- Sparse Kernel methods          PRML Ch. 7
- maximum margin classifier
- support vector machines

  - linearly separable data

  - overlapping data

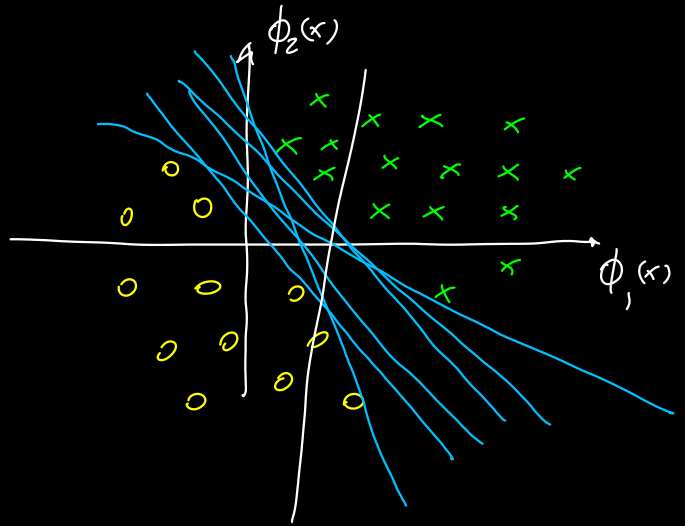# Maximum margin classifier

two-class problem

$$y(x) = w^T \phi(x) + b$$

training data

$$\{x_1 \cdots x_N \; ; \; x_n \in \mathbb{R}^D\}$$
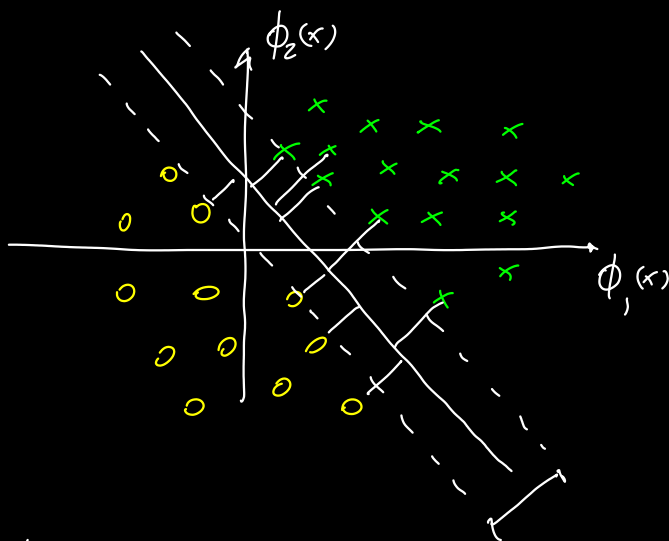$$\{t_1 \cdots t_N \; ; \; t_n \in \{-1, +1\}\}$$

$$\exists \; w, b \quad s.t. \quad \begin{array}{l} y(x_n) > 0 \Longleftrightarrow t_n = +1 \\ y(x_n) < 0 \Longleftrightarrow t_n = -1 \end{array} \Bigg\} \Rightarrow y(x_n) \cdot t_n > 0 \quad \forall_n$$

Goal: minimize generalization error

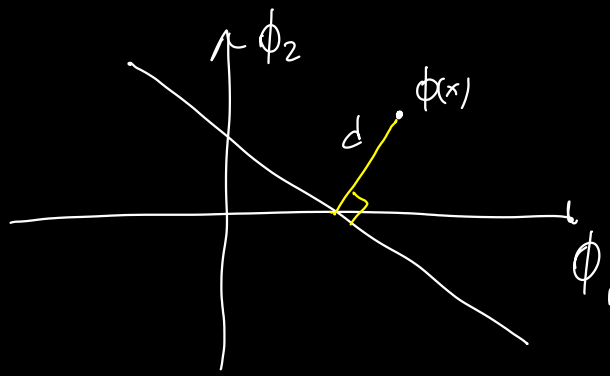We only have the training data (do not know the underlying distribution)
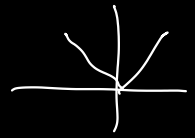
Heuristic solution: maximize margin

$w, b$

Recall:

$-$ if $x$ is on decision boundary $\Rightarrow y(x) = 0$

$$d = \frac{|y(x)|}{\|w\|}$$

— we are only interested in no error $\Rightarrow t_n\, y(x_n) > 0$

$$d = \frac{t_n\, y(x_n)}{\|w\|} = \frac{t_n\left(w^T \phi(x_n) + b\right)}{\|w\|} \qquad \forall n \in [1, N]$$

