

Global Optimization

complicated multivariate potential surface
local and global minima

earlier methods go to local ones, here: trials to find global one

Many methods: random (Monte Carlo) search, stochastic gradient methods, simulated annealing, genetic algorithm...

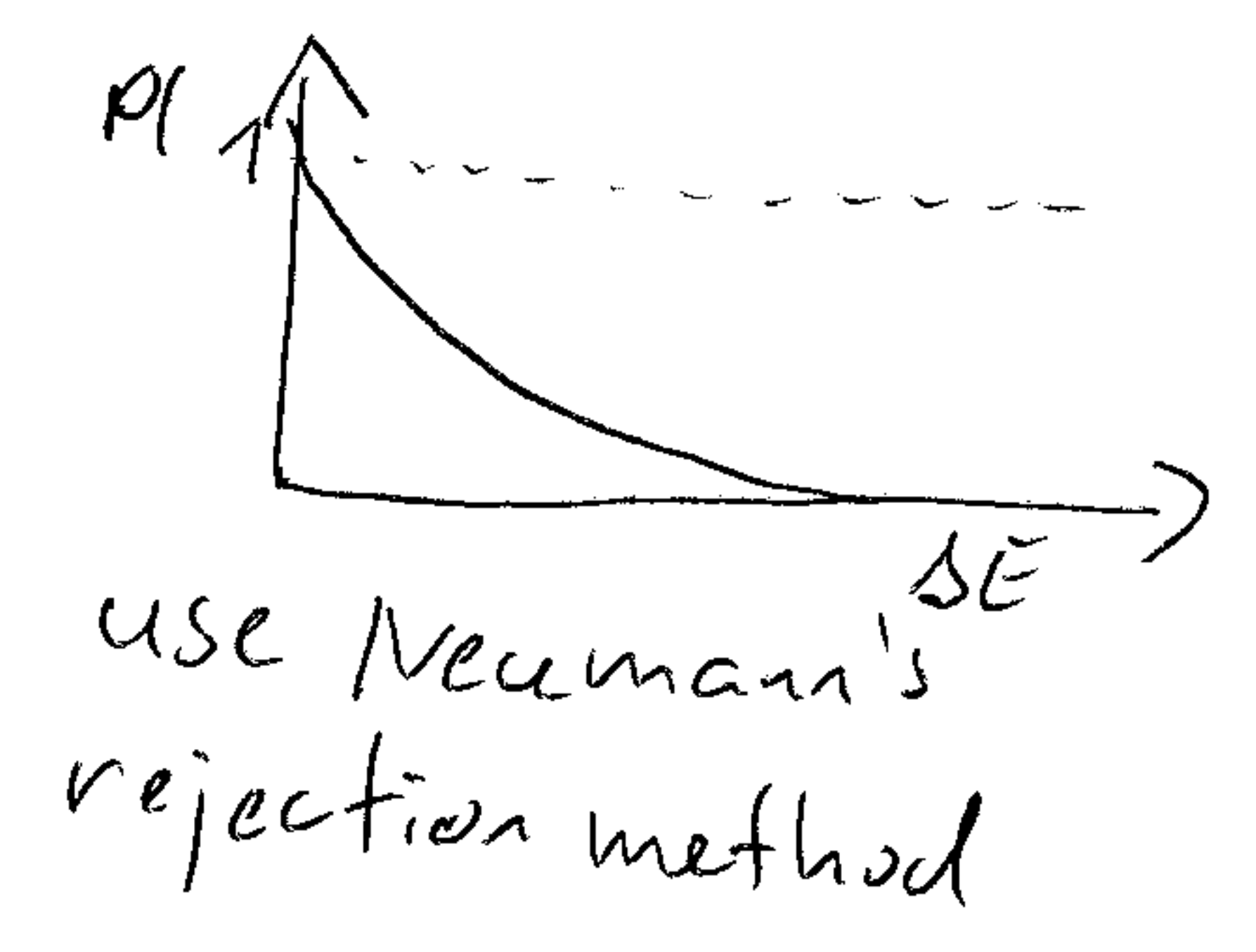
Simulated annealing

- define energy-like value $E(x)$ x can be independent variable, parameter...
- suppose $P(E) \sim \exp(-E/kT)$ \approx Boltzmann distributed
- define x_0 , calculate $E(x_0)$
- define T_0 ($\hat{=}$ temperature)

