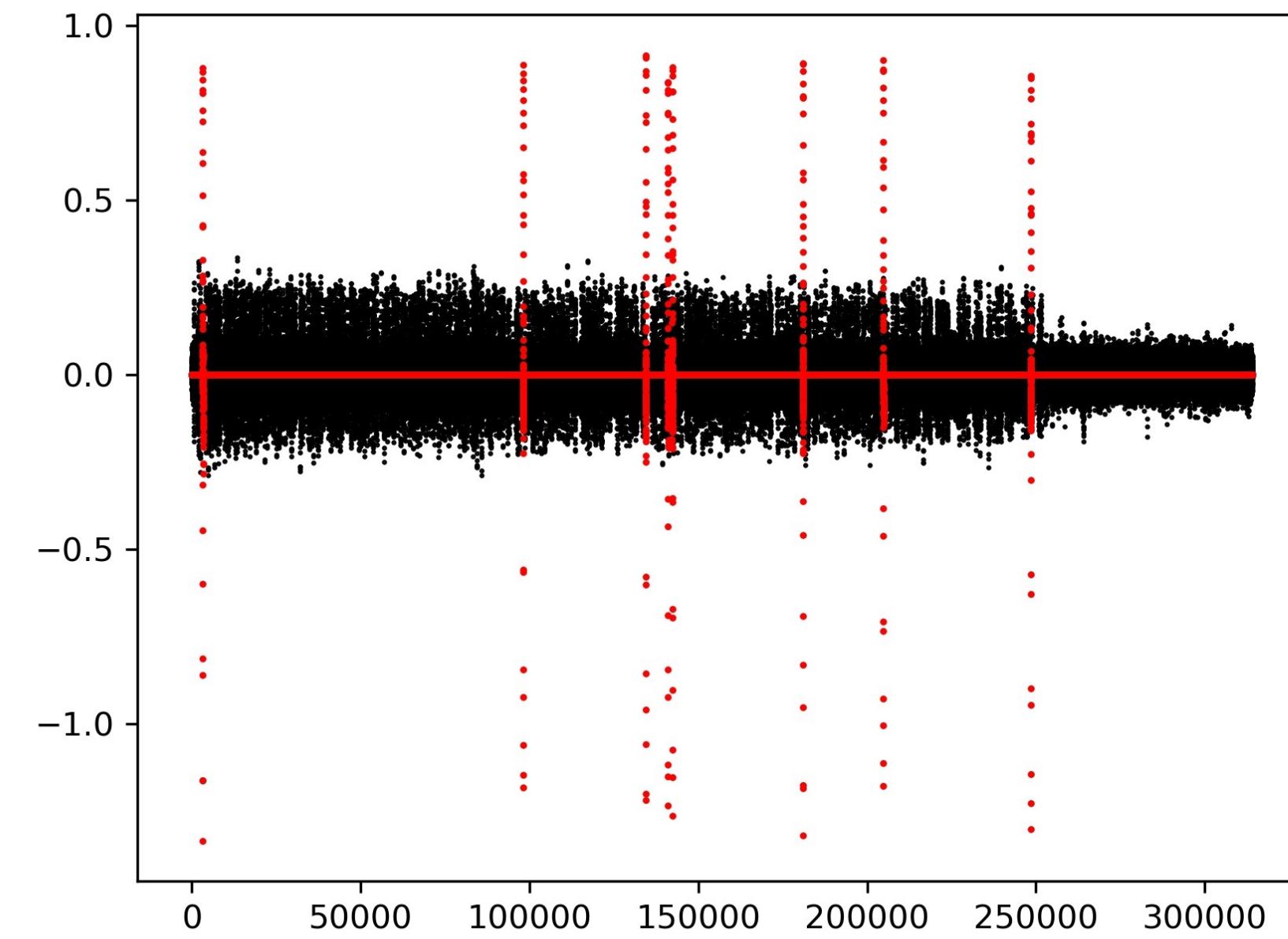


$\mathcal{H}\mathcal{C}$ quantities associated with the electrophysiological time series



$\mathcal{H}\mathcal{C}$ quantities clustered and thresholds associated with the clusters





Spike sorting

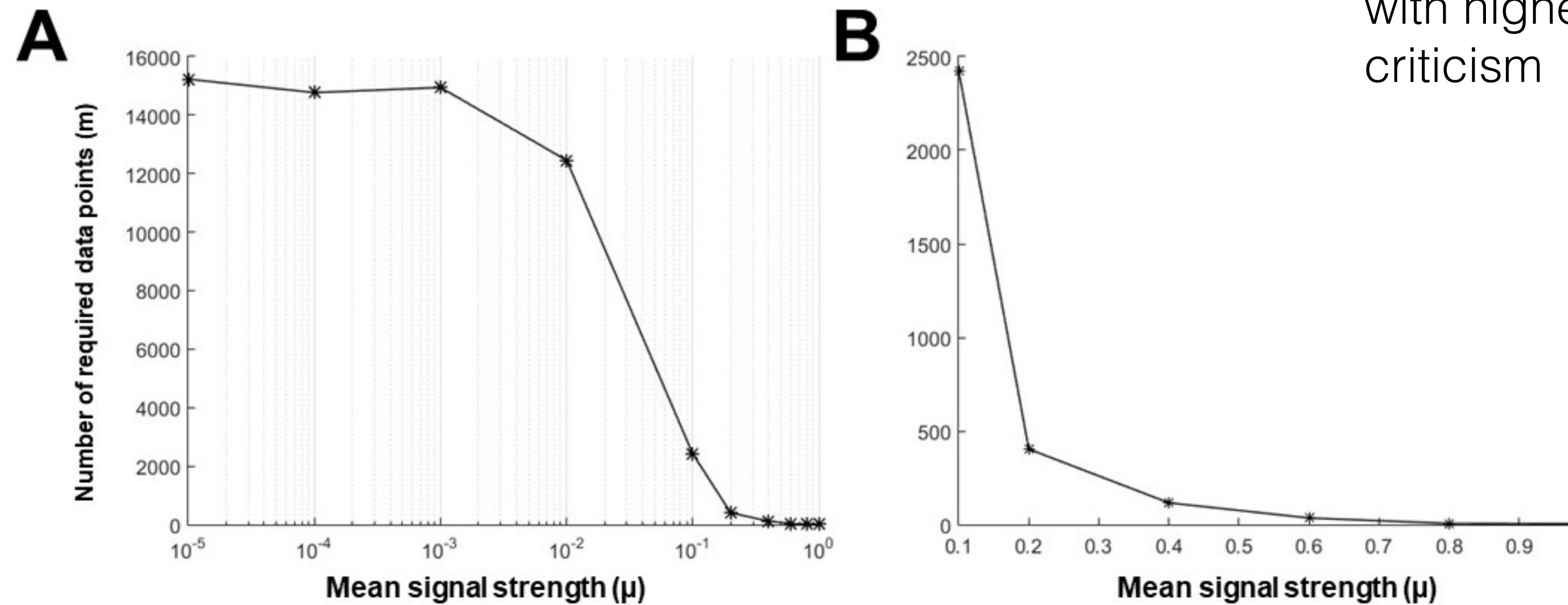
How many data points do we need to detect weak signals ?

