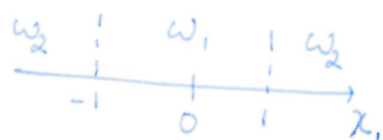


B Other mapping functions

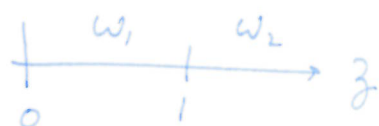
E4.13 : Feature mapping

(a)



not linearly separable

Let $z = x_1^2 \geq 0$



linearly separable

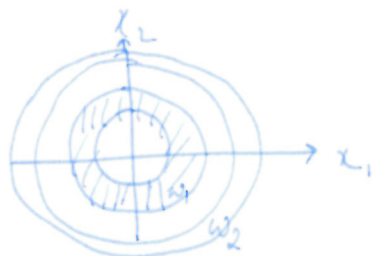
(b) dataset 3



$z = (x_1 - 6)^2$



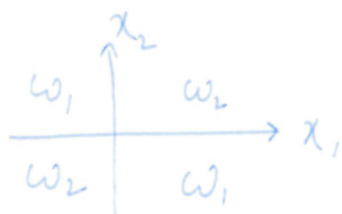
(c) dataset 2



$z = \sqrt{x_1^2 + x_2^2}$



(d) XOR



$z = x_1 x_2$



Pro and cons of feature engineering / mapping:

- +) same learning algorithm as for LDF in 3
- +) more general decision regions in \mathcal{X}
- difficult to find a good feature mapping for $d > 3$
 - often polynomial used
-) huge dimension d'

A polynomial discriminant function can be formulated as a LDF:

$$f_j(\underline{x}) = \underline{u}_j^T \underline{z} \text{ with}$$

$$\underline{z} = [1, \quad x_1, \dots, x_d, x_1^2, x_1 x_2, \dots]^T$$

$$\underline{u}_j = [w_{j,0}, w_{j,1}, \dots, w_{j,d}, w_{j,11}, w_{j,12}, \dots]^T$$

$\underline{x} \in \mathbb{R}^d$: original feature space

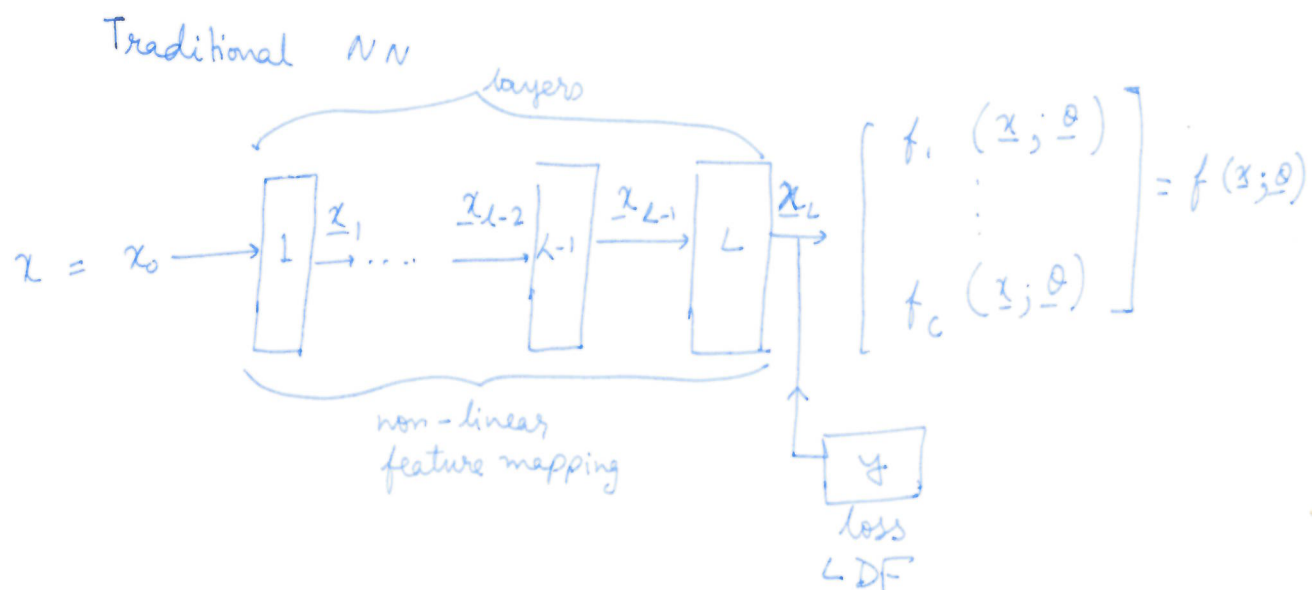
$\underline{z} \in \mathbb{R}^{d'}$: mapped feature space

eg $d = 100$

LDF	$x_k x_l$	$x_k x_l x_m$
$d' = 100 + 1 = 101$	$O(100^2)$	$O(100^3)$

- • higher computational complexity
- not enough training samples → overfitting

4.5.3 Neural Network (NN)



layers:

$$x_L = \phi_L \left(\underbrace{\sum_{l=1}^L w_l x_{l-1} + b_l}_{LDF} \right) ; 1 \leq l \leq L,$$

non-linear ~~do~~ activation f ?

