

FODA L28

Other Classifiers

SVMs, Neural Nets,

Decision Trees, etc.

Input $(X, y) \in \mathbb{R}^d \times \{-1, +1\}$
 $x_i \in \mathbb{R}^d$
 $y_i \in \{-1, +1\}$

Goal function

$$f: \mathbb{R}^d \rightarrow \{-1, +1\}$$

∴ $g: \mathbb{R}^d \rightarrow \mathbb{R}$

so linear $g(x_i) = \langle \alpha, (1, x_i) \rangle$

$$f(x_i) = y_i$$

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

or as many
 x_i, y_i as
possible

Classification

Supervised

Cross-Validation

$(X, y) \rightarrow$ Train
 \rightarrow Test

Goal

Generalize on
new data
rid new (x_{new})

Non-linear Kernel Classifiers

$$g(x) = \langle w, x \rangle = \sum_{i=1}^n \alpha_i y_i \langle x_i, x \rangle$$

α $x = (x_1, \dots)$ \swarrow mistake counter $\alpha \in \mathbb{R}^n$

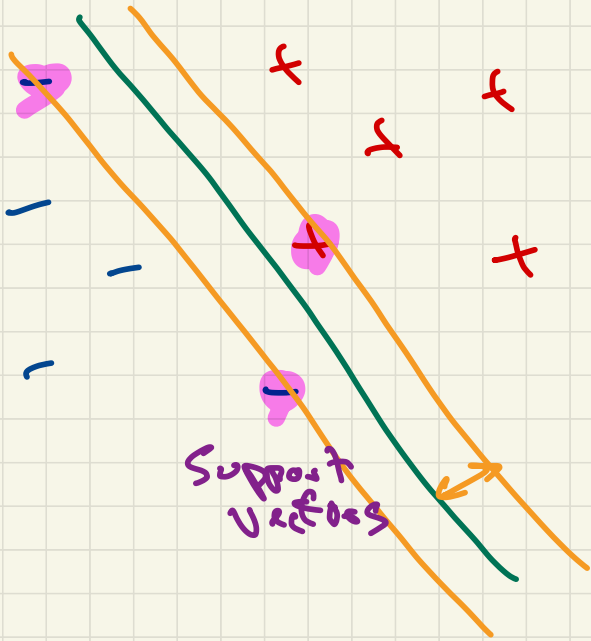
$$= \sum_{i=1}^n \alpha_i y_i K(x_i, x)$$

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$$

mostly $\alpha_j = 0$

Non-0 $\alpha_j =$ support vectors

$$\alpha = (\underbrace{2}_{i=1}, 0, 0, \underbrace{1}_{i=4}, 0, 0, \underbrace{2}_{i=7}, 0, \underbrace{1}_{i=9})$$



Cost function

$$h(\alpha) = \sum_{i=1}^n \ell(y_i; g_{\alpha}(x_i))$$

$$g_{\alpha}(x_i) = \sum_{j=1}^n \alpha_j y_j K(x_j, x_i)$$

$\alpha \in \mathbb{R}^n$, but can restrict to k support vectors.