Optimization:

$$\underset{w,b}{\arg\max} \left\{ \frac{1}{\|w\|} \min_n \left[ t_n\left(w^T \phi(x_n) + b\right)\right]\right\} \qquad \begin{array}{l}\text{hard to}\\ \text{solve}\\ \text{directly}\end{array}$$

$$\uparrow$$
$$\text{independent of } n$$

$$w, b \rightarrow \kappa w, \kappa b$$

$$\frac{t_n\left(\kappa w^T \phi(x) + \kappa b\right)}{\|\kappa w\|} = \frac{t_n\left(w^T \phi(x) + b\right)}{\|w\|}$$

Canonical representation of the decision hyperplane

$$t_n\left(w^T \phi(x_n) + b\right) = 1 \qquad \text{for the point that is}$$
$$\text{closest to the hyperplane}$$

Then
$$t_n\left(w^T \phi(x_n) + b\right) \geq 1 \qquad \forall n \in [1, N]$$

$$\underset{w,b}{\text{argmax}} \left\{ \frac{1}{\|w\|} \underbrace{\min_n \left[ t_n \left( w^T \phi(x) + b \right) \right]}_{1} \right\} =$$

$$= \underset{w,b}{\text{argmax}} \left\{ \frac{1}{\|w\|} \right\} = \underset{w,b}{\text{argmin}} \frac{1}{2} \|w\|^2$$

subject to the constraint $\quad \underline{t_n \left( w^T \phi(x) + b \right) \geq 1} \quad \forall n \in [1,N]$



$$t_n \left( w^T \phi(x) + b \right) = 1$$

$$t_n \left( w^T \phi(x) + b \right) > 1$$

# Appendix E

quadratic programming

Lagrange multipliers $\quad a = (a_1 \cdots a_N)^T$

$\quad a_n \geq 0 \quad$ 1 for each data point

$$L(w, b, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n \left( w^T \phi(x_n) + b \right) - 1 \right\}$$

minimize $\qquad\qquad \geq 0$

$$\frac{\partial L}{\partial w} = 0 \qquad w = \sum_{n=1}^{N} a_n t_n \phi(x_n) \quad \longleftarrow$$

$$\frac{\partial L}{\partial b} = 0 \qquad \boxed{0 = \sum_{n=1}^{N} a_n t_n}$$

$$\tilde{L}(a) = \frac{1}{2} \left\| \sum_{n=1}^{N} a_n t_n \phi(x_n) \right\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n \overbrace{\left( \sum_{m=1}^{N} a_m t_m \phi^T(x_m) \right)}^{w^T} \phi(x_n) \right\}$$

$$- b \sum_{n=1}^{N} a_n t_n + \sum_{n=1}^{N} a_n =$$

$$= \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m \phi^T(x_n) \phi(x_m) - \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m \phi^T(x_n) \phi(x_m)$$

$$\tilde{J}(a) \qquad + \sum_{n=1}^{N} a_n = -\frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m \underbrace{\phi^T(x_n) \phi(x_m)}_{K(x_m, x_n)} + \sum_{n=1}^{N} a_n$$

