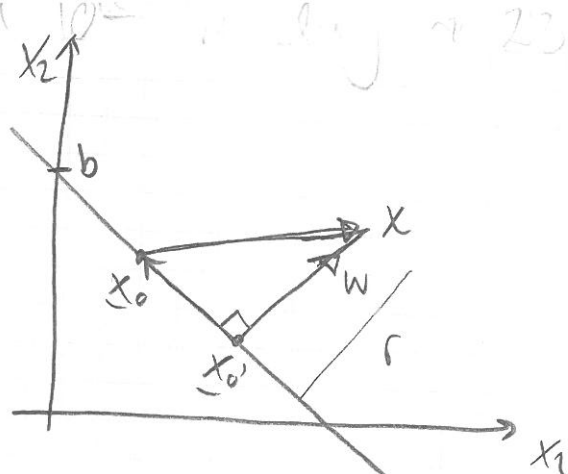


# SVM



$$L = \{x : x^T w + b = 0\}$$

i) For  $x_0 \in L$  har vi at  $x_0^T w = -b$

ii) For  $x'_0 \in L$  og  $x_0 \in L$  gælder:

$$(x_0 - x'_0)^T w = -b - (-b) = 0$$

Dette medfører at  $w$  er normalvektor til  $L$

ii') Specielt er  $w^* = \frac{w}{\|w\|}$  også normal (og her enheds længde)

iii) Afstanden (med fortegn) fra punkt  $x$  til  $L$  er givet ved skalar projektionen:

$$(x - x_0)^T w^* = (x - x_0)^T \frac{w}{\|w\|} = \frac{1}{\|w\|} (x^T w - x_0^T w)$$

idet  $x_0 \in L \rightarrow \frac{1}{\|w\|} (x^T w + b)$

Definer  $g(x) = x^T w + b$ . Da er  $g'(x) = \frac{\partial}{\partial x} g(x) = w$ .

$$\text{Da } (x - x_0)^T w^* = \frac{g(x)}{\|g'(x)\|} = r$$

Ydermere  $\rho = 2|r|$ . Som vi sawe kalder for marginen.

Antag vi har data således  $(Y, X)$  udgør data for  $i=1, \dots, n$ .

Vi kan vælge (for to klasser) at omkode  $y$  til

$y_i = 1$  for klasse 1

$y_i = -1$  for klasse 2.

### Ønske

Maximere afstande mellem hypereplan og punkter. således alle klasse 1 observationer er på én side af planen, mens alle klasse 2 obs. er på modsat side:

$$\max_{w, b, \|w\|=1} \rho \quad \text{således}$$

$$\begin{aligned} x_i^T w + b &\geq r \quad y_i = 1 \\ x_i^T w + b &\leq -r \quad y_i = -1 \end{aligned}$$

↑  
Husk  $r$   
med fortegn.

Bruger nu at  $y_i = 1$  eller  $y_i = -1$ :

$$\left. \begin{aligned} y_i = 1 & \quad x_i^T w + b \geq r \\ y_i = -1 & \quad x_i^T w + b \leq -r \end{aligned} \right\} \quad y_i (x_i^T w + b) \geq r$$

$$\text{Dvs} \quad \max_{w, b, \|w\|=1} \rho \quad \text{således} \quad y_i (x_i^T w + b) \geq r$$

~~Vælges  $\rho$  til~~ Som er ækvivalent med

$$\max_{w, b} \rho \quad \text{således} \quad \frac{1}{\|w\|} y_i (x_i^T w + b) \geq r$$

Vælg nu  $\|w\| = \frac{1}{r}$  for  $v_i$

3/9

$$\max_{\underline{w}, \underline{b}} \rho = 2/r = \frac{2}{\|w\|} \quad \text{således } y_i(x_i^T w + b) \geq 1$$

Dette problem har samme fix punkt (løsning) som

$$\min_{\underline{w}, \underline{b}} \frac{1}{2} \|w\|^2 \quad \text{således } y_i(x_i^T w + b) \geq 1, \\ \text{Ovs } 1 - y_i(x_i^T w + b) \leq 0$$

Sam er et Quadratic Programming Problem (QP)

For store dataset/mængder er QP dog ikke beregningsmæssigt hensigtsmæssig.

Løsning v Lagrange Multiplier eller rettere "Karush-Kuhn-Tucker" ~~metode~~ betragtes. // Wolfe dual

Def primære Lagrange funktion som skal minimeres mht.  $\underline{w}$  og  $\underline{b}$  e således:

burde være "+", men vi har vendt uligheder:  $y_i(x_i^T w + b) - 1 \leq 0$

$$L_P = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i(x_i^T w + b) - 1]$$

Differentier mht  $\underline{w}$  og  $\underline{b}$  giver det hhv (husk  $\|w\|^2 = w^T w$ )

$$\frac{\partial}{\partial \underline{w}} L_P = \underline{w} - \sum_{i=1}^n \alpha_i y_i x_i = \underline{0} \Rightarrow \underline{\hat{w}} = \sum_{i=1}^n \alpha_i y_i x_i$$
$$\frac{\partial}{\partial \underline{b}} L_P = - \sum_{i=1}^n \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$\alpha_i \geq 0$  for alle  $i=1, \dots, n$ .