$$X \sim \mathcal{N}(\mu, 1)$$

To detect small
non-zero

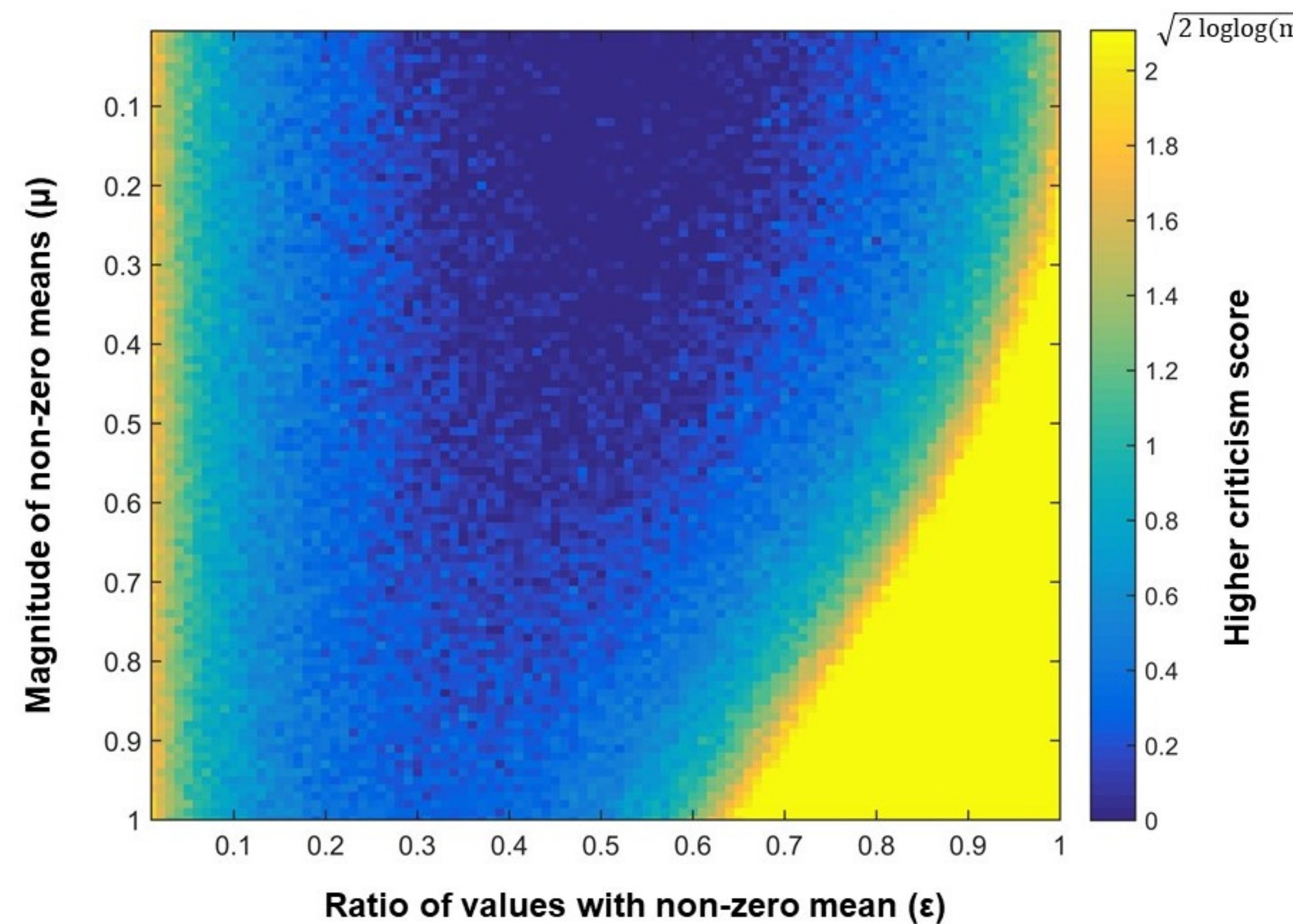
$$\mu$$

with higher
criticism

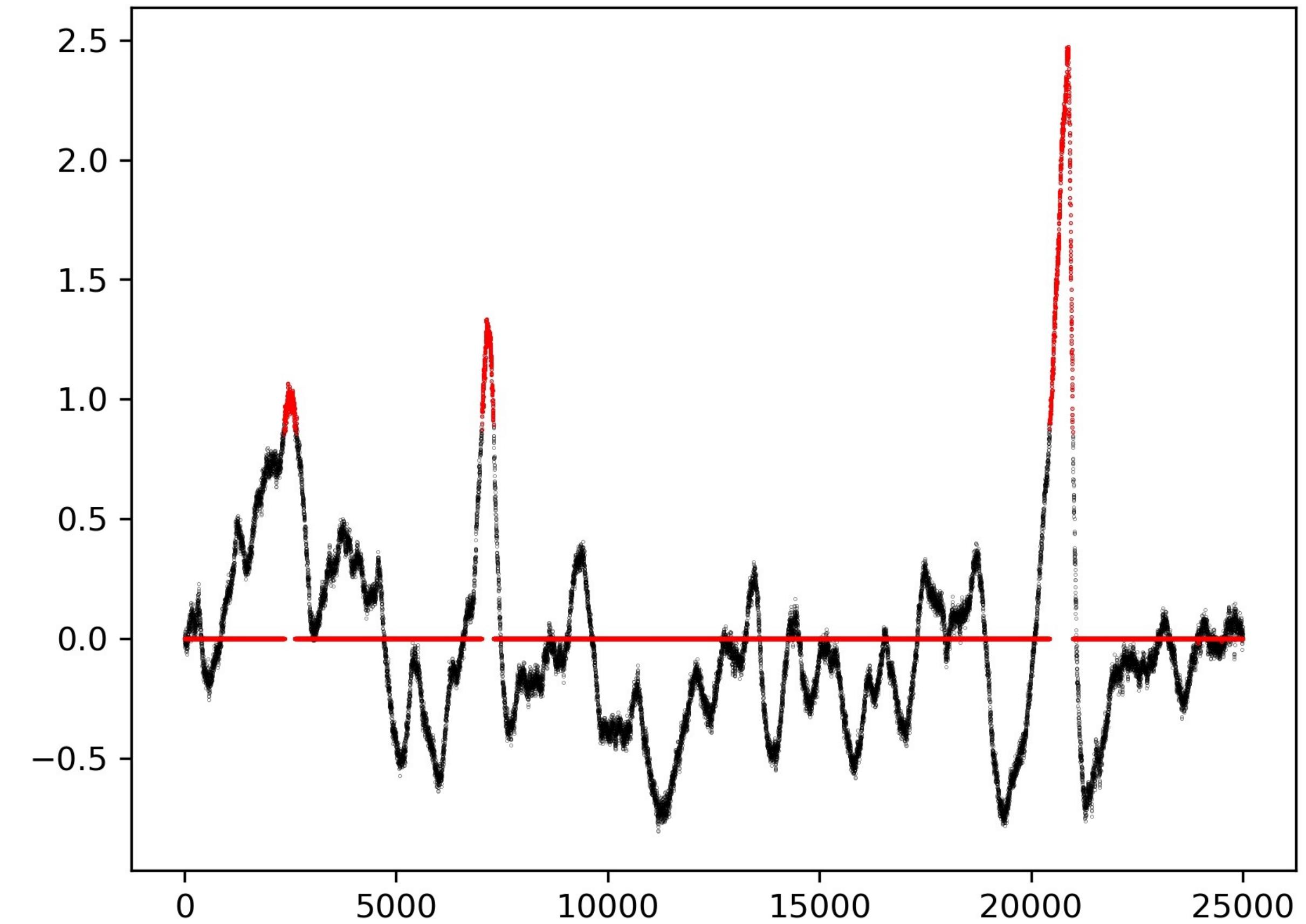
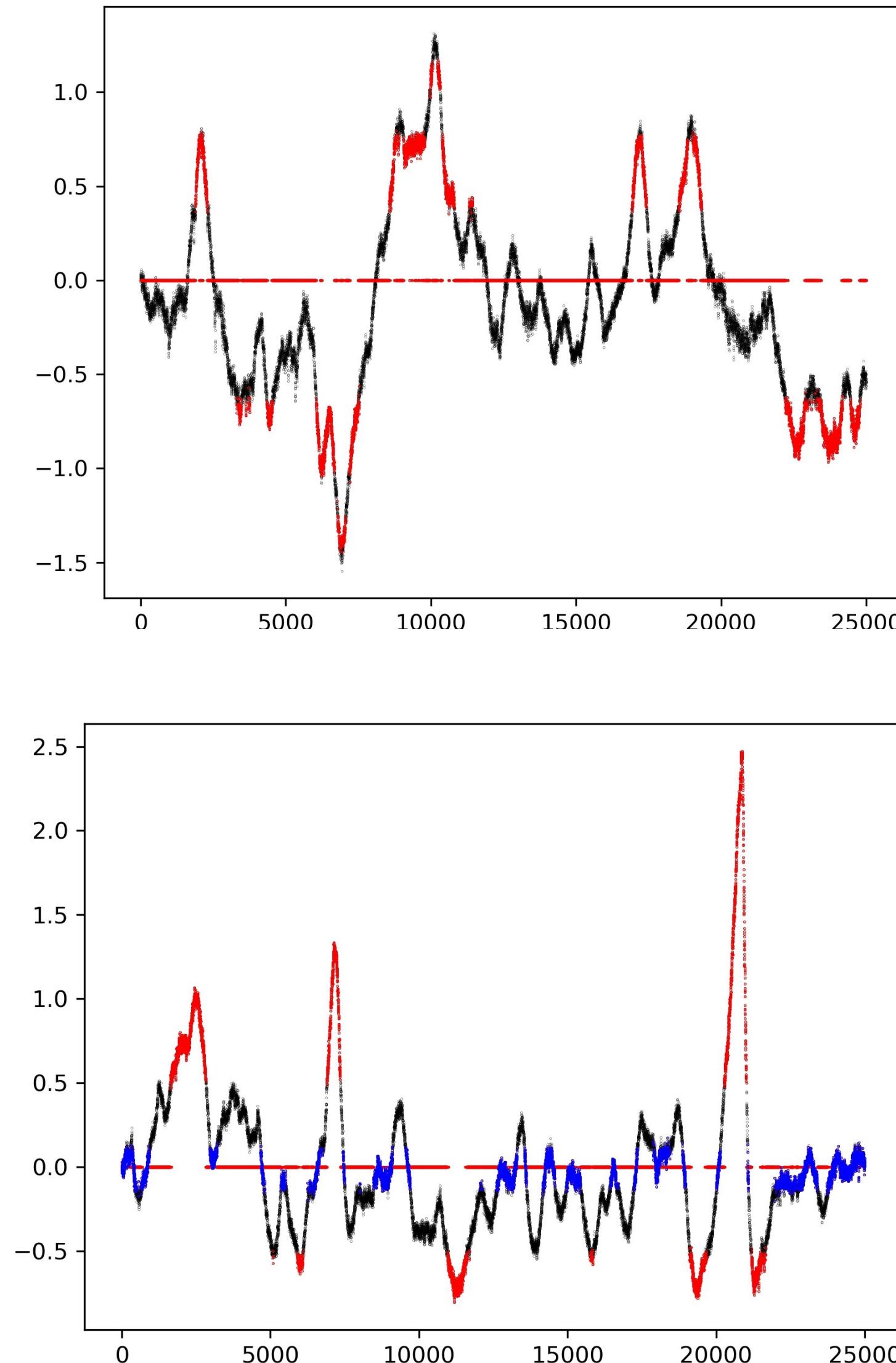


How many data points do we need to detect weak sparse signals ?

$$X \sim \varepsilon \mathcal{N}(\mu, 1) + (1 - \varepsilon) \mathcal{N}(0, 1)$$



Multi-unit activity data peaks detected with higher criticism



Multi-modal correlated time series and correlated noise

Note: HC is not sensitive to shuffling of coordinates, hence ignores the correlation

How to make the performance of higher criticism better in this situation ?

The covariance matrix of the data is a positive definite matrix: apply matrices using linear algebra to turn this matrix to identity (in sparse cases), or to decompose this matrix to the product of simpler matrices. Then use the HC on the result of application of such matrices to the coordinates.

