

FoDA

L28

- Neural Networks
- + other classifiers

Classification

Input $X \subset \mathbb{R}^d$ labels $y \in \{-1, +1\}^n$
 x_1, \dots, x_n

Goal function $g: \mathbb{R}^d \rightarrow \mathbb{R}$

want $\text{sign}(g(x_i)) = y_i$

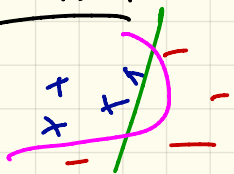
linear

$$g_{\alpha}(x) = \langle \alpha, x \rangle + \alpha_0$$

(non-linear)

kernel

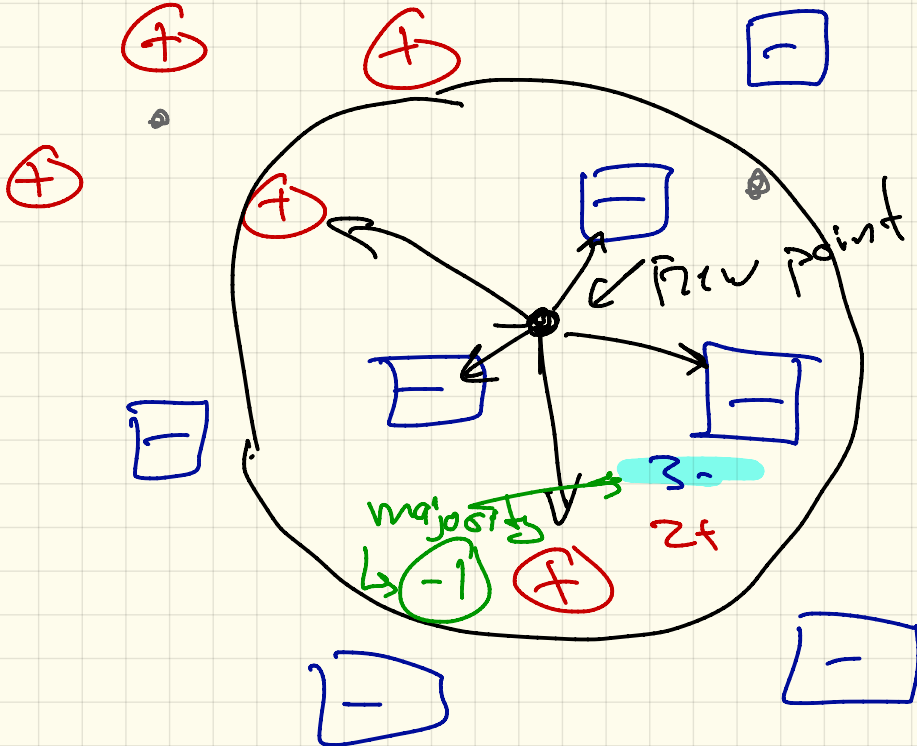
$$g_{\alpha}(p) = \sum_{j=1}^n \alpha_j k(x_j, p) + \alpha_0$$