dual representation for the
optimality criterion

$$a_n \geq 0$$

$$\sum_{n=1}^{N} a_n t_n = 0 \qquad n \in [1, N]$$

Quadratic problem in $M$ variables $\Rightarrow \underline{O(M^3)}$

$$L(\overbrace{w}^{M}, b, a)$$
$$\uparrow$$

$$\tilde{L}(a) \qquad N \text{ variables} \Rightarrow O(N^3)$$

$$y(x) = w^T \phi(x) + b \stackrel{=}{_{\uparrow}} \sum_{n=1}^{N} a_n t_n \underbrace{\phi^T(x_n) \phi(x)}_{K(x, x_n)} + b$$
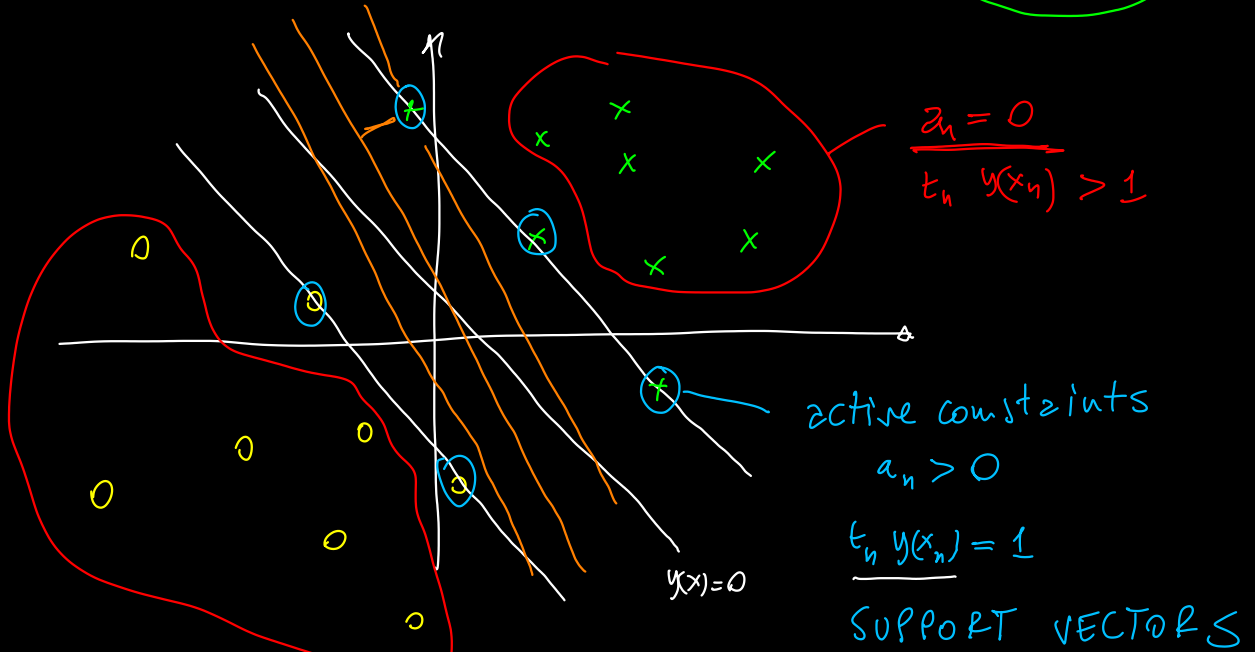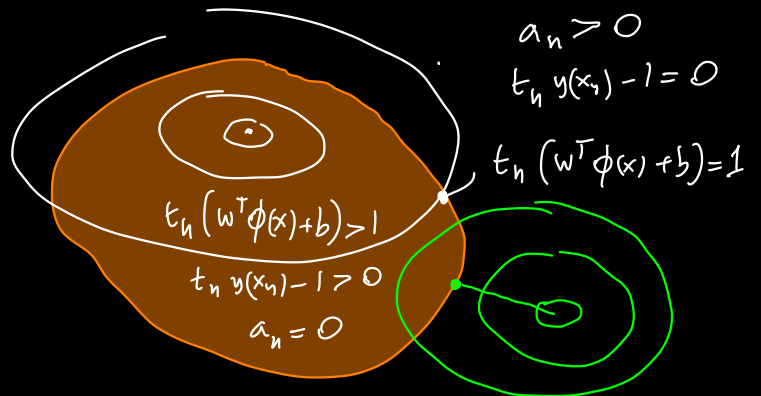
$$= \sum_{n \in N_S} a_n t_n K(x, x_n) + b$$

# Karush-Kuhn-Tucker KKT conditions:

$$a_n \geq 0$$

$$t_n \, y(x_n) - 1 \geq 0$$

$$a_n \{ t_n \, y(x_n) - 1 \} = 0$$

$a_n > 0$
$t_n \, y(x_n) - 1 = 0$

$t_n (w^T \phi(x) + b) = 1$

$t_n (w^T \phi(x) + b) > 1$
$t_n \, y(x_n) - 1 > 0$
$a_n = 0$



$a_n = 0$
$t_n \, y(x_n) > 1$

active constraints
$a_n > 0$

$t_n \, y(x_n) = 1$

SUPPORT VECTORS

$y(x) = 0$

if $x_n$ is support vector

$$\Rightarrow \quad 1 = t_n \, y(x_n) = t_n \left( \sum_{m \in S} a_m \, t_m \, K(x_n, x_m) + b \right)$$

$+1, -1$

$$t_n = t_n^2 \left( \phantom{\sum_{m \in S}} \right)$$

$\| \atop 1$

$N_s \, b$

$$\sum_{n \in S} t_n = \sum_{n \in S} \sum_{m \in S} a_m \, t_m \, K(x_n, x_m) + \sum_{n \in S} b$$

$$w = \sum_{n \in S} a_n t_n \phi(x_n)$$

$$b = \frac{1}{N_S} \sum_{n \in S} \left[ t_n - \sum_{m \in S} a_m t_m K(x_n, x_m) \right]$$