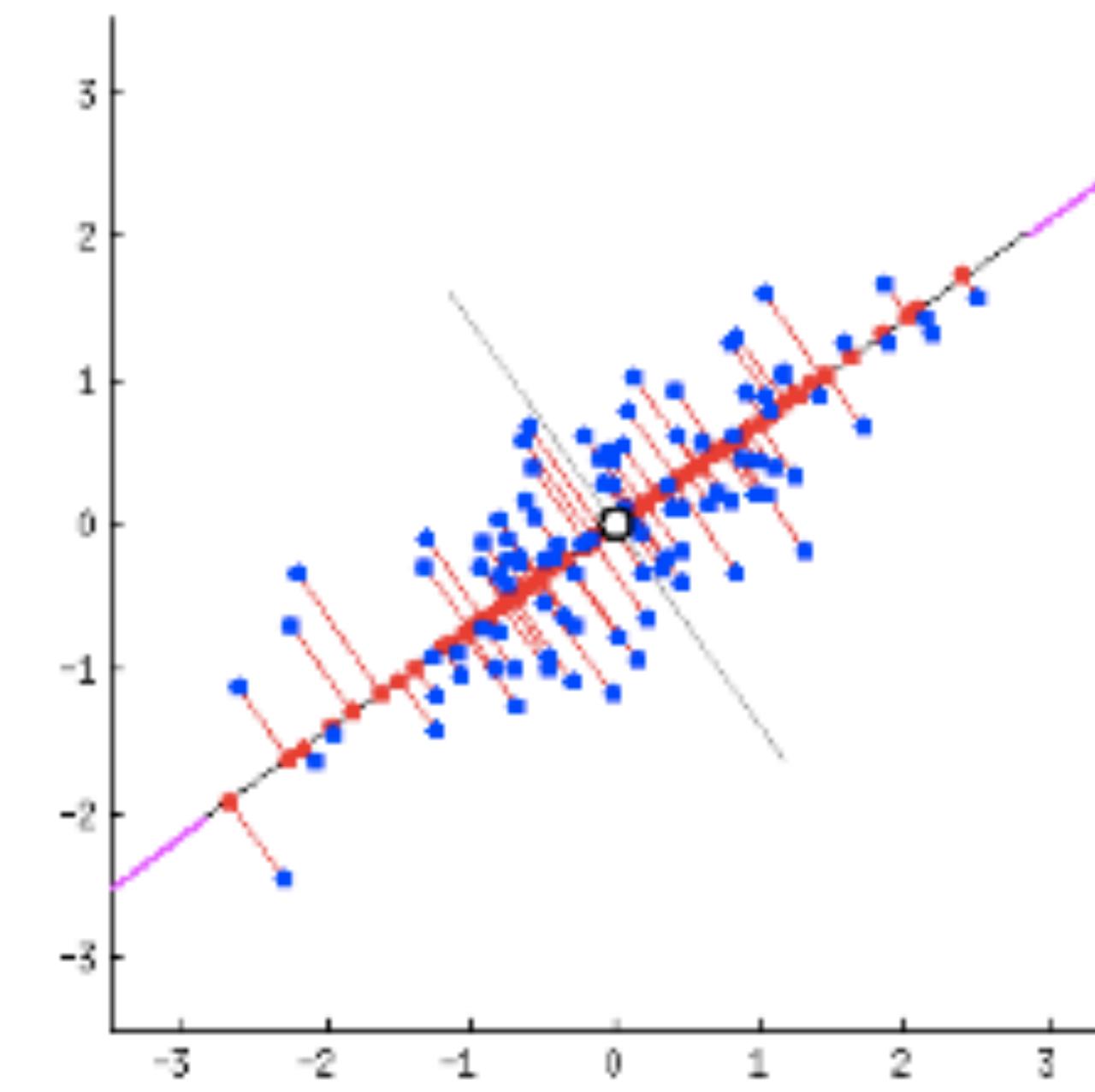
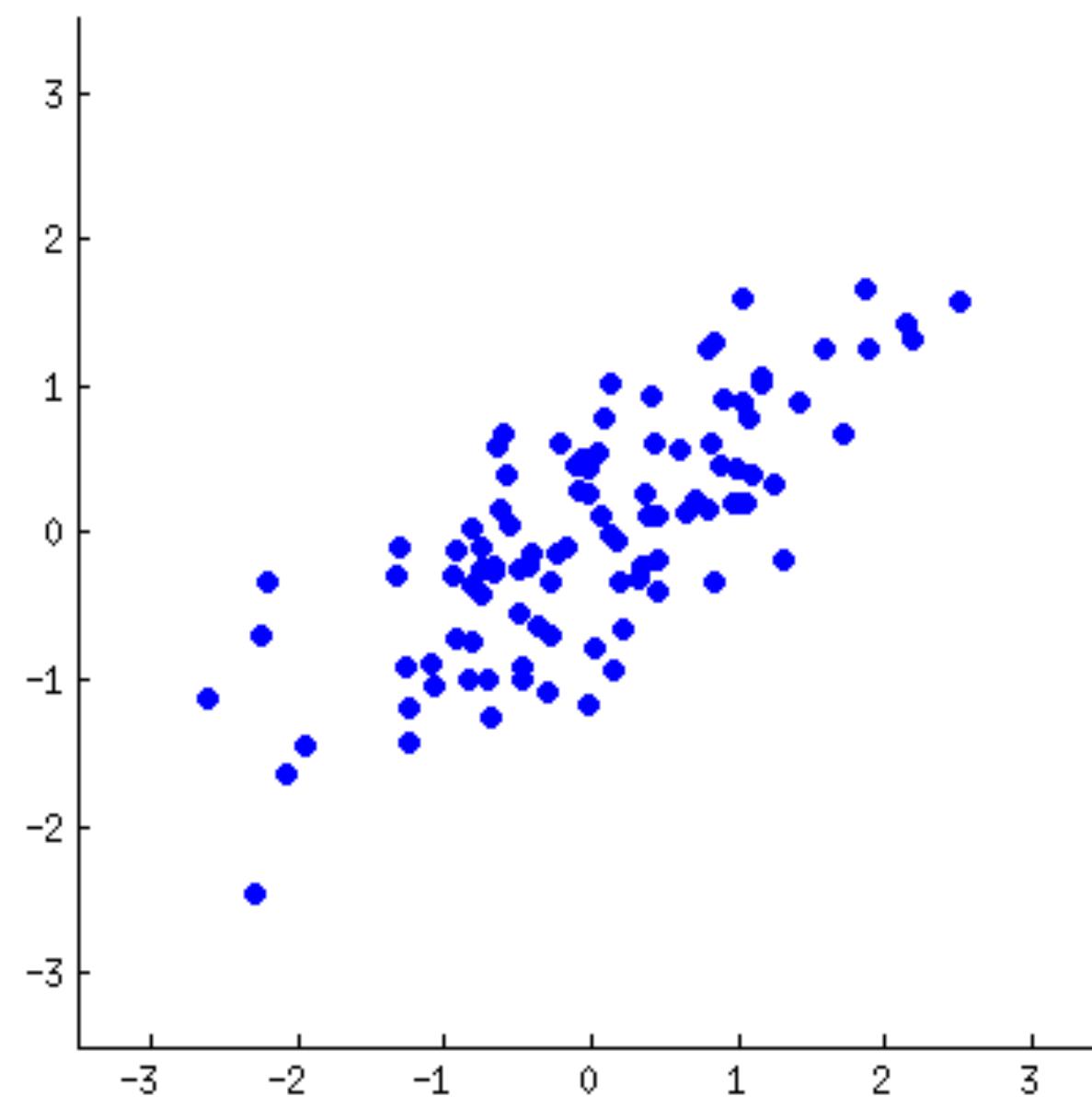
Analysis of movement and behavioral pattern recognition

We analyze the data from the smart housing facilities to understand stochastic properties of the animal behavior in their social life and to find relations between their behavior and biological states related to health, disease, gene modification, hormonal imbalances, etc.

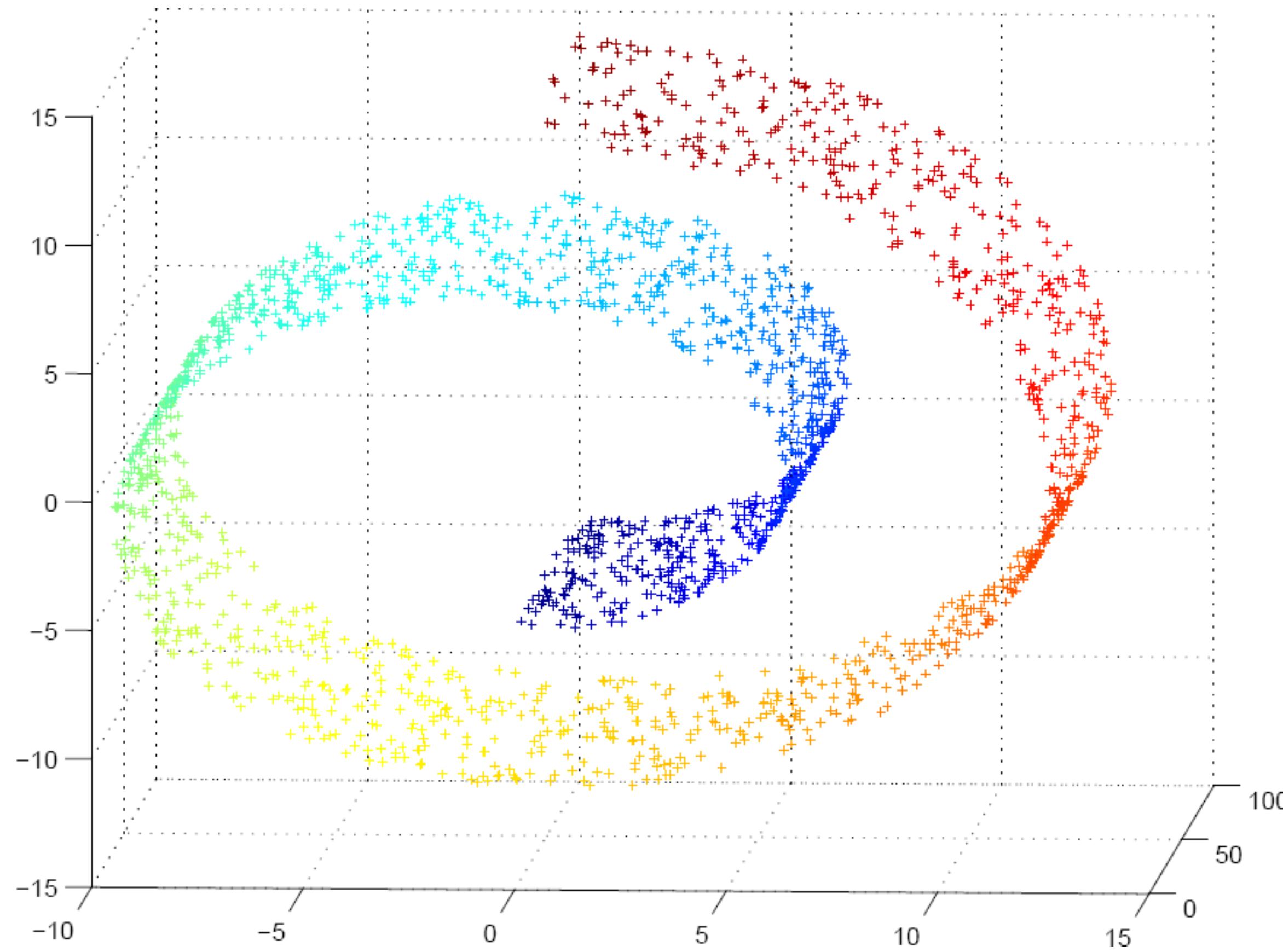
Geometric learning will be useful for encoding the correlations especially in the movements and for seeing patterns of social behavior in the geometric properties of the underlying manifold (non-linear structure).

Deep learning will be used to characterize animal behavior based on finger/hand motion, body statue, gait, facial mimics and complex interactive processes such as touching and grooming. Subsequently, we will study correlations between characteristic behavior and biological time-series (such as body temperature, cerebral activity, metabolism, and blood flow) as well as relations with the underlying neural and genetic mechanisms, and predict statistically disease-onset and dynamics.

Principal component analysis



Geometric learning: correlations preserved in the geometry of the data while dimension reduction is achieved



Geometric learning based on heat kernel and Laplacian

M. Belkin, P. Niyogi, Laplacian Eigenmaps and Spectral Techniques for Embedding and Clustering, in
“Advances in Neural Information Processing Systems 14” pp.585–591, MIT Press, 2001.

M. Belkin, P. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Comput.* 15 (6) (2003) 1373–1396.

M. Belkin, P. Niyogi, Towards a theoretical foundation for Laplacian-based manifold methods, *J. Comput. System Sci.* 74 (2008), no. 8, 1289–1308.

