

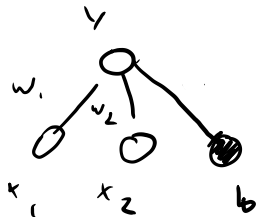
Klassisk maskininläring

• regression, statistiska metoder

Neurala nätverk

50-talet: baserat på idag något föräldrad syn på hjärnan.
Ändå så abstrakt att "det funkar ändå".

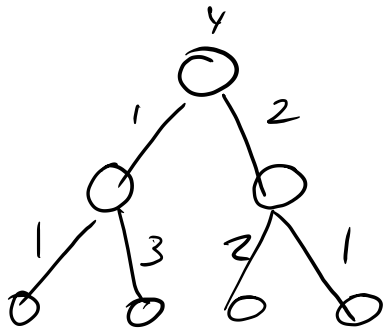
Perceptron



$$y = w_1 x_1 + w_2 x_2 + b$$

$$y = w \cdot x + b$$

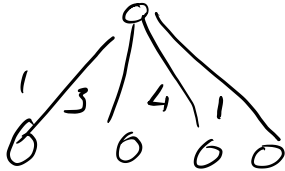
"2D linje" - plan som spänns upp av w



$$y = 1(1x_1 + 3x_2) + 2(2x_3 + 1x_4)$$

$$y = 1x_1 + 3x_2 + 4x_3 + 1x_4$$

$$L(a+b) = L(a) + L(b)$$



Linjärd, alltså inte komponerbard.
Vi behöf någon icke-linjäritet!

Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$

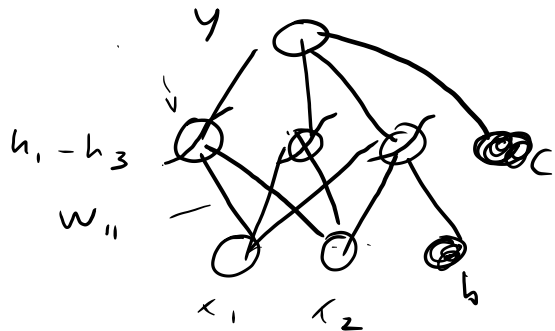
ReLU: $v(x) = \begin{cases} x & \text{om } x > 0 \\ 0 & \text{annars} \end{cases}$



MLP Multilayered perceptron



50-80 talet, teorin för hur man tränar djupa nätverk fanns inte. Dessutom beräkningsmässigt dyrt.

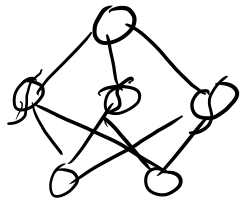


$$f(x) = \sigma(Wx + b) + c$$

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b = \begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix}$$

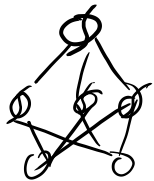
$$\sigma \Rightarrow \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$

$$\begin{bmatrix} v_1 & v_2 & v_3 \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} + c = y$$



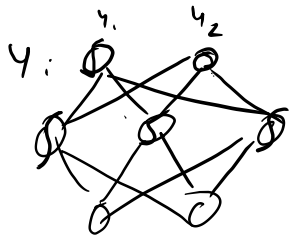
← regression

← feature expansion



sigmoid på $y \rightarrow [0, 1]$

tolka som sannolikheter \rightarrow binär klassificering



⑤ σ

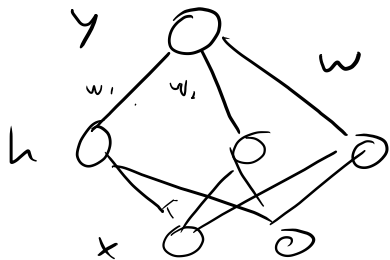
↑

softmax

$$y_i = \frac{e^{l_i}}{\sum_j e^{l_j}}$$

$$\sum_{i=1}^r y_i^2 = 1$$

\rightarrow multinomial



lärja slumpmässiga vikt

$y - \hat{y}$ - justera varje vikt "lite grann"
mot en bättre lösning.