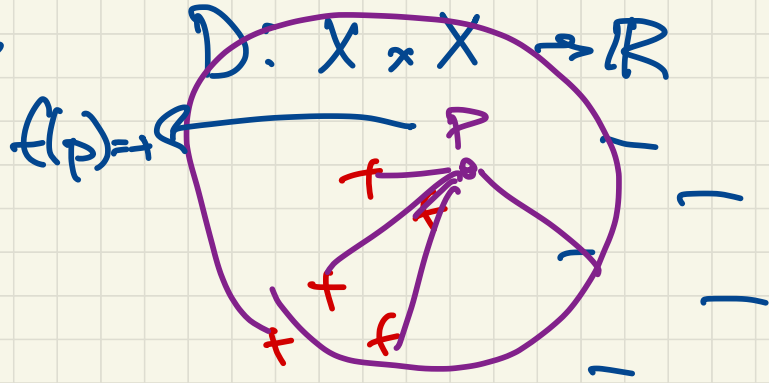
Run gradient descent on $h(\alpha)$

Find ~~(a)~~ Perceptron $z_i = y_i g_{\alpha}(x_i)$
 (b) SGD w/ hinge loss

KNN classifiers

Uses distance

No training

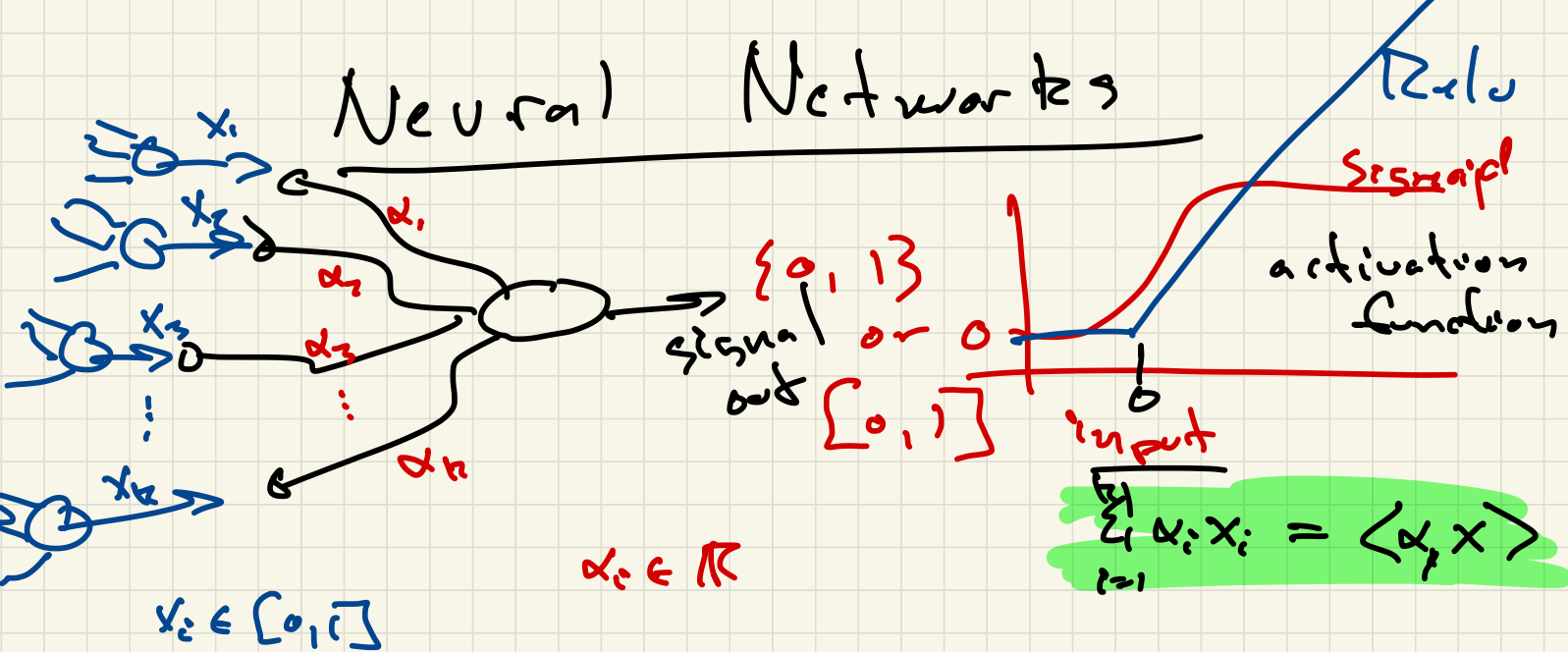


$$f(p) \rightarrow \{-1, +1\}$$

= Find the k closest points
 $x_1, x_2, \dots, x_k \in X$ to p .

Then vote, which sign is
most prevalent $\rightarrow \{-, +\}$

Neural Networks



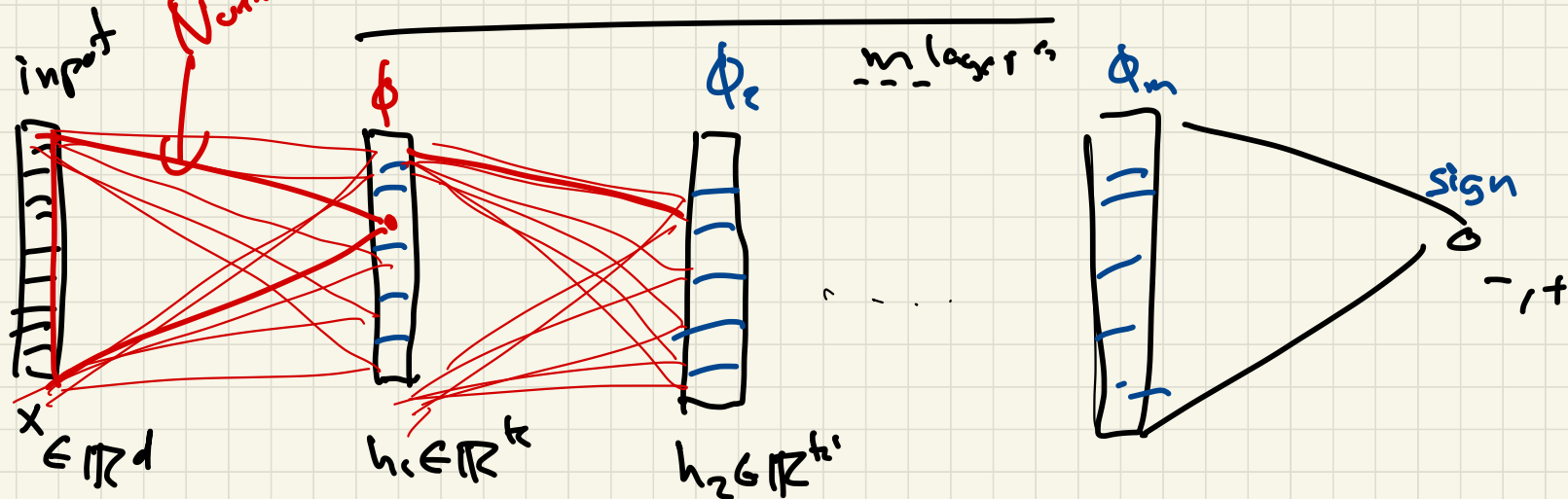
Activation Function

$$\phi: \mathbb{R} \rightarrow [0, 1] \quad \text{or ... } [-1, +1], \quad \text{or } [0, \infty)$$