$$\phi_L(a) = a$$

Training :

loss $\ell(f(x; \theta), y)$ eg L_2 norm $\hat{=}$ LS

cost f : $L(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(f(x^{(n)}; \theta); y^{(n)})$

$$\min_{\theta} L(\theta)$$

NN: automatically learnt non-linear feature mapping
+
learned LDF.

4.5.4 Support Vector Machine (SVM)

(4)

Vapnik, 1996

Basic idea :

- a binary classifier : $y_n \in \{-1, 1\}$ instead of w_1, w_2
- use a LDF $f(\underline{x}) = \underline{w}^T \underline{x} + w_0$
- estimated class $\hat{y}(\underline{x}) = \text{sign}(f(\underline{x})) \in \{-1, 1\}$

Advanced ideas

- non-linear-feature mapping

$$\underline{z} = \phi(\underline{x}) \text{ instead of } \phi_L(\cdot) \text{ in NN}$$

- kernel trick

- training

- maximum margin instead of LS in NN

- convex optimization problem instead of non-convex in NN

4.5.4.1 Maximum margin

(5)

Assume $S = \{(\underline{x}_n, y_n), 1 \leq n \leq N\}$ linearly separable
i.e. $y_n > 0$

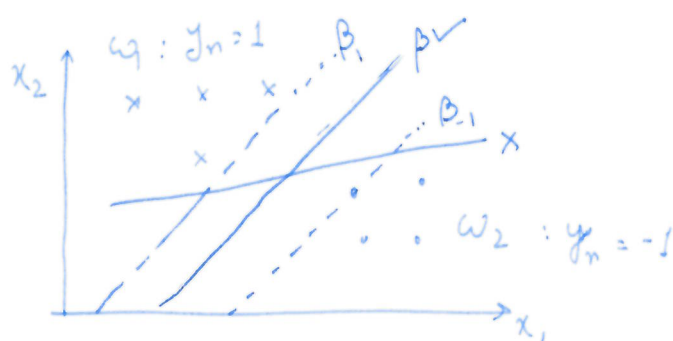
\Rightarrow infinite number of LDFs $f(\underline{x})$ with
 $\gamma(f(\underline{x}), S) > 0$

Which one?

NN: $\angle S$ - Least Squares

SVM: maximum margin

$$\max_{\underline{w}, w_0} \gamma(f(\underline{x}), S)$$



$\beta_1, \beta, \beta_{-1}$: 3 parallel hyperplanes

$\beta = \{ \underline{x} \mid f(\underline{x}) = 0 \}$: desired decision boundary

β_1, β_{-1} : 2 limit hyperplane with the same margin
 $r > 0$ to β

β_1 : passes through the closest \underline{x}_n from w_1
" " " " " w_2
 β_{-1} : " " " " " w_2

The pts on β_1, β_{-1} are called support vectors

$$SV = \{n \mid \underline{x}_n \text{ on } \beta_1 \text{ or } \beta_{-1}\}$$

Scaling of $f(\underline{x})$:

If $f(\underline{x})$ is a LDF, then $\alpha f(\underline{x})$ with $\alpha > 0$ is also a LDF

Choose α such that $|f(\underline{x}_n)| = 1 \quad \forall n \in SV$

$$\bullet \quad n \in SV : \text{margin } \gamma_n = y_n \cdot \frac{f(\underline{x}_n)}{\|\underline{w}\|} = \frac{1}{\|\underline{w}\|}$$

$$\bullet \quad n \notin SV : \text{margin } \gamma_n > \frac{1}{\|\underline{w}\|}$$

$$\bullet \quad \gamma(f(\underline{x}), S) = \min_n \gamma_n = \frac{1}{\|\underline{w}\|}$$

Maximum margin Training:

$$\max \gamma \text{ or } \min \|\underline{w}\| \text{ or } \min_{\underline{w}, w_0} \frac{1}{2} \|\underline{w}\|^2$$

$$\text{s.t. } y_n f(\underline{x}_n) = y_n (\underline{w}^T \underline{x}_n + w_0) \geq 1 \quad \forall n$$

in comparison to LS learning in NN:

+ SVM less sensitive to statistical variations in \underline{x}
 + " " " to noise in \underline{x}

\Rightarrow lower generalization error rate

+ perfect classification for linearly separable data

+ Convex optimization; no local minima, global convergence

+ only support vectors determine $f(x)$,

"forget" all non-support vectors

→ reduced complexity

4-62

4-63