\rightarrow create new x (close to x_0), randomly, calculate $E(x)$ and $\Delta E = E(x) - E(x_0)$

if $\Delta E \leq 0$ x is the new accepted minimum point

$\Delta E > 0$ x is accepted with $P(\exp(-\Delta E/kT))$ probability
 x_0 remains, if x is rejected



- repeat
- reduce sometimes $T \rightarrow 0$
- $T \approx 0$ converged

\approx Monte Carlo simulation of condensed matter is very similar

Global optimization

GO-2

Genetic algorithm (evolutionary algorithm)

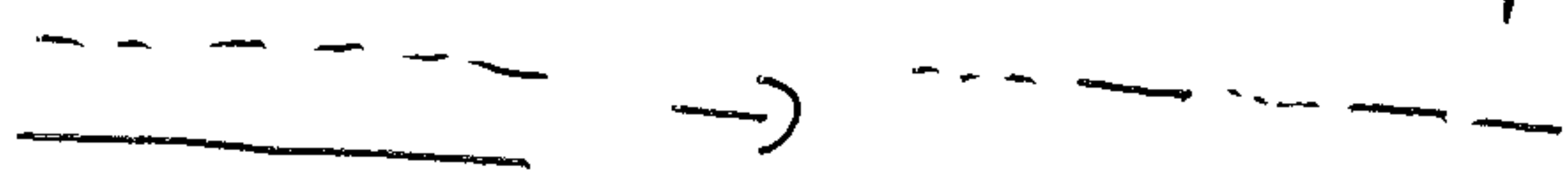
Previous min-max search: one starting point (except simplex)
(one of biologically inspired methods)

Here: more starting points \rightarrow result is also more points (simplex)

Nomenclature: - points at a given iterative sets: population
- function to optimize: fitness (large is good)
- variables \rightarrow stored binary, like genetic information

Next generation (population) is created by:

- reproduction: most fit points survive to the next
- crossing: new member is created by interchanging variables of two ones from the previous generation



- mutation: small changes in the values

Iteration stops, when fitness (average) does not ~~change~~ become better.

Result: ~~set of~~ best collection of elements with best fitness or elements of the best population