Bilbetingelser  
ligninger.



Indsætter (8) i  $L_D$  får vi det såkaldte Wolfe dual

$$L_D = \frac{1}{2} \hat{W}^T \hat{W} - \sum_{i=1}^n \alpha_i [y_i (X_i^T \hat{W} + b) - 1]$$

$\alpha [y_i x_i^T w + y_i b - 1] = \alpha y_i x_i^T w + \alpha y_i b - \alpha$

$$= \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i x_i \right)^T \left( \sum_{j=1}^n \alpha_j y_j x_j \right) - \sum_{i=1}^n \alpha_i \left[ y_i (X_i^T \left\{ \sum_{j=1}^n \alpha_j y_j x_j \right\} + \underbrace{\sum_{j=1}^n \alpha_j y_j b}_{=0}) - 1 \right]$$

$$= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i \neq j} \alpha_i y_i \alpha_j y_j x_i^T x_j$$

$$- \underbrace{\left( \sum_{i=1}^n \alpha_i y_i \right) \left( \sum_{j=1}^n \alpha_j y_j \right)}_{=0} + \sum_{i=1}^n \alpha_i$$

$$= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

således  $\left\{ \begin{array}{l} \sum_{i=1}^n \alpha_i y_i = 0 \\ \alpha_i \geq 0 \end{array} \right.$

!!  $\left\{ \begin{array}{l} \alpha_i [y_i (X_i^T w + b) - 1] = 0 \\ \forall i=1, \dots, n \end{array} \right.$

Her af ser vi at (!!):

• Hvis  $\alpha_i > 0$  så må  $y_i (X_i^T w + b) = 1$ . Dvs at  $x_i$  ligger på randen/kanten af vores område

• Hvis  $y_i (X_i^T w + b) > 1$  så  $\alpha_i = 0$ .

husk:  $y_i (X_i^T w + b) \geq 1$

DERFOR: Punkter med  $\alpha_i > 0$  er således eneste bidrag til

$$\hat{W} = \sum_{i=1}^n \alpha_i y_i x_i$$

og kaldes derfor SUPPORT VECTORS

Bemærk

$\hat{g}(x) = x^T \hat{w} + \hat{b}$  giver antegnning til en  
~~klassifikator~~ :  
 klassifikator :  $g(x) = \text{sign}\{\hat{g}(x)\}$

Sidebemærkning:

Logistisk regression giver tilsvarende hyperplan i det separable tilfælde:

Fit :  $\text{glm}(\underbrace{\text{class}}_y \sim \cdot, \text{family} = \text{"binomial"})$

Planen er givet for  $P(\text{klasse } i | x) = 0.5$

$$0 = \logit 0.5 = \log \frac{0.5}{1-0.5} = \logit P(\text{klasse } i | x) = \beta^T x + \beta_0.$$

Så for  $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  har vi med  $\beta = (\beta_1, \beta_2)$ :

$$0 = \beta^T x + \beta_0 = \beta_1 x_1 + \beta_2 x_2 + \beta_0 \Rightarrow x_2 = -\frac{\beta_0 + \beta_1 x_1}{\beta_2}.$$

## Det ikke-separable tilfælde

6/9

En mulig løsning er at indføre "slack" variable, som tillader at alle observationer er på den korrekte side af hyperplanen.

"A few rotten apples in the basket",

sidebetingelser bliver nu:

$$y_i(x_i^T w + b) \geq 1 - \xi_i$$

hvor  $\xi_i \geq 0$  samt  $\sum_{i=1}^n \xi_i \leq C$ , hvor  $C$  er en "tuning-parameter". I R-implementation kaldes  $C$  "cost".

$\xi_i$  kan fortolkes som afstanden punkt  $i$  ligger på den forkerte side af  $x^T w + b = y_i$ .