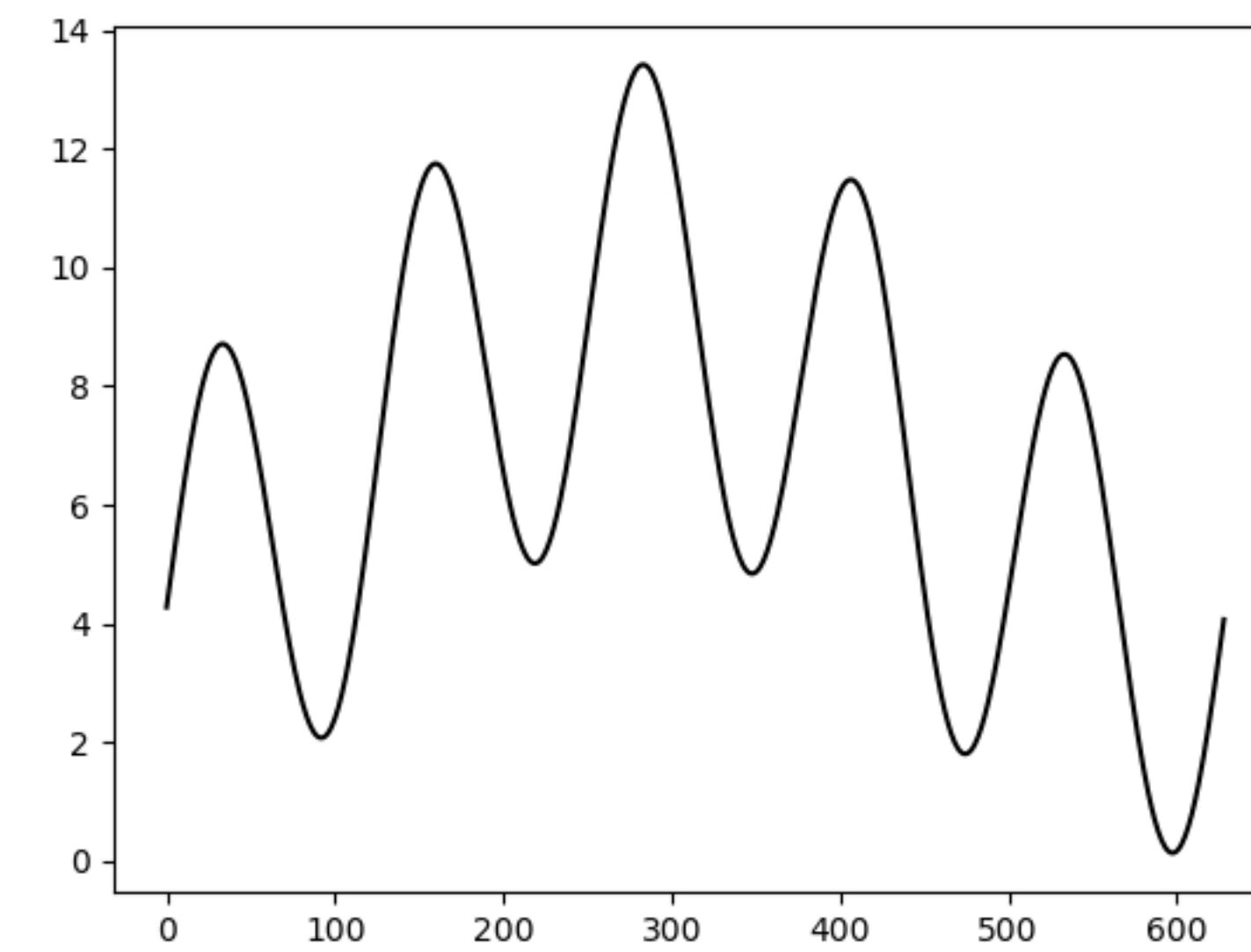
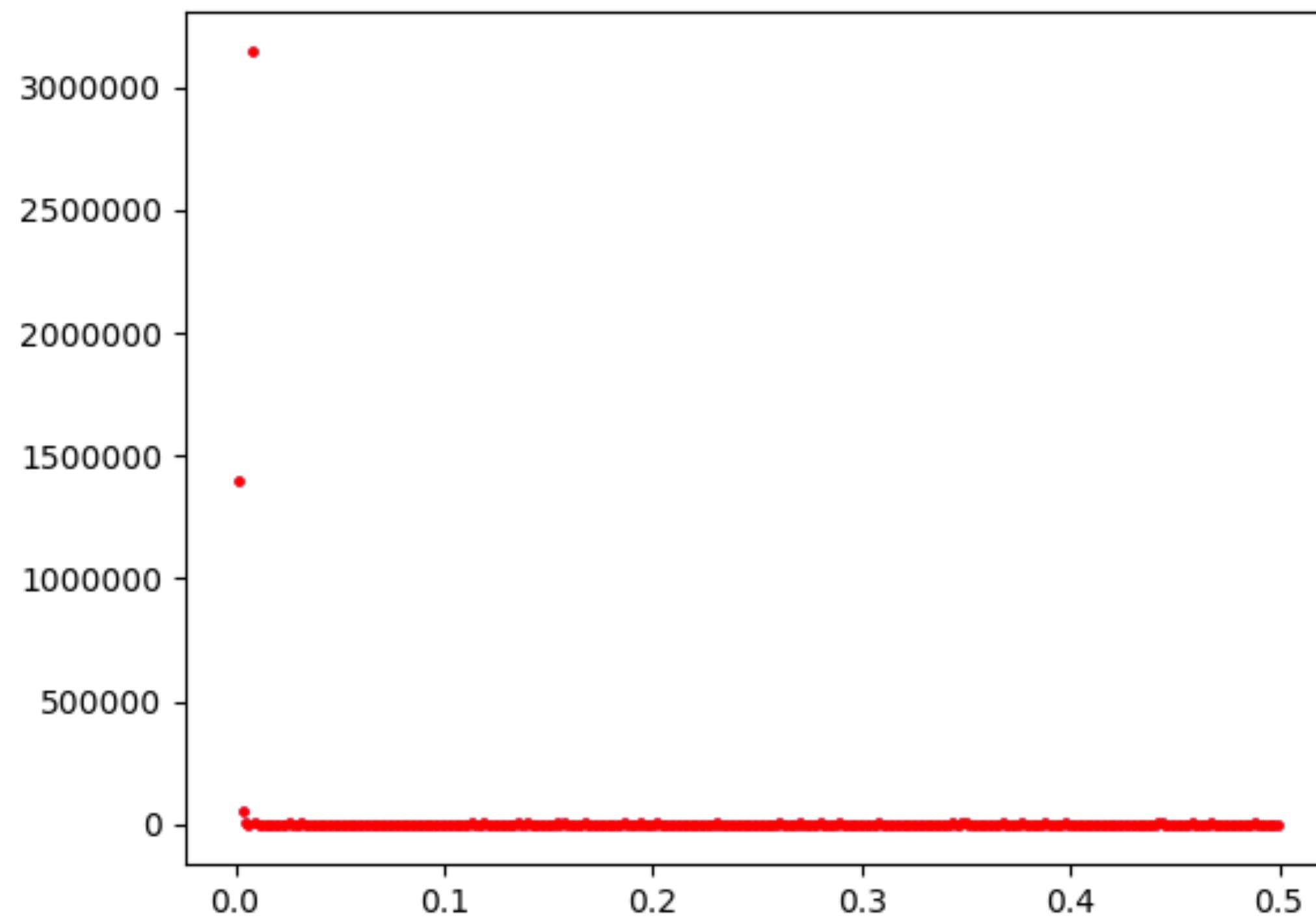
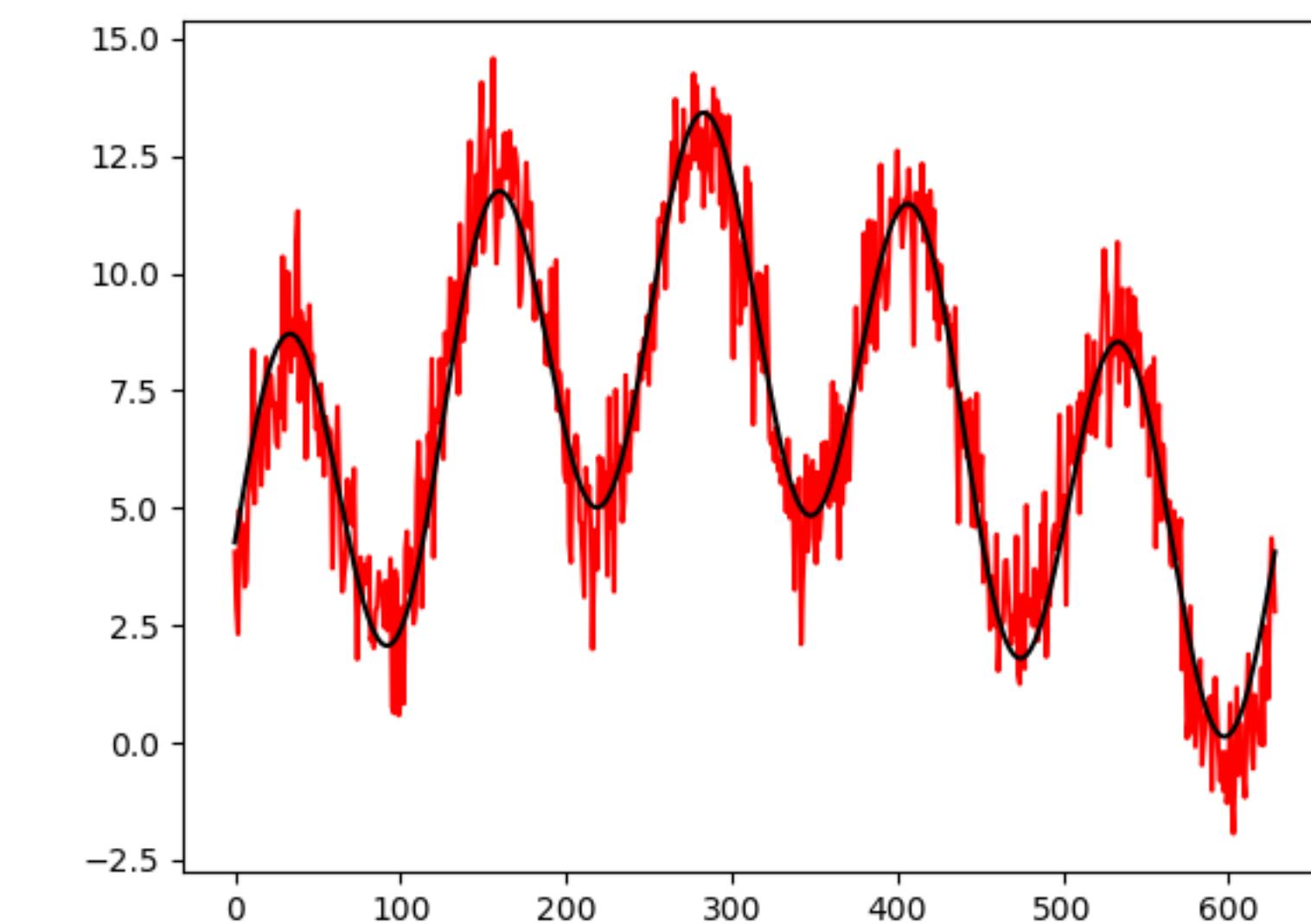
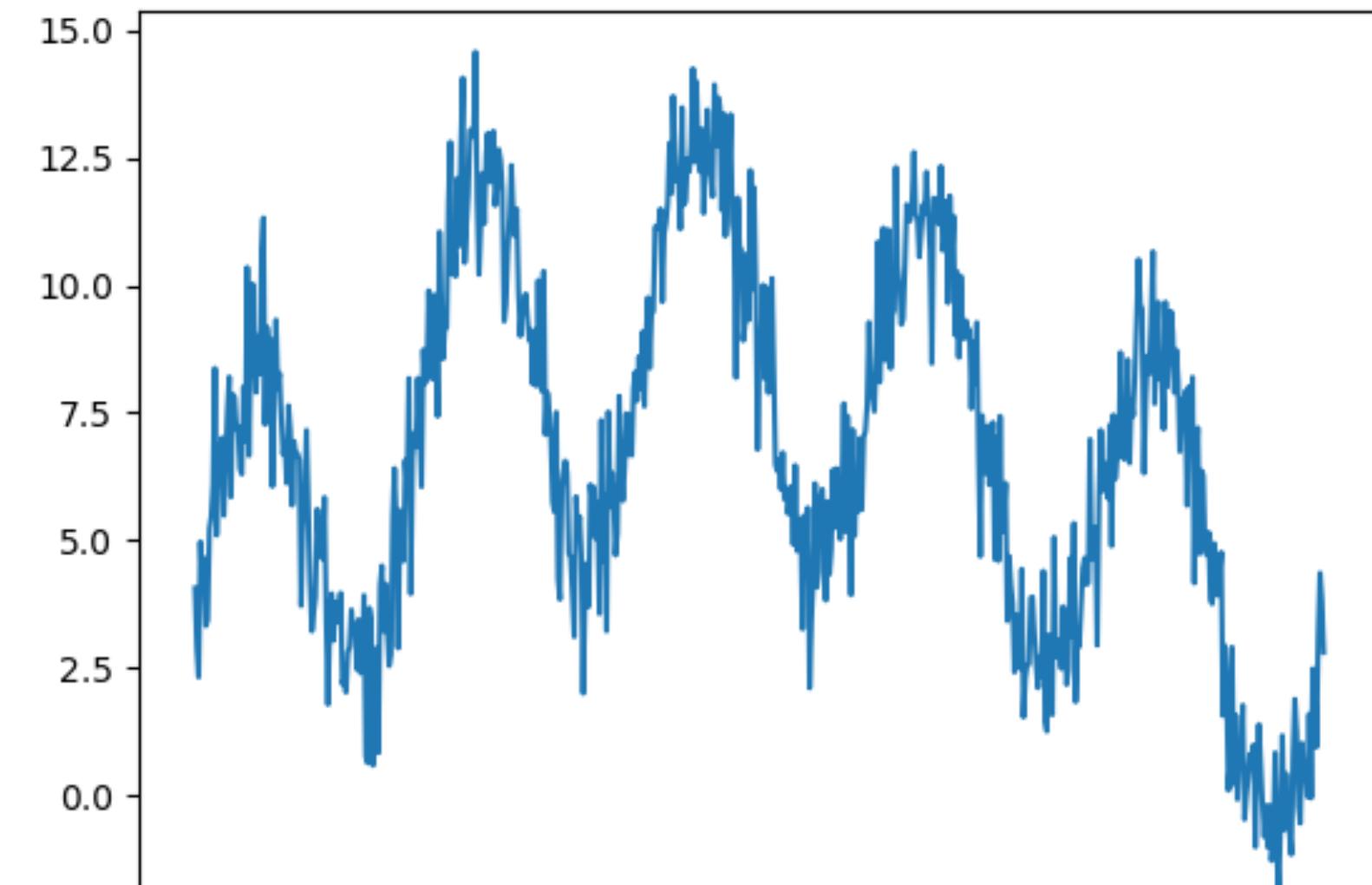
Instantaneous time frequency analysis