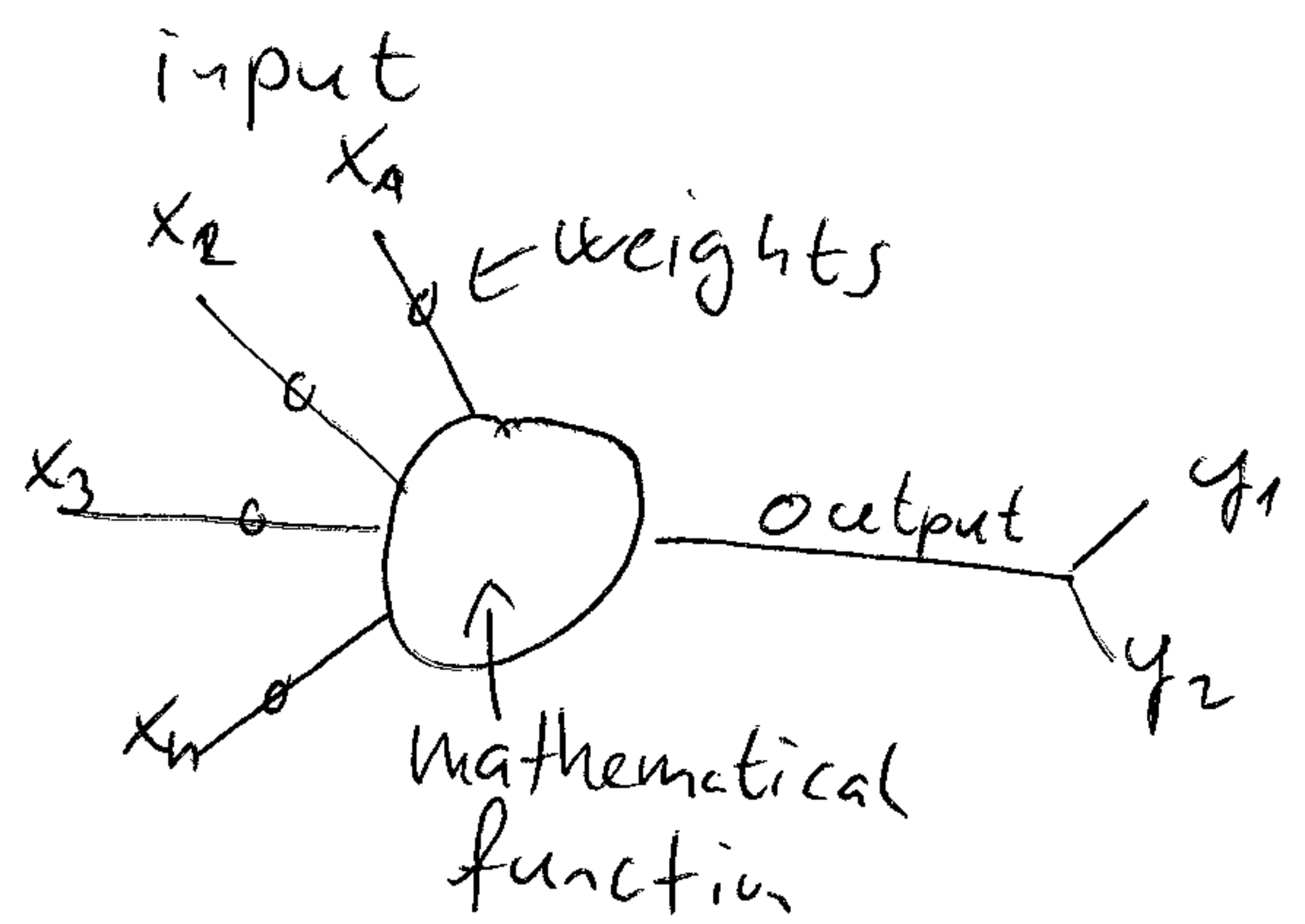