Igen får vi et kvadratisk programmerings problem (QP)

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

hvilket

$$\xi_i \geq 0 \text{ og } y_i(x_i^T w + b) \geq 1 - \xi_i.$$

Note  
 $C = \infty$   $\approx$  separable

Det primære Lagrange funktional

$$L_p = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i(x_i^T w + b) - (1 - \xi_i)] - \sum_{i=1}^n \mu_i \xi_i$$

minimer mht  $w, b$  og  $\xi_i$



7/9

Differentierbar mht  $\underline{w}$ ,  $b$  og  $\underline{\xi}$  giv:

$$\frac{\partial}{\partial \underline{w}} L_P = \underline{w} - \sum_{i=1}^n \alpha_i y_i x_i = 0 \Rightarrow \hat{\underline{w}} = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial}{\partial b} L_P = - \sum_{i=1}^n \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial}{\partial \xi_i} L_P = C - \alpha_i - \mu_i = 0 \Rightarrow \alpha_i = C - \mu_i \quad i=1, \dots, n. \quad (**)$$

Indsættes (\*\*) samt  $\alpha_i \geq 0$ ,  $\mu_i \geq 0$  og  $\xi_i \geq 0$   $i=1, \dots, n$  i  $L_P$  får vi igen Wolfe dual:

$$L_D = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \quad \text{således at}$$

$0 \leq \alpha_i \leq C$  og  $\sum_{i=1}^n \alpha_i y_i = 0$ , hvor maximering af

$L_D$  også skal opfylde Karush-Kuhn-Tucker betingelser:

$$\left. \begin{aligned} \alpha_i [y_i (x_i^T \underline{w} + b) - (1 - \xi_i)] &= 0 \\ \mu_i \xi_i &= 0 \\ y_i (x_i^T \underline{w} + b) &\geq (1 - \xi_i) \end{aligned} \right\} i=1, \dots, n \quad (*)$$

Igen her vi ~~ikke~~ ~~ikke~~

•  $\alpha_i > 0$  må medføre  $y_i (x_i^T \underline{w} + b) = 1 - \xi_i$

Der på "random" vil  $\xi_i$  forskydes

~~alle  $\alpha_i \neq 0$~~   $\left\{ \begin{array}{l} \text{Nogle vil have } \xi_i = 0 \text{ og være på margin.} \\ \text{Der } 0 < \alpha_i < C \text{ for disse} \end{array} \right.$

$\left\{ \begin{array}{l} \text{For } \alpha_i > 0 \text{ og } \xi_i > 0 \text{ må } \mu_i = 0 \text{ så} \\ (*) \text{ og } (**) \text{ giver at } \alpha_i = C. \end{array} \right.$

Hermed ser vi at

8/9

$$L_0 = \sum x_i - \frac{1}{2} \sum_{i,j} x_i x_j y_i y_j x_i^T x_j$$

$$= \sum x_i - \frac{1}{2} \sum_{i,j} x_i x_j y_i y_j \langle x_i, x_j \rangle$$

↑ indek produkt.

### Kernel trick

Består i at afbillede  $x$  ind i højere dimensionalt rum  $H$  via  $\phi: X \rightarrow H$

Der

$$L_0 = \sum x_i - \frac{1}{2} \sum_{i,j} x_i x_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle$$

Der

g(x) som før var  $g(x) = x^T w + b$

nu bliver til

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

Her vi hellere fix  $\hat{w} = \sum_{i=1}^n x_i y_i x_i$  så har vi nu

$$\hat{w} = \sum_{i=1}^n x_i y_i \phi(x_i)$$

$$\tilde{g}(x) = \phi(x)^T \sum_{i=1}^n x_i y_i \phi(x_i) + b$$

$$= \sum_{i=1}^n x_i y_i \langle \phi(x), \phi(x_i) \rangle + b$$



Lad  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  da her vi <sup>9/9</sup>  
 altså

$$L_D = \sum_{i=1}^n x_i - \frac{1}{2} \sum_{i,j} x_i x_j y_i y_j K(x_i, x_j).$$

~~altså~~ altså

for d-grads polynomium:

$$K(x, x') = (1 + \langle x, x' \rangle)^d.$$

For  $d=2$ : og  $x = (x_1, x_2)$  og  $z = (z_1, z_2)$

$$\begin{aligned} (1 + \langle x, z \rangle)^d &= (1 + x_1 z_1 + x_2 z_2)^2 \\ &= (1 + (x_1 z_1)^2 + (x_2 z_2)^2 + 2x_1 z_1 + 2x_2 z_2 + 2x_1 x_2 z_1 z_2) \end{aligned}$$

Dette giver os at vi afbilder  $\phi: \mathbb{X} \rightarrow \mathbb{H}$ , hvor

$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6$ , hvor  $\phi(x) = (\phi_1(x), \dots, \phi_6(x))$  hvor

$$\phi_1(x) = 1$$

$$\phi_2(x) = \sqrt{2} x_1$$

$$\phi_3(x) = \sqrt{2} x_2$$

$$\phi_4(x) = x_1^2$$

$$\phi_5(x) = x_2^2$$

$$\phi_6(x) = \sqrt{2} x_1 x_2$$

$$K(\phi(x), \phi(z))$$

$$= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2.$$