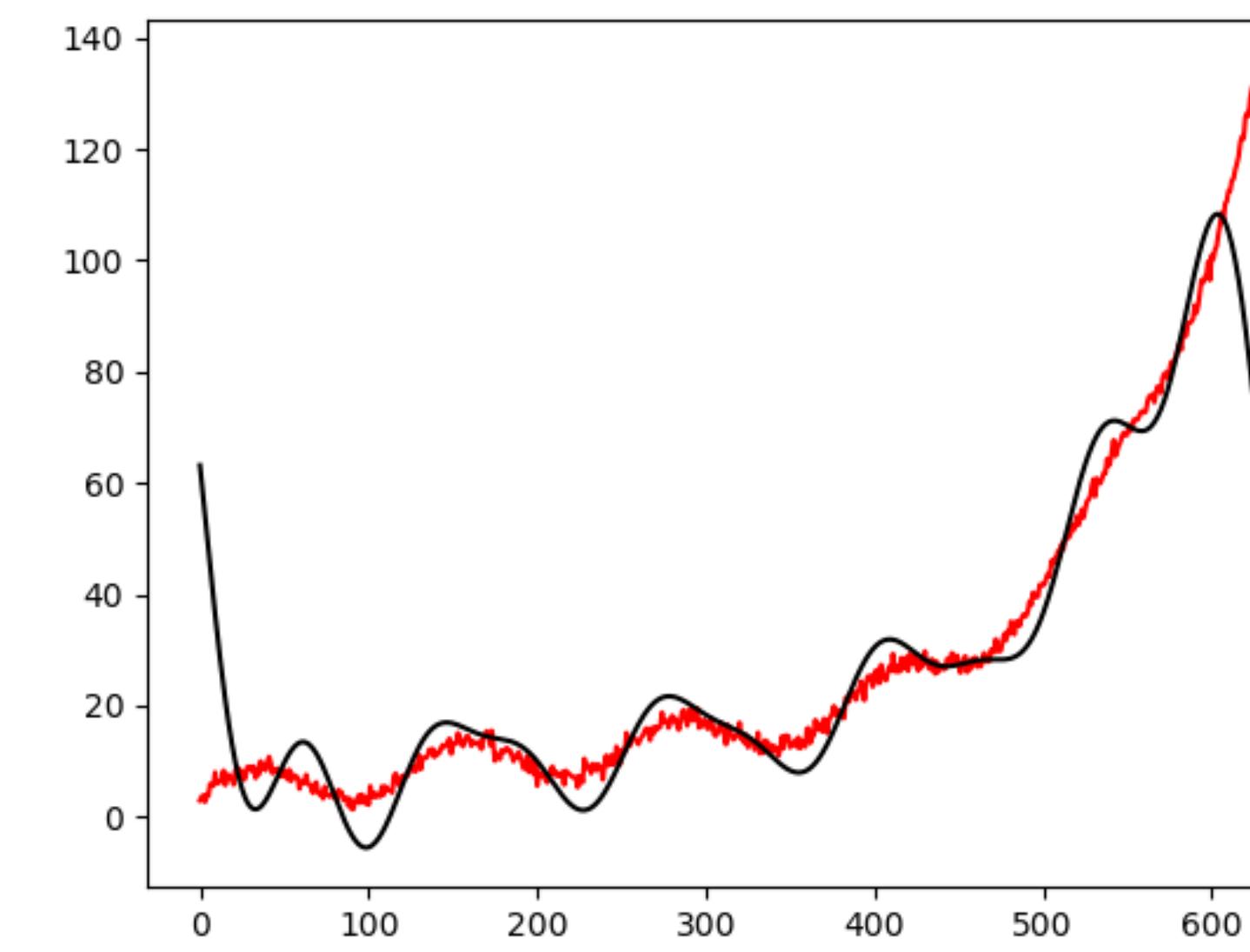
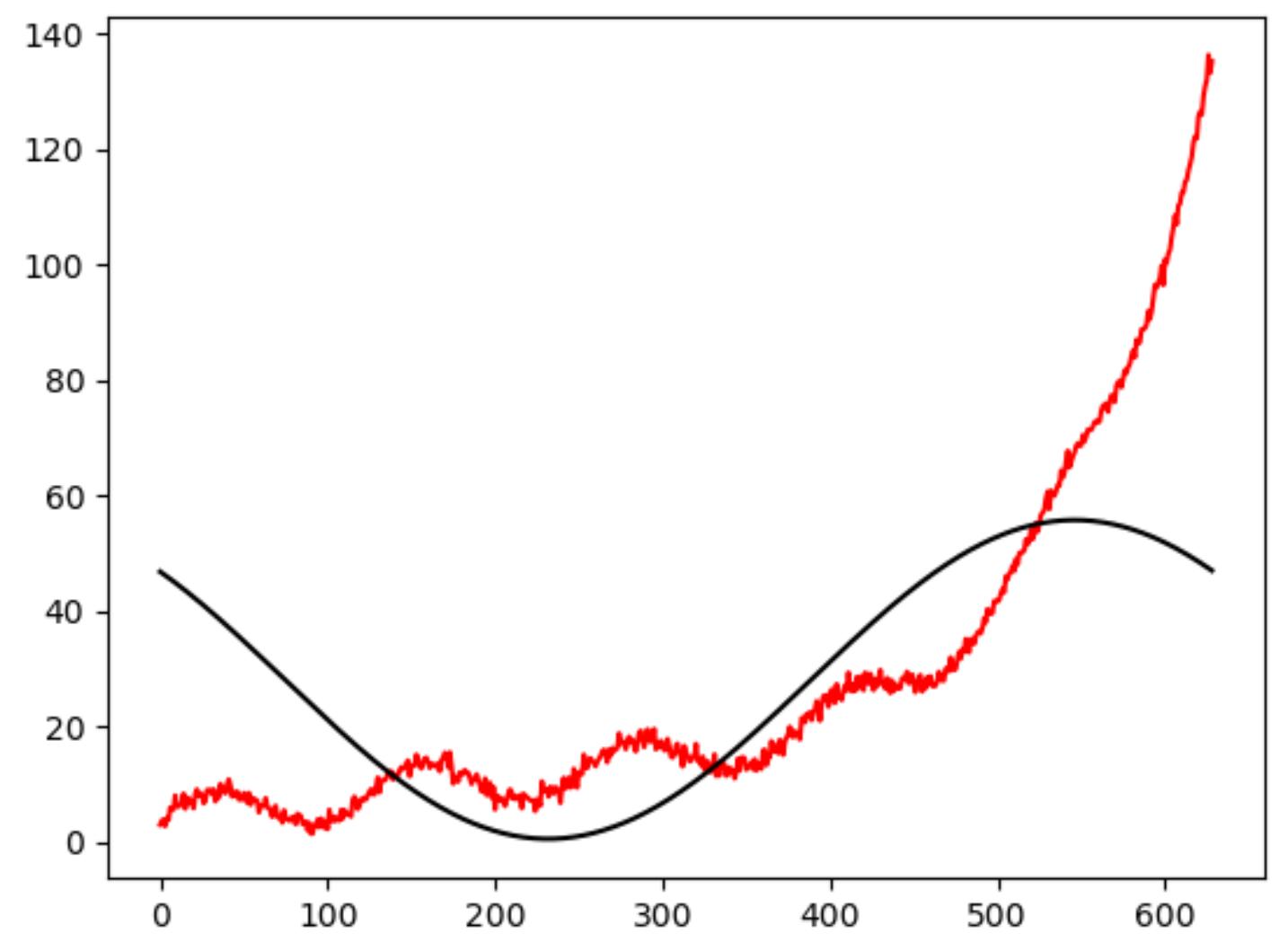
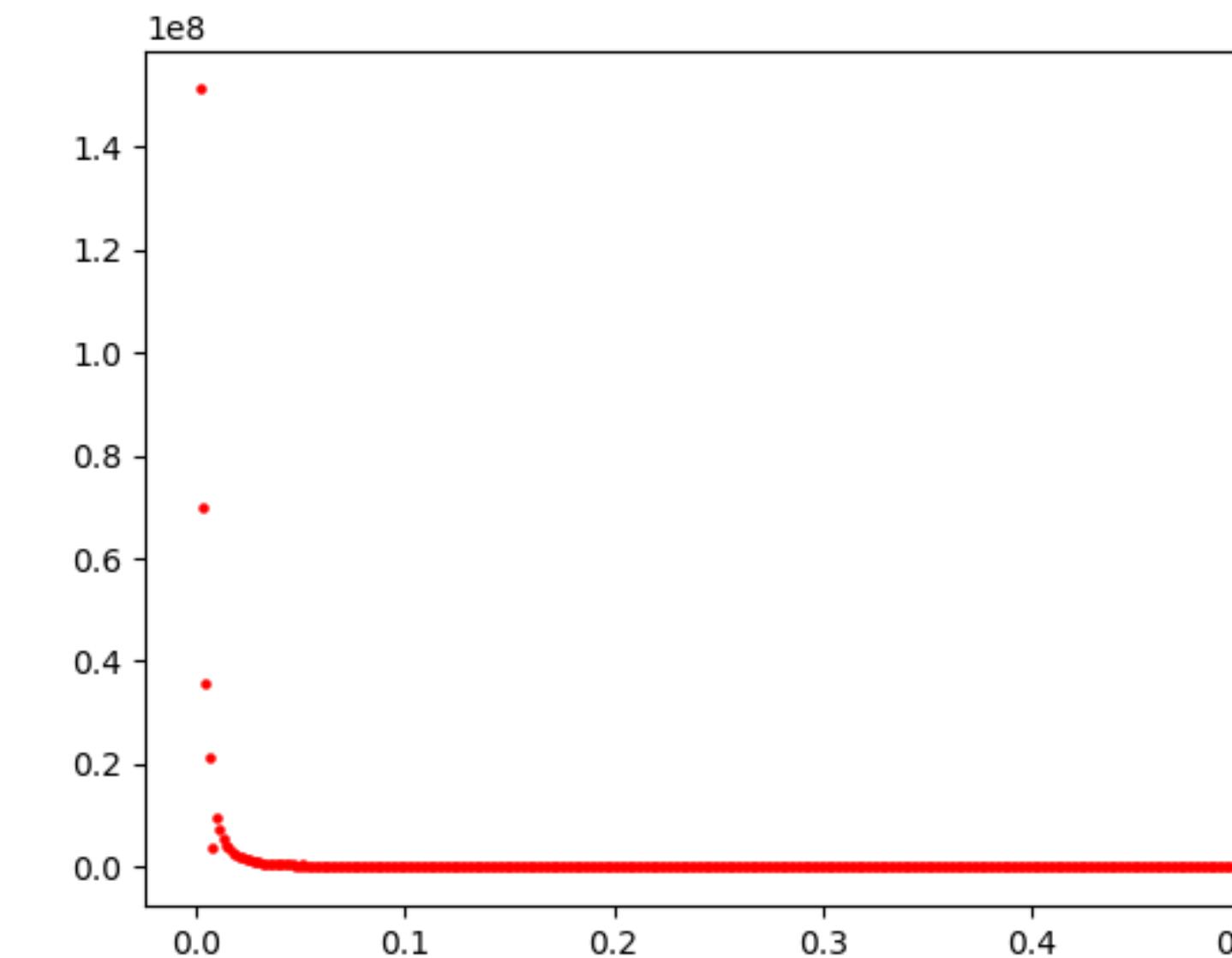
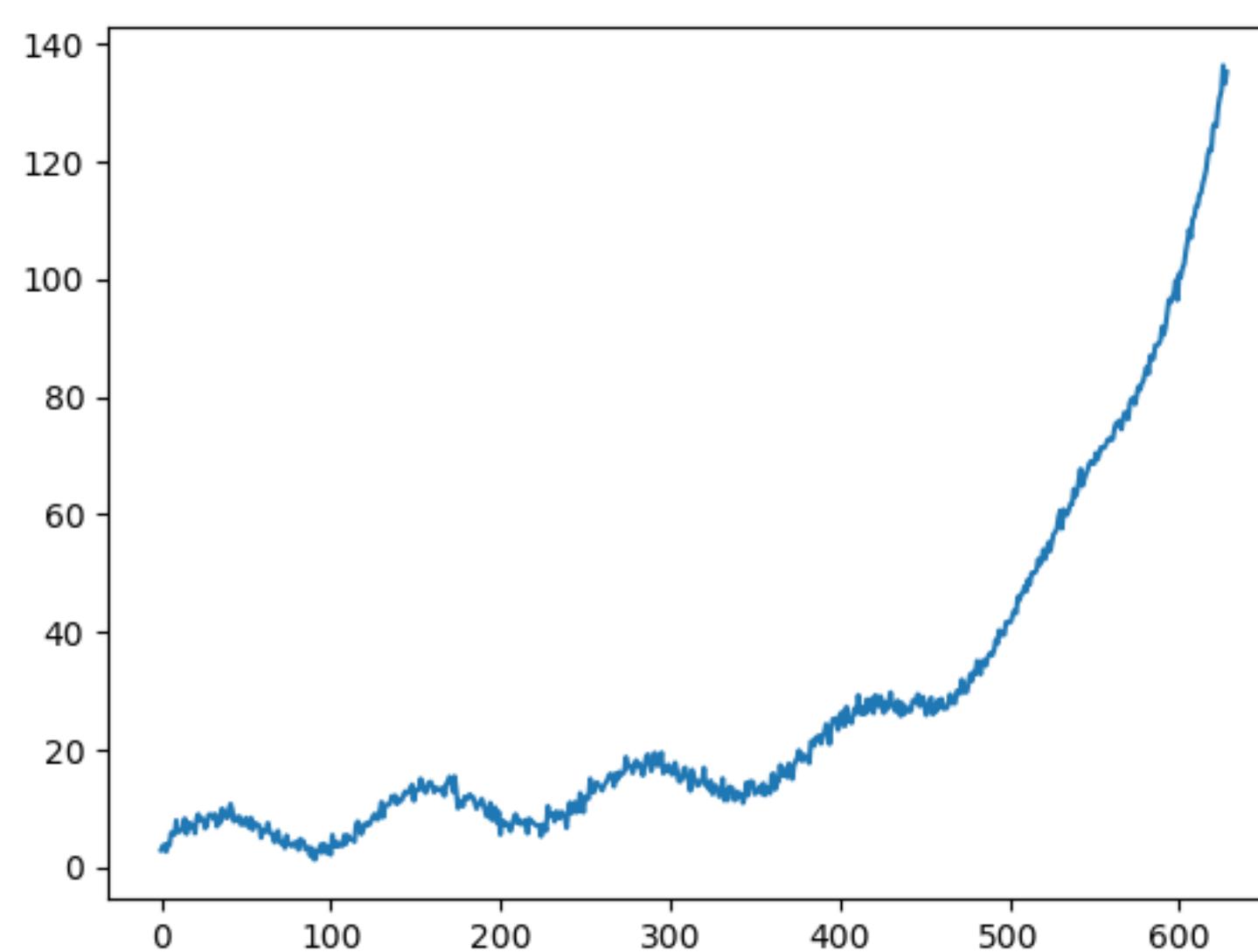
Fourier analysis:

$$x(t) = \sum_j a_j e^{i\omega_j t} = \sum_j a_j (\cos(\omega_j t) + i \sin(\omega_j t))$$

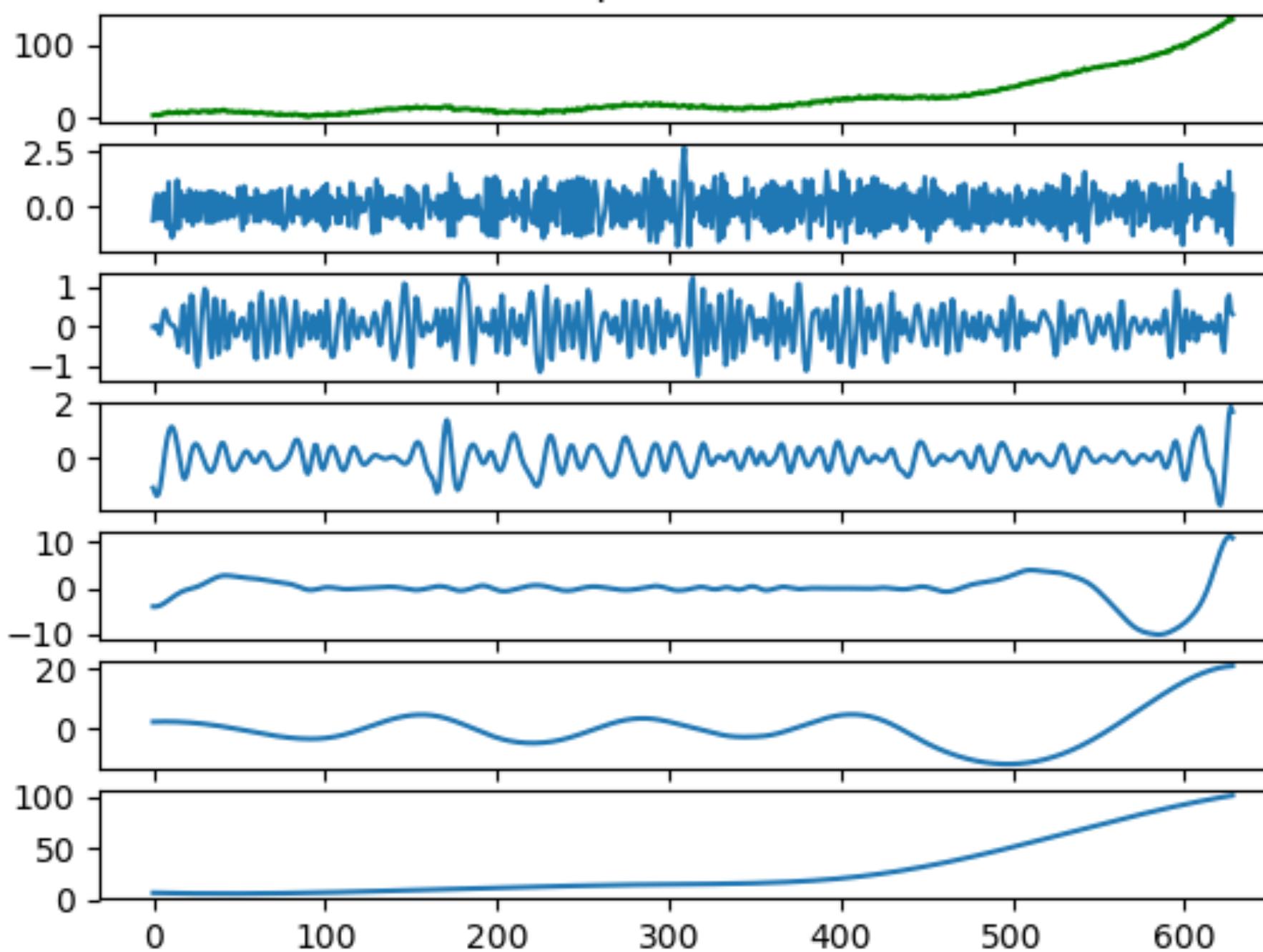
Hilbert spectral analysis: $x^*(t) = P.V. \int \frac{x(\tau)}{t - \tau} d\tau$

$$x(t) = \sum_j a_j(t) e^{i\omega_j(t)} = \sum_j a_j(t) (\cos(\omega_j(t)) + i \sin(\omega_j(t)))$$