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

ReLU $\phi(y) = \max(0, y) \in [0, \infty)$

Neural Networks

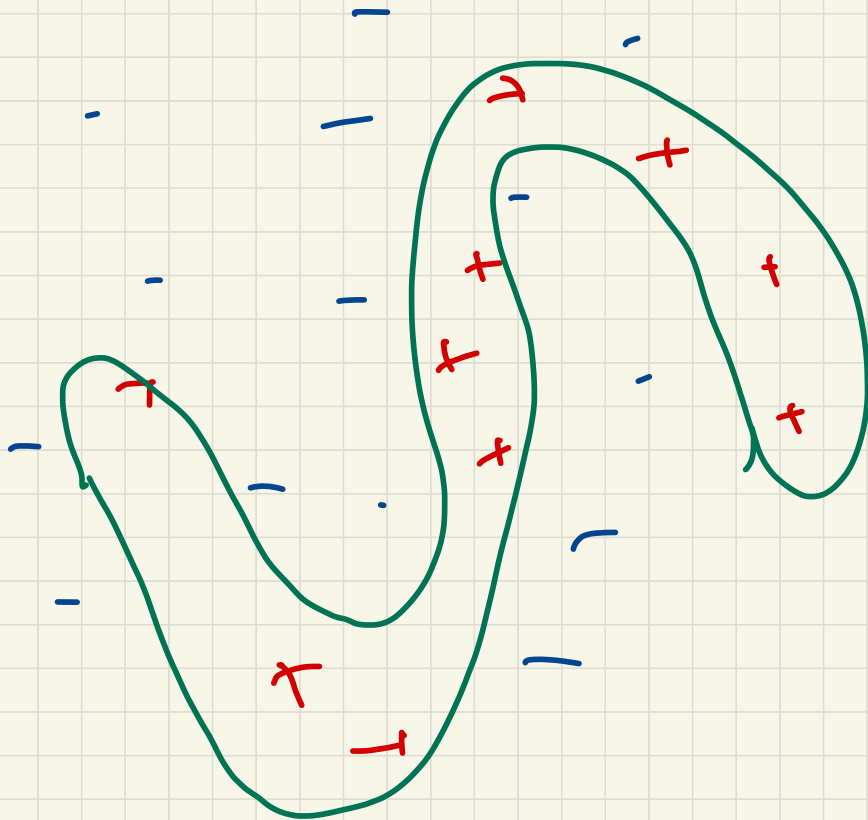


Neuron, α_j, ϕ_j
 $h_{i,j} = \phi_j(\langle \alpha_j, x \rangle)$

Parameters

$\alpha_1, \dots, \alpha_{k_1}$ layer 1
 $\alpha_1, \dots, \alpha_{k_2}$ layer 2
 \vdots
 $\alpha_1, \dots, \alpha_{k_m}$ layer m

$\approx m \cdot k^2$
 parameters.

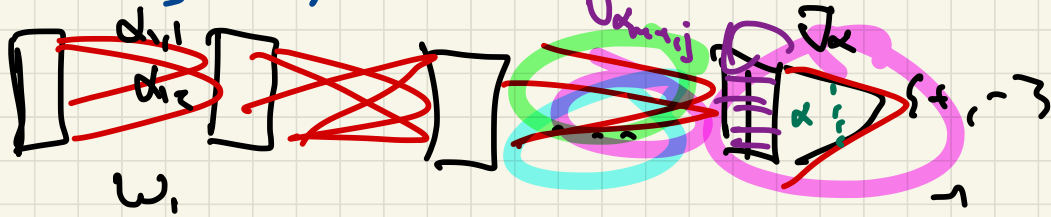


Train (Deep) NN w/ GD.

Can efficiently compute gradient.

$\approx k^2 m$ parameters

by layers from last backwards.



Backpropagation

x_i

w_i

$$w_i = \begin{matrix} \alpha_{i,1} \\ \alpha_{i,2} \\ \alpha_{i,\dots,k} \end{matrix}$$

h_m

$$(y_i - \hat{y}_i)^2$$

$$g(x_i) = w^m \phi^m (w^{m-1} \phi^{m-1} (\dots \phi'(w^1 x) \dots))$$

Decision Trees

$X \in \mathbb{R}^d$ or $\{A, \dots, C\}$