↳ still be optimized w/ GD
convex

KNN Classifier

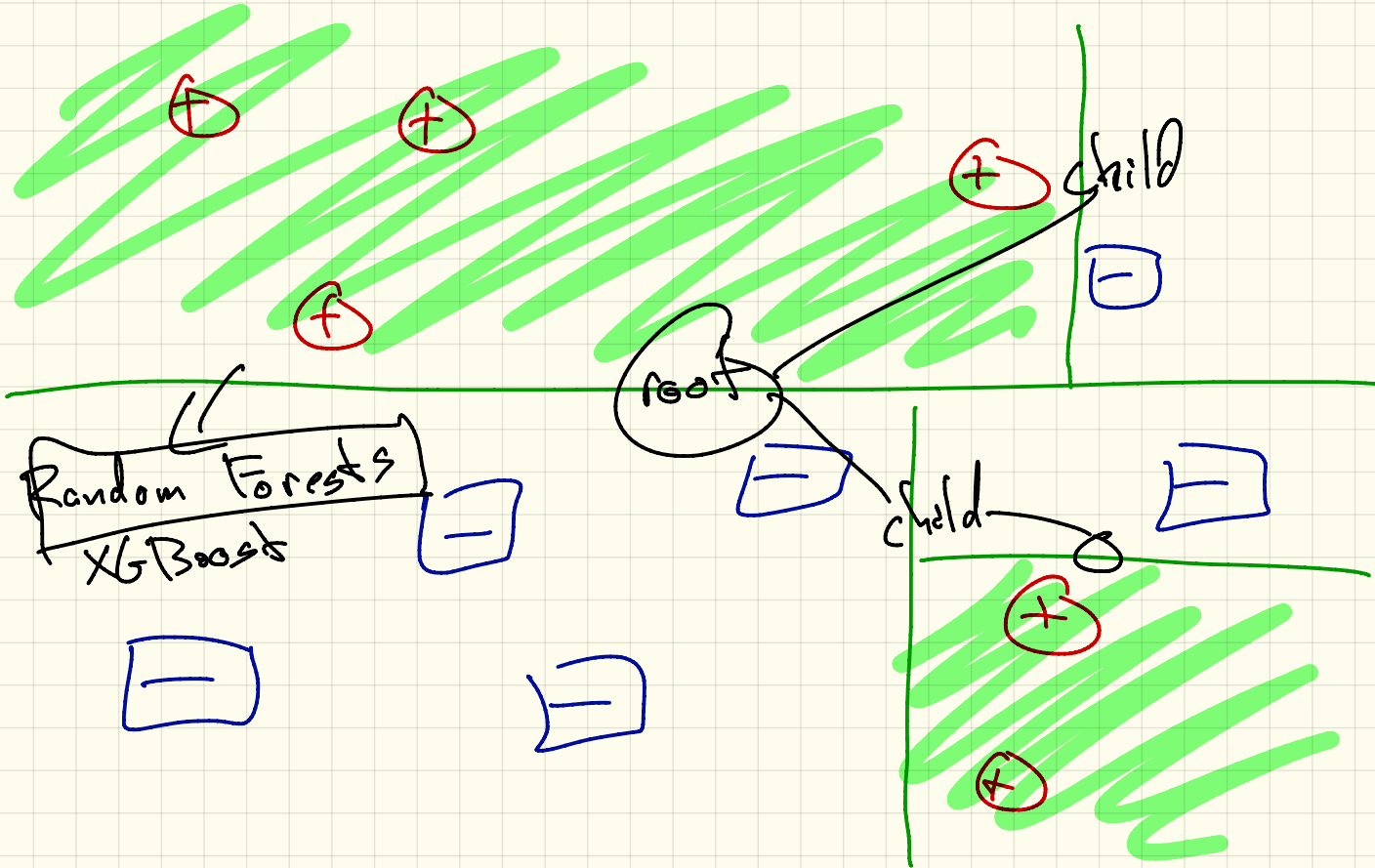
\mathbb{R}^d



Input (x_{ig})
& distance
 $d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

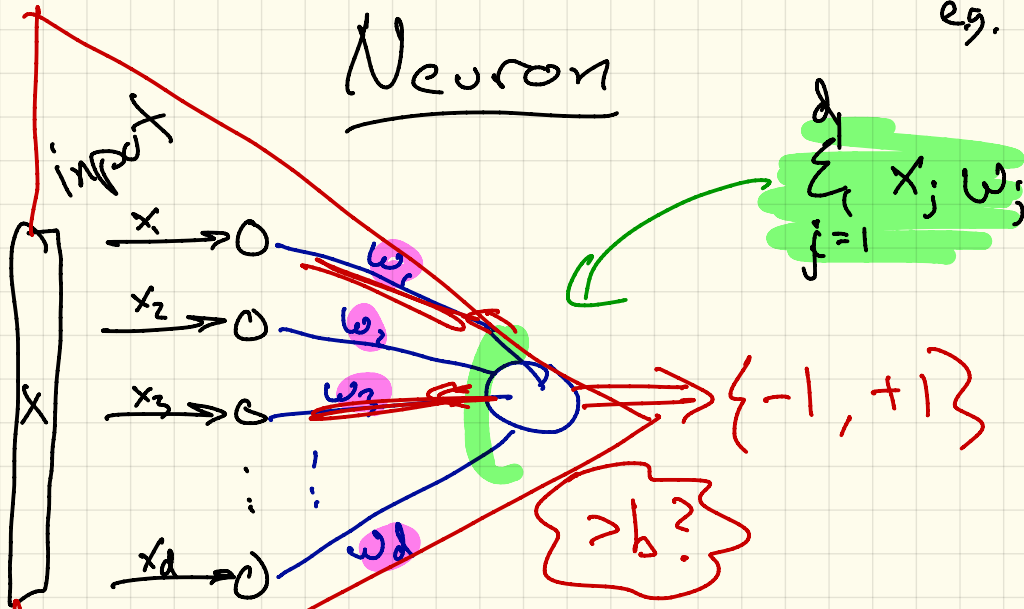
Find closest k neighbors
($k=5$)

Decision Trees



Neuron

eg. uniform $w_j = \frac{1}{d}$



?
 $> b$

$> b?$

linear classifier

$$\sum_{j=1}^d x_j w_j > b \iff \langle x, w \rangle - b > 0$$

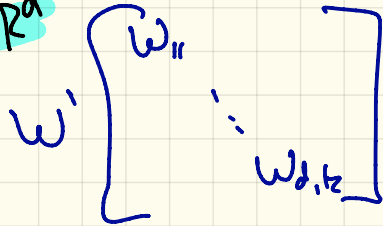
Neural Network

input

b.d weights



$\in \mathbb{R}^d$

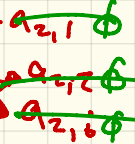


$$a_i = w x$$

$\in \mathbb{R}^{k_1}$



$W^2 \phi$



$h_2 \in \mathbb{R}^{k_2}$

$\dots W^3 \phi, W^4 \phi, \dots$

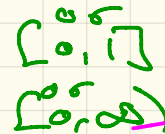


Final classifier

$\{ -1, +1 \}$

ϕ activation

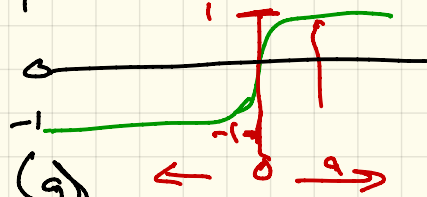
$$\phi: \mathbb{R} \rightarrow [-1, +1]$$