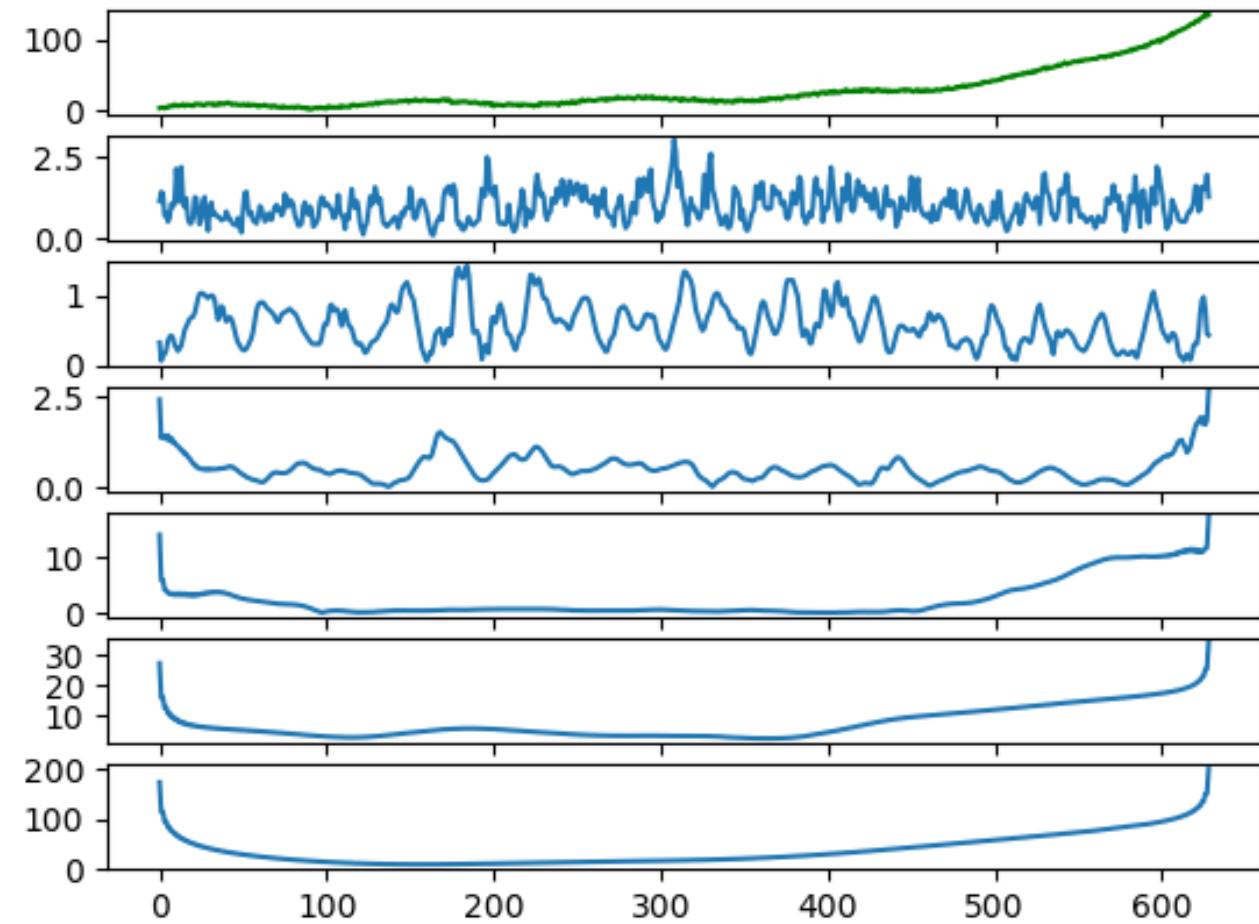




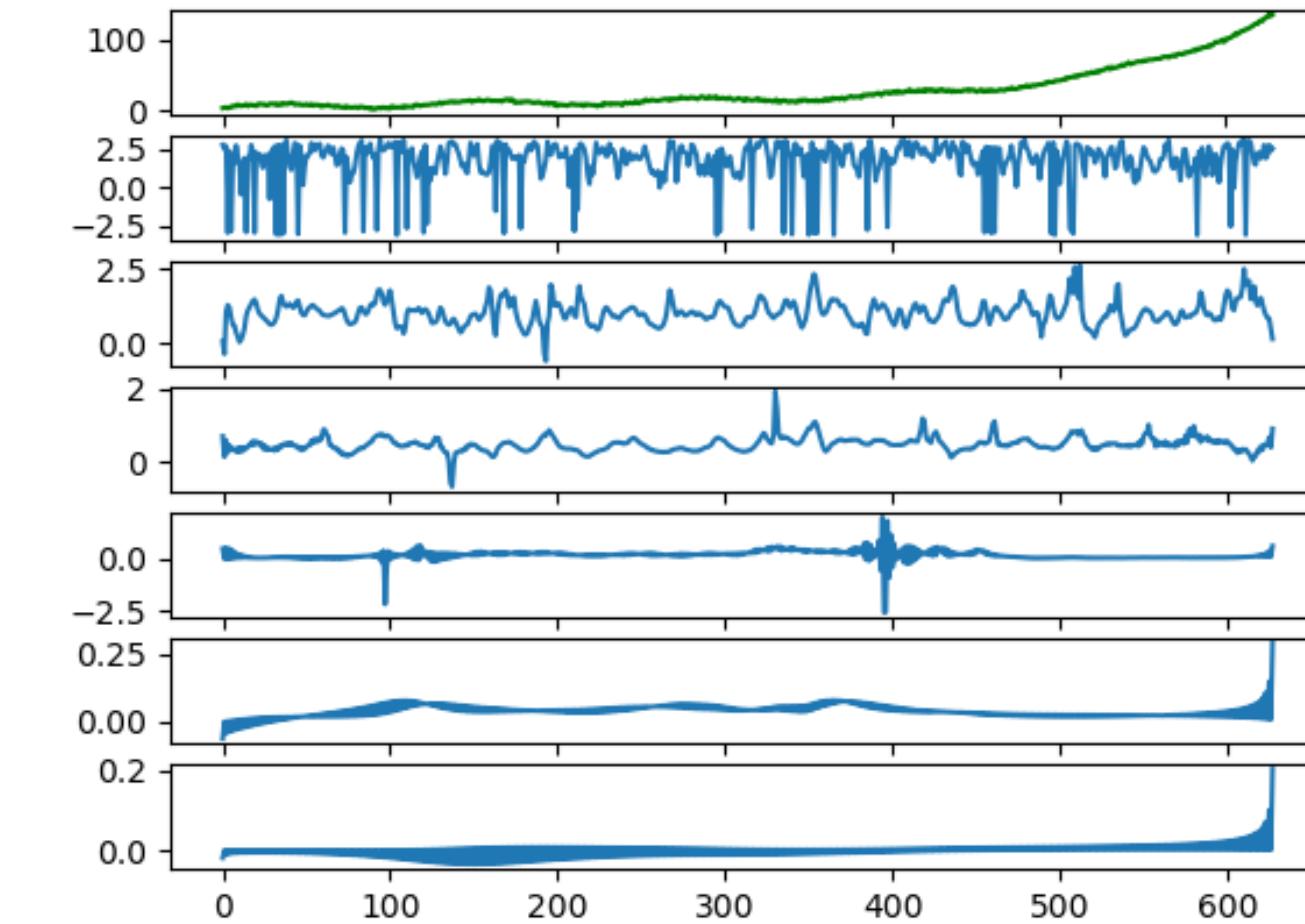
The input and its IMFs

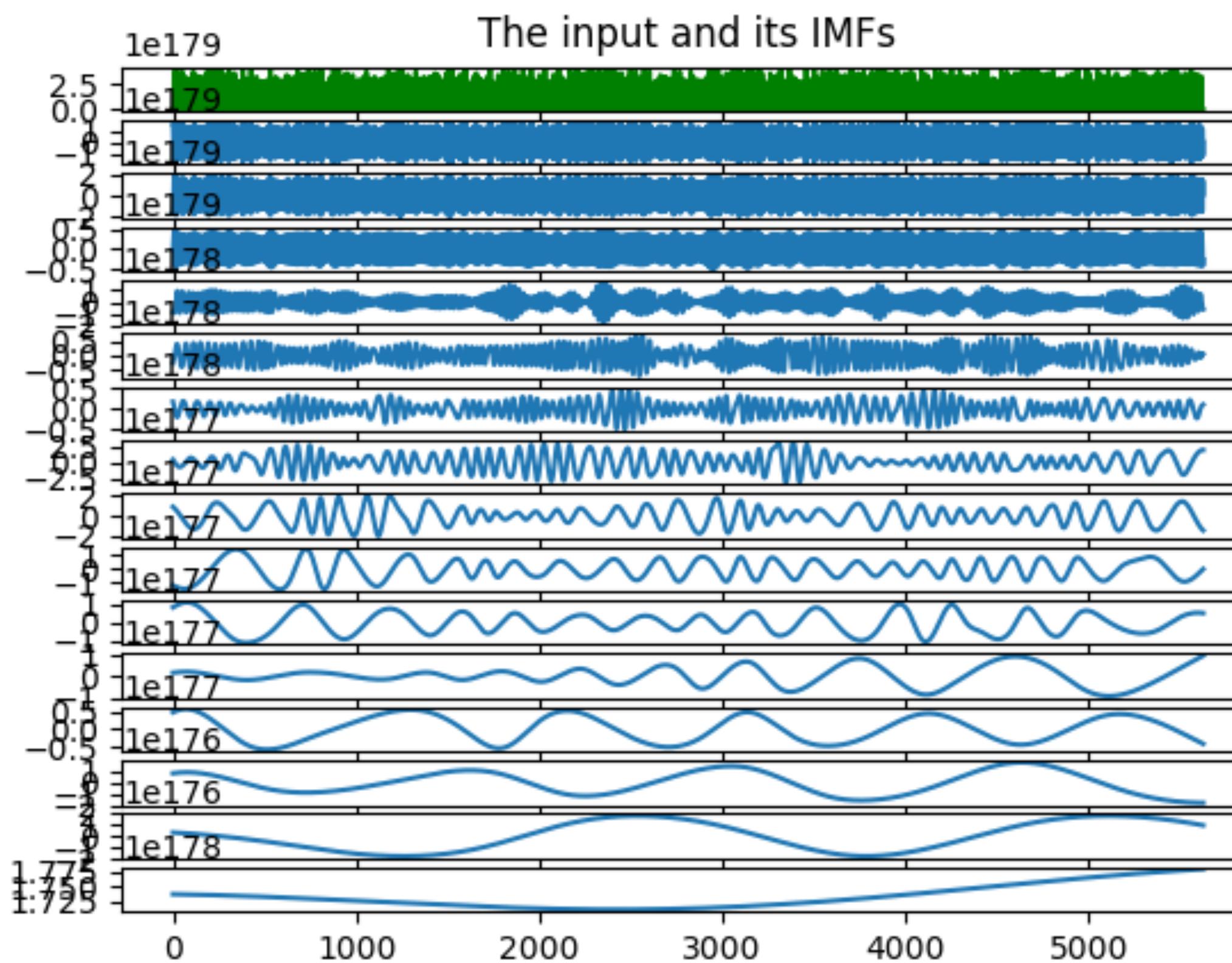
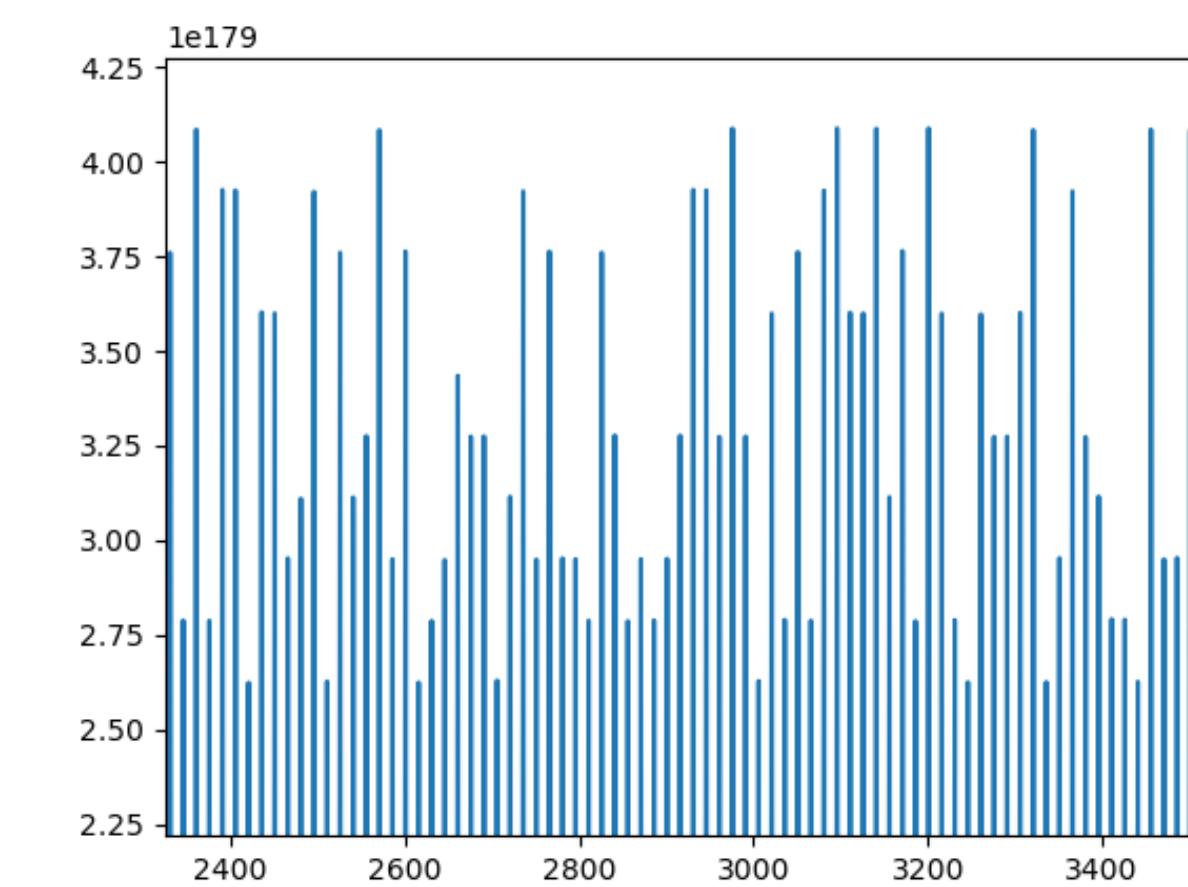
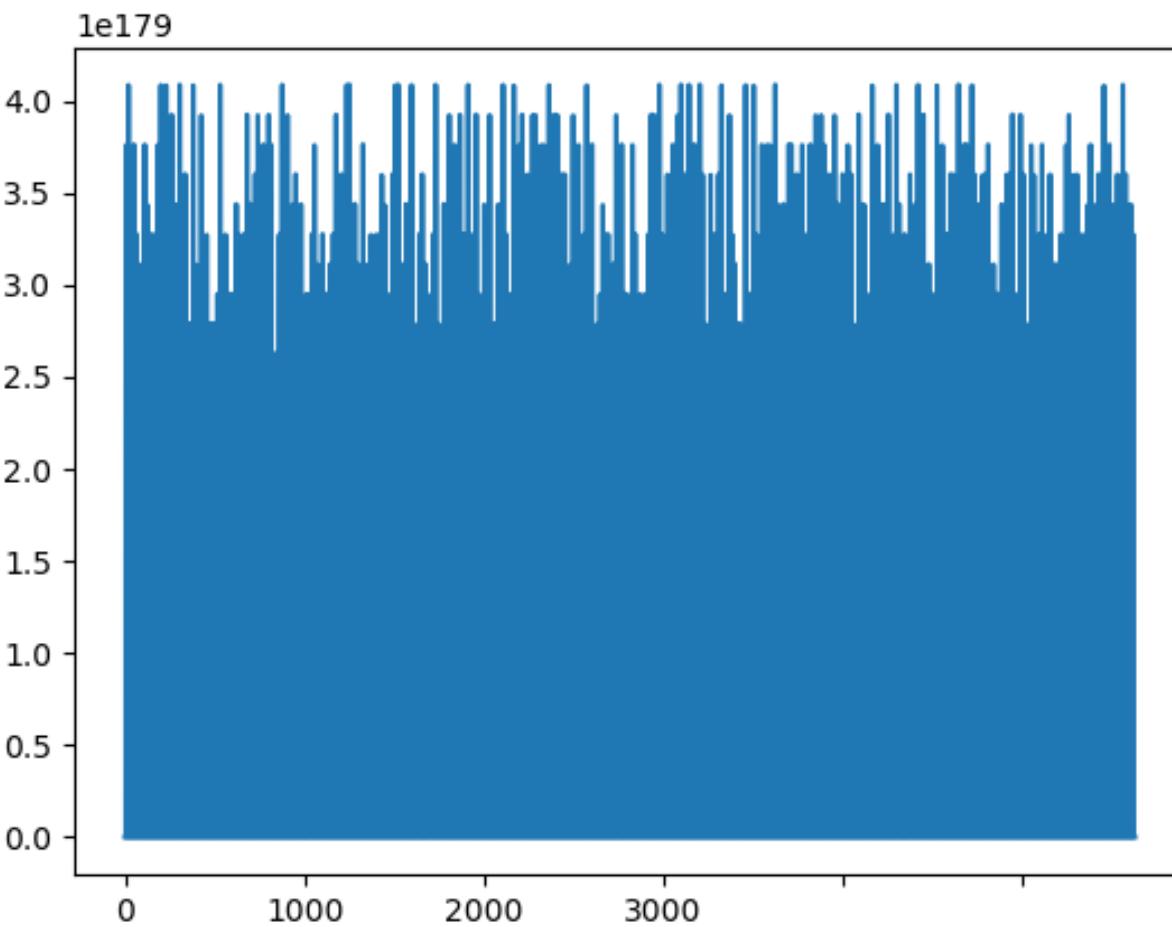