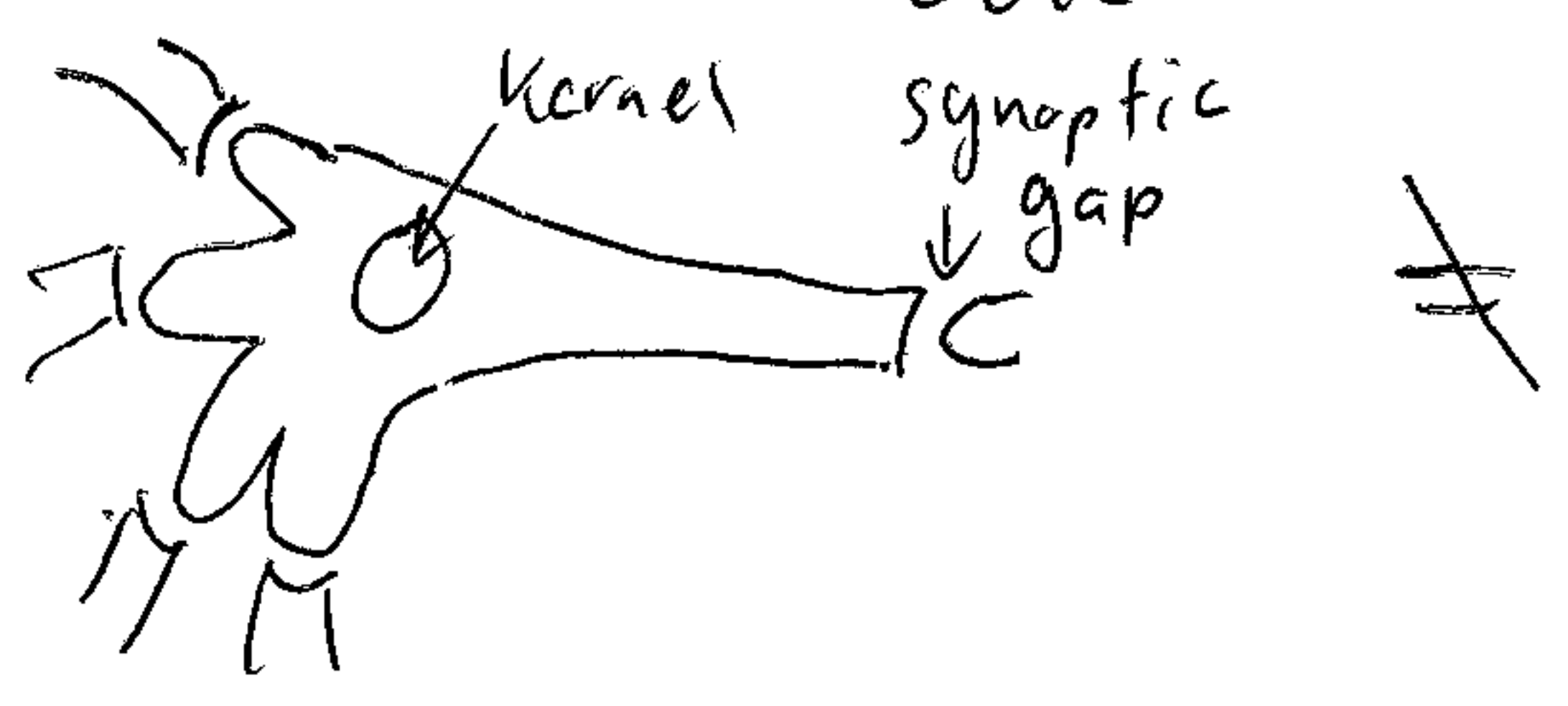
- Examples:

- minimization, parameter estimation
- function approximation
- program writing

Neural networks

NN-1

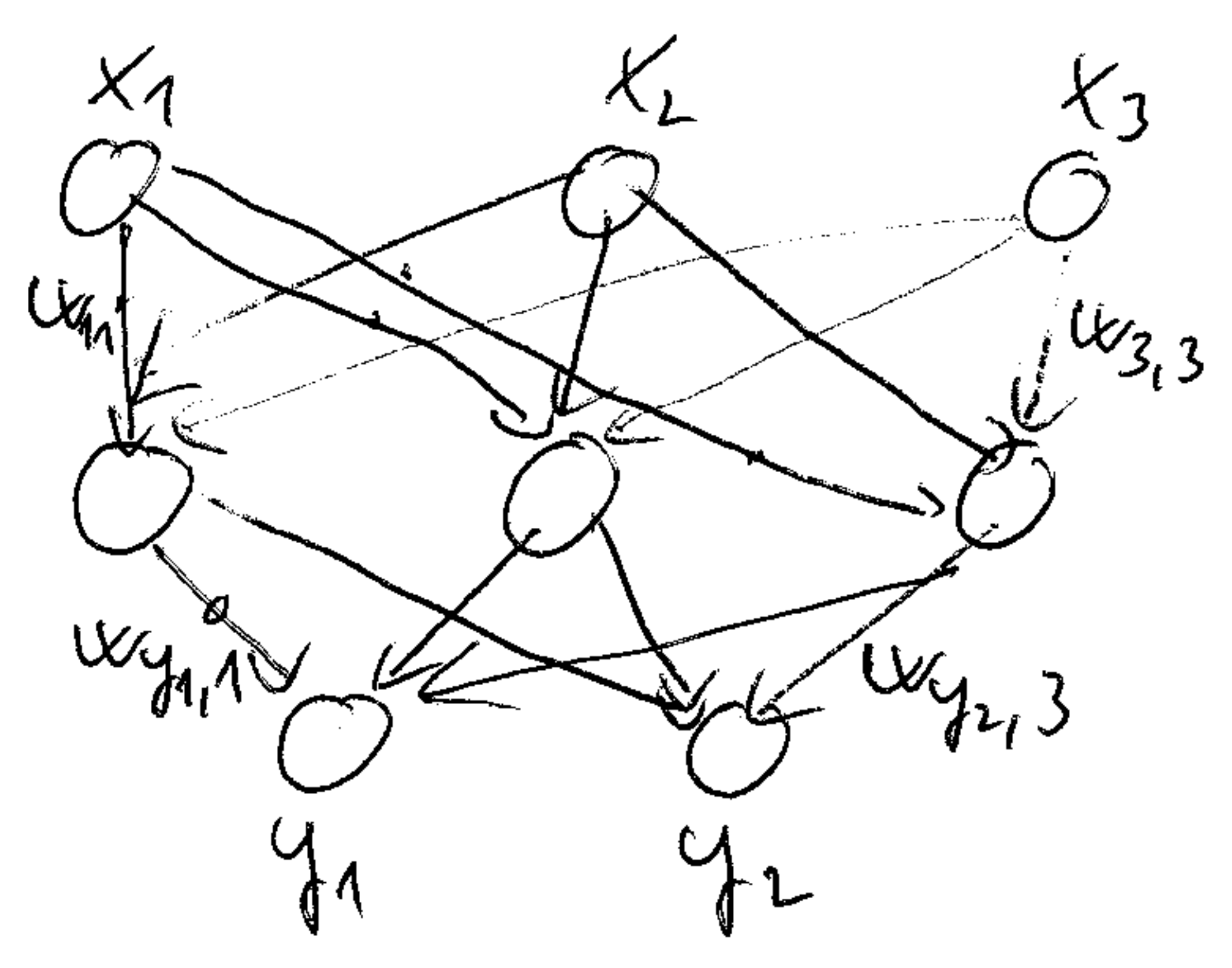
biology: neural cell



$$y_j = f_j \left(w_{j0} + \sum_{m=1}^M w_{jm} x_m \right)$$

Simplest structure:

input layer
hidden layer
output layer



functions:

sigmoid

$$f(z) = \frac{1}{1 + e^{-z/a}}$$

tangent hyperbolic

$$f(z) = \frac{e^{z/a} - e^{-z/a}}{e^{z/a} + e^{-z/a}}$$

Supervised learning:

N observation, M input variables

J points in the hidden layer, R number of hidden layers

K number of outputs

$R=1$ all continuous functions can be approximated
 $R=2$ also non continuous ones! (derivable ones)

- 1.) Training: N cases where x and y are known \Rightarrow determine all weights
- 2.) Use it when only x -s are known, to get y

Dataset size for training:

$$\frac{J(M+K)+1}{N} < 0.1$$

to avoid overfitting

Deep-learning: more hidden layers

1/1-2

How to find W-s?

X input matrix (N x M)

y output matrix (N x K)

weight matrix

minimise: $E = \|Y - F(W, X)\|$

- gradient method by layers
iterative

$W \rightarrow \underline{w} \in$ row vector in W

$$\underline{w}(t+1) = \underline{w}(t) - \beta \frac{\partial E^2(W, t)}{\partial \underline{w} t} = \underline{w}(t) - \beta \underset{\substack{\uparrow \\ \text{constant}}}{d(t)}$$

mostly by least squares:

$$\left(\frac{\partial E^2}{\partial \underline{w}}\right) = 2E \frac{\partial E}{\partial \underline{w}} = -2E \frac{\partial F(W, X)}{\partial \underline{w}} X$$

- back propagation algorithm

for y answer on x input

target } $E = \frac{1}{2} \cdot (y - f(z(x, w)))^2$
min }

$$\frac{\partial E}{\partial w_i} = \frac{dE}{df} \cdot \frac{df}{dz} \cdot \frac{\partial z}{\partial w_i}$$

(if $z = \sum w_i x_i$
 $\frac{\partial z}{\partial w_i} = x_i$ $\frac{\partial z}{\partial x_i} = w_i$)

$$f(z) = \frac{1}{1 + e^{-z}} \text{ (sigmoid)}$$

$$\frac{df}{dz} = \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} = \frac{1}{(1 + e^{-z})} \cdot \left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}}\right) = f \cdot (1 - f)$$