Activation Functions ϕ

goal: ϕ proxy $f: \mathbb{R} \rightarrow \begin{cases} +1 & \text{if } \geq 0 \\ -1 & \text{if } < 0 \end{cases}$

② differentiable

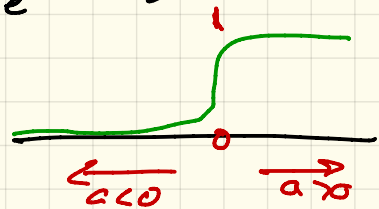


hyperbolic tangent

$$\begin{aligned} \phi(a) &= \tanh(a) \\ &= \frac{e^a - e^{-a}}{e^a + e^{-a}} \in [-1, +1] \end{aligned}$$

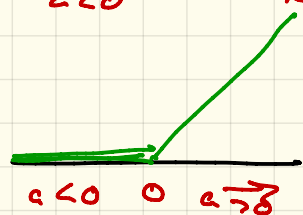
sigmoid

$$\phi(a) = \frac{e^a}{e^a + 1} \in [0, 1]$$



ReLU

$$\phi(a) = \max\{0, a\} \in [0, \infty)$$



Neural Network Function

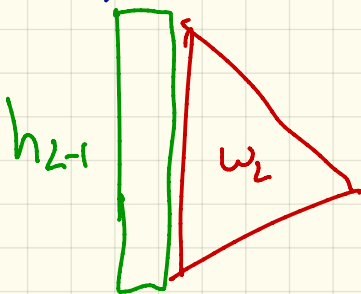
$$g(x) = g(x) = -b + w^L \phi \left(w^{L-1} \dots \phi \left(w^2 \phi \left(w^1 x \right) \right) \dots \right)$$

parameters $(1 + k_{L-1} + k_{L-1} \cdot k_{L-2} + \dots + k_1 \cdot k_2 + d \cdot k_1)$ $\rightarrow \alpha \in \mathbb{R}^{\#A}$

Back Propagation

$$\nabla_{g(x)}$$

$$\nabla_{w^L(x)}$$



h_{L-1}

$l(z_i)$

$z_i = y_i g(x_i)$

$g(x_i)$
 y_i

1. compute $g(x)$

1a $h_1(x)$

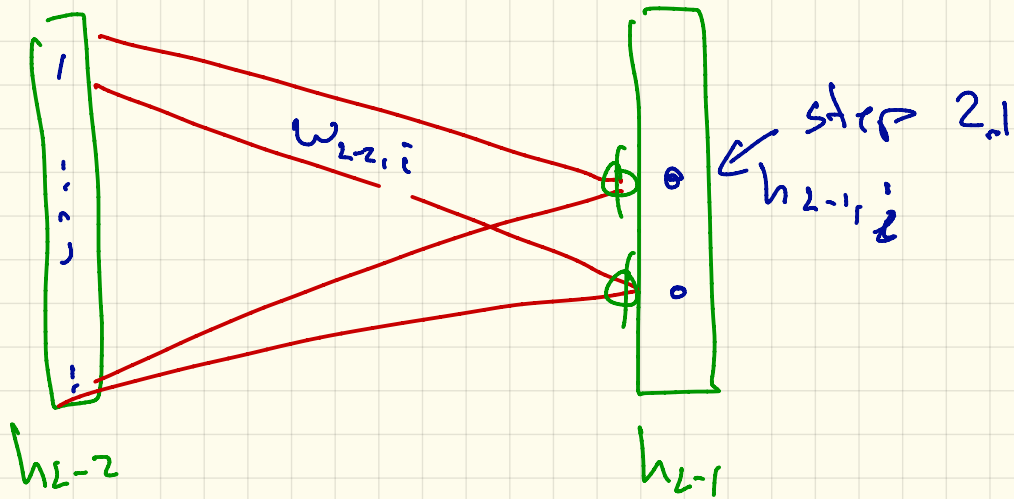
1b $h_2(x) = h_2(h_1(x))$

$h_{L-1}(x) = h_{L-1}(\dots h_1(x))$

2. $\nabla g(x)$

$z_1 \nabla_{w_2}(x, y, h_{L-1})$

\vdots



$$\nabla_{w_{L-2,i}} (x_i, y) =$$

instead of $l(z_i)$
 $\phi(\nabla h_{L-1,i})$

chain rule
 \hookrightarrow dynamic programming