Iterativt för all träningsdata:

$$w \leftarrow w - \eta \nabla C(w)$$

kedjeregeln över funktionella moduler

moduler:

$$f(x) = \frac{1}{1 - e^{-x}}$$

$$a(x) = -x$$

$$b(a(x)) = e^{-x}$$

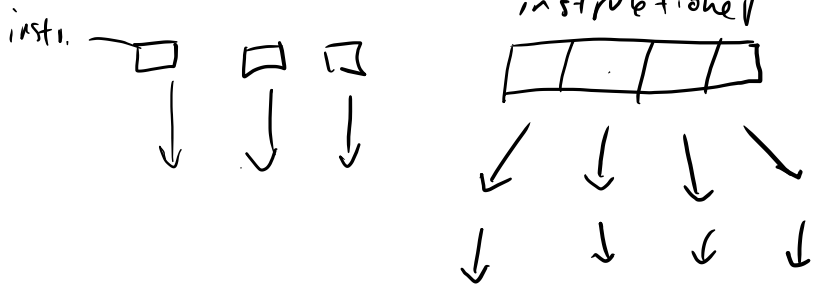
$$c(b(a(x))) = 1 - e^{-x}$$

$$d(c(b(a(x)))) = \frac{1}{1 - e^{-x}}$$

→ forward: symbolisk derivering

← backward: numeriskt

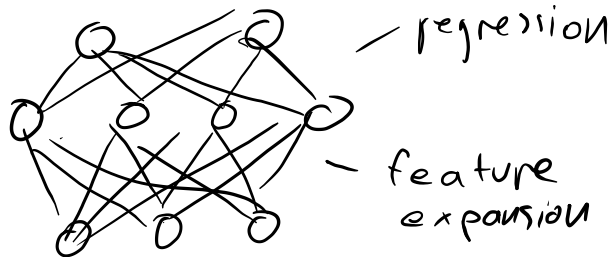
parallelism (på två sätt)



kompiatorer för vektorinst.
våldigt svåra att skriva.

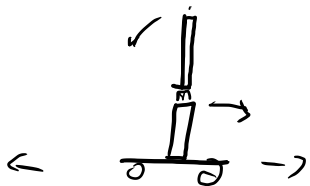
Men AA fungerar som vektorprocessorer.

NN



lär sig features

βX



two-pole balancing

SVM

$$K(x_i, x_j)$$

expanderar featuresymden
våldigt mycket, kollapsar
dimensionen för att undvika
overfit

βX

NLP

Klassiska tekniker (innan 2017) räcker inte.

GPT: generative pre-trained transformers.

Transform: Encoder \rightarrow dolda lager \rightarrow Decoder
GPT 1

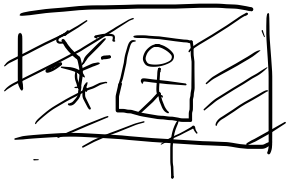
Embedding:

omvandlar ord till tal. Vi kallar dem då "tokens"



hidden layers

Attention



Du kanske inte är så dum ändå.

• skalär produkt mellan inbäddade ord.

problem med syftningar; höger eller vänster association.

imp (x, y)

$$1 + 2 * 3 \div 4$$

self-attention (masked self attention)

titta endast åt ett håll.

Du kanske inte är så dum ändå. $\rightarrow [0, 1, 0, 3, \dots]$

Optimeringsvillkor för oövervakad inlärning
(pre-training)

$$\max \sum_i^n \log P(u_i | u_{i-k}, \dots, u_{i-1}; \theta)$$

Vi lär alltså systemet att lösa först.

Sedan tränar vi övervakat:

<sentence> | <answer>

På löp vi oss domänkunskap genom att förutsäga svaren.
(tex fakta om Mozart)

$$A(e) = (Q, K, V)$$

\downarrow \downarrow \downarrow
 \downarrow \downarrow \downarrow
 termerna nycklar värden



Multi-head, multi-cell self attention