$$\frac{\partial E}{df} = f - y$$

$$\frac{\partial E}{\partial w_i} = (f - y) \cdot f \cdot (1 - f) \cdot x_i$$

$$\Delta w_i = - \frac{\partial E}{\partial w_i}$$

\in step in gradient direction, for all i

η - learning coefficient \in effect of new

α - suspension

\in old

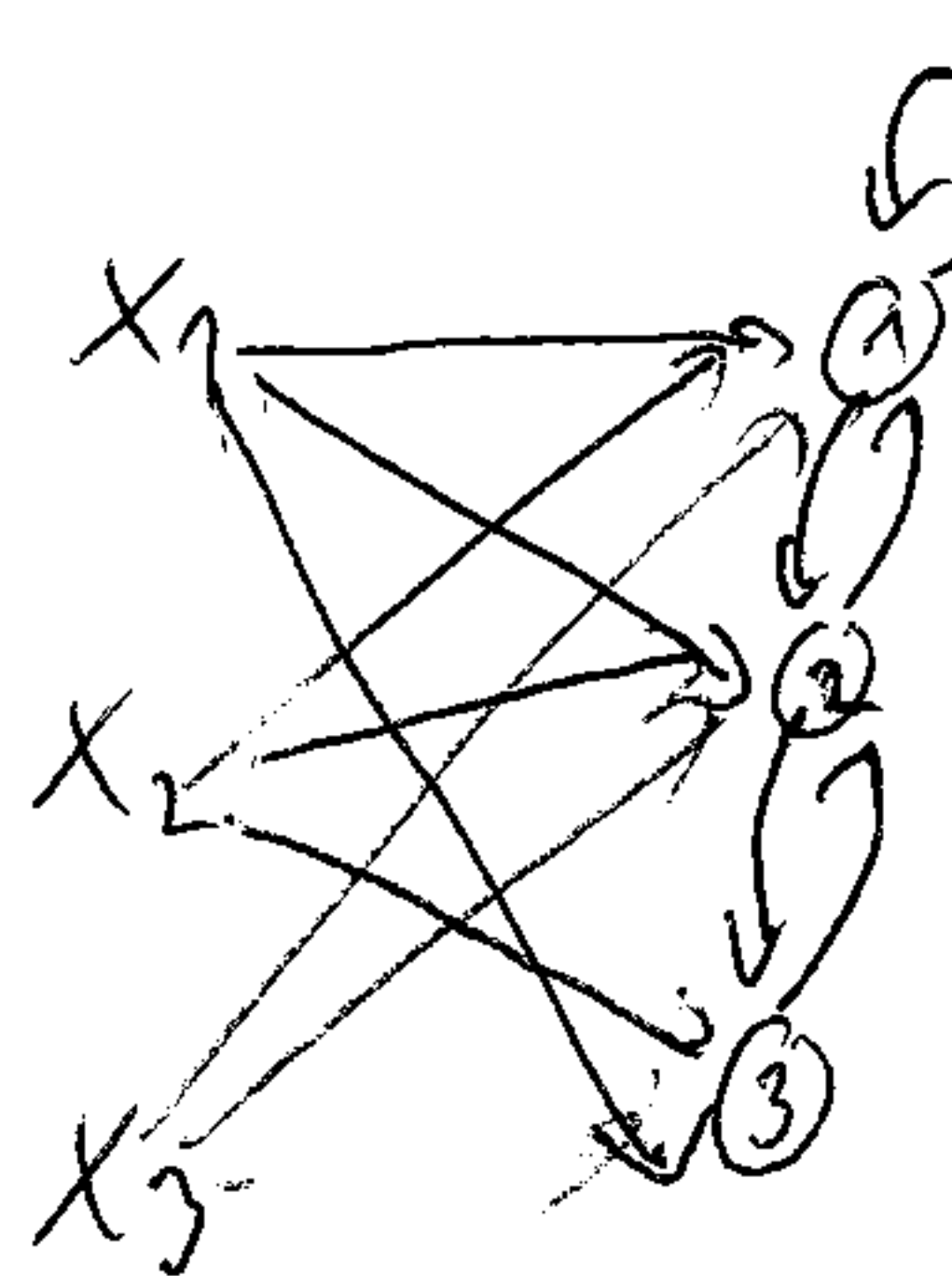
Unsupervised learning

NN-3

output is not important

mapping of input onto a lower dimension \rightarrow clustering of objects
clustering of variables

Kohonen network \approx memory in brain data site reduction



functions on a grid, eg. 2D grid r_i, r_j, \dots
change $\underline{w}_i \rightarrow$

- select an $\underline{x}_i \rightarrow$ search the most similar $\underline{w}_k \Rightarrow \underline{w}_k^*$

$$\underline{w}_k(t+1) = \underline{w}_k + \eta \phi(r_{kn}) (\underline{x} - \underline{w}_k(t))$$

learning speed

physical distance, eg. gaussian

\Rightarrow objects are mapped on a grid

similar objects on the same / neighboring grid

many different versions for NN

- telecommunication connection

- ...