The input and instantaneous AMPLITUDES of its IMFs

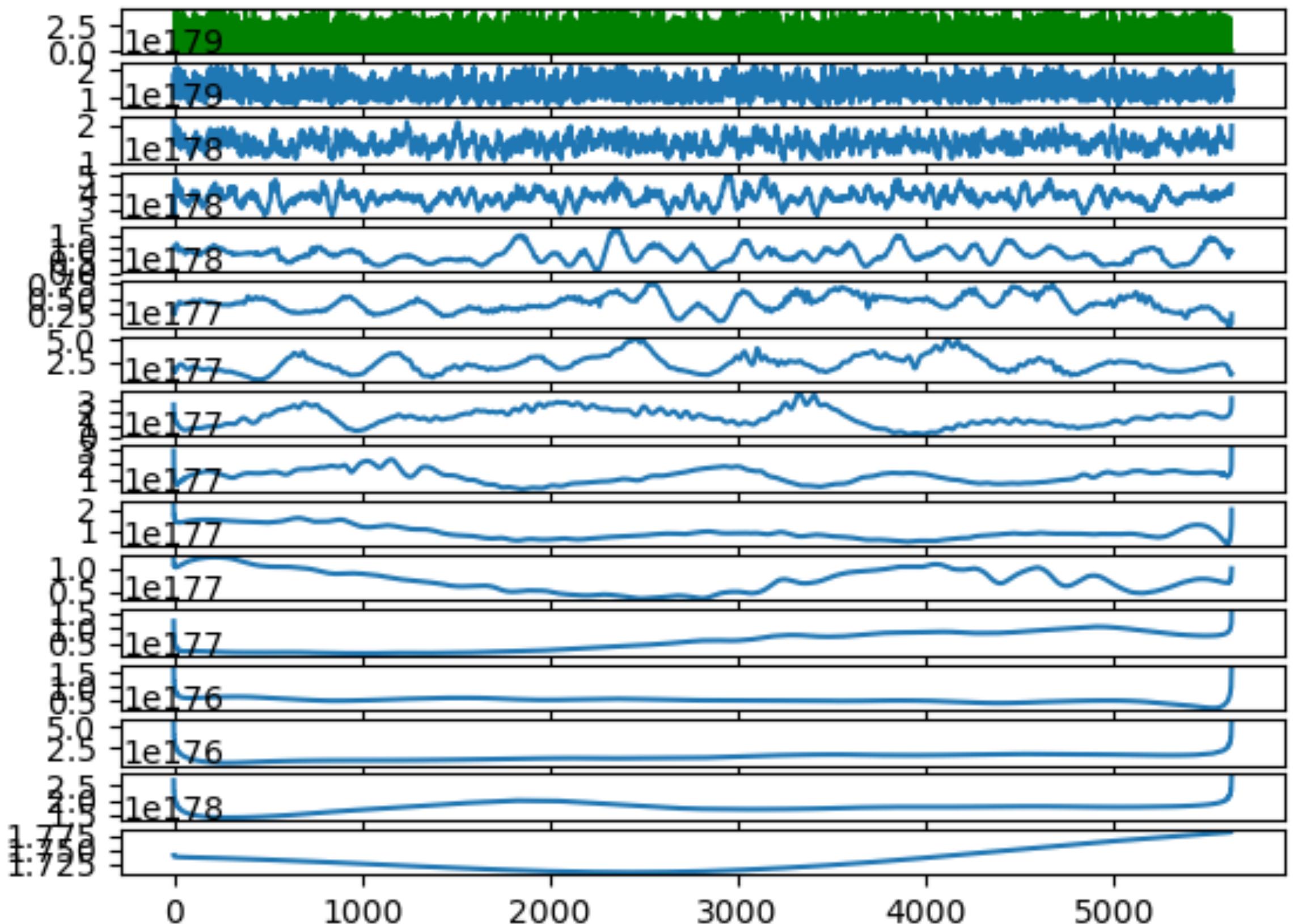


The input and instantaneous FREQUENCIES of its IMFs





The input and instantaneous AMPLITUDES of its IMFs



The input and instantaneous FREQUENCIES of its IMFs

