

1 Úvod a motivace

Pro klasifikaci do dvou tříd můžeme použít klasifikátoru založeném na perceptronu. Jednou z nevýhod takového klasifikátoru je nalezení rozdělovací nadplochy bez dalšího kritéria, výsledek je v podstatě závislý na průběhu učení, jakým způsobem jsou vybírány prvky z trénovací množiny. V případě klasifikátorů SVM je tato rozdělovací nadplocha volena podle optimalizujícího kritéria.

V případě neseparabilních dat, lze použít stejného triku jako u učení klasifikátorů na principu perceptronu. Přejdem do vyšší dimenze prostoru, kde jsou data již separabilní.

Základní rozdělení SVM:

Lineární

- separabilní data
- neseparabilní data

Nelineární SVM

Další vlastností SVM je možnost vyčíslit chybu na testovací množině, pomocí „VC confidence“, z Vapnikovy teorie statistického učení.

2 Chyby klasifikátoru

Problém klasifikace do dvou tříd si můžeme představit jako nalezení funkce, která přiřadí naměřeným příznakům správnou třídu.

Množinu příznaků $\mathbf{x}_i \in R^n$, kde $i = 1, \dots, l$ je počet prvků v trénovací množině klasifikujeme do dvou tříd $-1, 1$. Pro každé \mathbf{x}_i známe y_i , což je třída, do které patří. Předpokládáme, že jednotlivá měření příznaků jsou nezávislá.

Úkolem učení klasifikátoru je nalezení transformační funkce z prostoru \mathbf{x}_i do prostoru y_i .

Úkolem učení klasifikátoru je nalezení transformační funkce z prostoru Transformace $\mathbf{x}_i \rightarrow y_i$ je závislá na obecné hodnotě parametru α . Tento parametr je závislý na volbě klasifikátoru a také na metodě učení klasifikátoru.

$$\mathbf{x}_i \rightarrow f(\mathbf{x}_i, \alpha)$$

Bez znalosti chyby klasifikátorů je jeho použití obtížné, neboť nevíme jak moc se můžeme na výsledky spolehnout.

Očekávanou chybu klasifikátoru na skutečných datech po naučení lze vyjádřit:

$$R(\alpha) = \int \frac{1}{2} |y - f(\mathbf{x}, \alpha)| dP(\mathbf{x}, y) \quad (1)$$

Jesliže existuje sdružená hustota pravděpodobnosti pak

$$dP(\mathbf{x}, y) = p(\mathbf{x}, y) d\mathbf{x} dy$$

$R(\alpha)$ z (1) se nazývá očekávané riziko. Tato hodnota uvádí skutečnou chybu, ale pravděpodobnosti $P(\mathbf{x}, y)$ nebo její odhad, lze získat velmi těžko.

$$R_{emp}(\alpha) = \frac{1}{2l} \sum_{i=1}^l |y_i - f(\mathbf{x}_i, \alpha)|. \quad (2)$$

$R_{emp}(\alpha)$ z (2) se nazývá empirická ztráta a udává chybu na trenovací množině, pro konečný počet pozorování. Hodnota $\frac{1}{2} |y_i - f(\mathbf{x}_i, \alpha)|$ se nazývá ztráta. Protože klasifikaci provádíme do tříd $-1, 1$ může nabývat pouze hodnot 0 nebo 1. Zvolíme-li η takové, že $0 \leq \eta \leq 1$, nabývá ztráta hodnoty η s pravděpodobností $1 - \eta$.

Samotná chyba na trenovací množině mnoho nevypovídá o chybě na testovacích (reálných) datech. Následující vztah dává odhad skutečné chyby větší než vztah (1), zato jej můžeme vyčíslit aniž bychom znali $P(\mathbf{x}, y)$

$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\left(\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l} \right)} \quad (3)$$

Kde h je celé nezáporné číslo nazývané VC (Vapnik-Chervonenkis) dimenze. Odmocnina ve vztahu (3) se nazývá „VC confidence“ (Vapnik 1995).

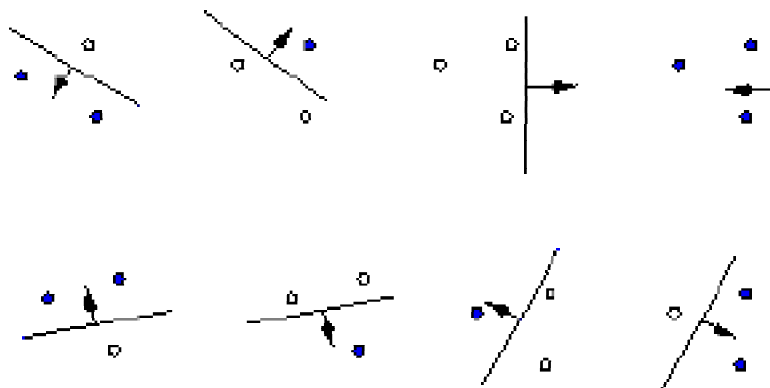
2.1 VC dimenze

VC dimenze je vlastnost množiny funkcí $f(\alpha)$; α představuje obecné parametry v závislosti na volbě funkce f . Pro klasifikaci do dvou tříd je VC dimenze definována jako maximální počet bodů, jež lze rozbít funkcí $f(\alpha)$.

Mějme množinu l bodů, kterým přiřazujeme označení $-1, 1$ existuje 2^l možných přiřazení (označení).

Jestliže pro každé označení l bodů existuje funkce z $f(\alpha)$ pak $f(\alpha)$ rozbíjí množinu bodů l .

VC dimenze h zaručuje, že existuje alespoň jedna množina h bodů, které lze rozbít.



Obrázek 1: příklad rozbití množiny 3 bodů v rovině

2.2 VC dimenze a počet parametrů

Mohlo by se zdát, že hodnota VC dimenze závisí na počtu parametrů α rozbíjející funkce. Následující příklad ukáže funkci s jediným parametrem α a s hodnotu VC dimenze nekonečno.

Mějme prahovou funkci :

$$\theta(x), x \in (R) : \theta(x) = 1 \forall x > 0; \theta(x) = -1 \forall x \leq 0$$

a jedno parametrickou funkci

$$f(x, \alpha) = \theta(\sin(\alpha x)), x, \alpha \in \mathbf{R}.$$

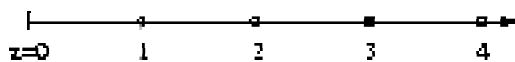
Zvolíme nějaká čísla l a současně najdeme l bodů, které mohou být rozbity. Body vybereme podle předpisů:

$$x_i = 10^{-i}, i = 1, \dots, l.$$

Bodům přiřadíme třídy $y_1, y_2, \dots, y_l; y_i \in -1, 1$.

$f(x, \alpha)$ přiřazující bodům příslušnost ke třídě, a pro $\alpha = \pi(1 + \sum_{i=1}^l \frac{(1-y_i)10^i}{2})$.

Takovýto klasifikátor rozbíjí nekonečně mnoho bodů a má tedy VC dimenzi nekonečno. Jistě lze najít množinu čísel, která touto funkcí nelze rozbít, bez ohledu na nekonečnou VC dimenzi. Ta nám však zaručuje, že existje alespoň jedna množina bodů jež lze rozbít.



Obrázek 2: příklad množiny čísel jež nelze rozbít funkcí $f(x, \alpha) = \theta(\sin(\alpha x)), x, \alpha \in \mathbf{R}$.

2.3 VC dimenze orientované nadroviny

Pro m bodů z \mathbf{R}^n zvolme jeden bod za počátek souřadnic. m bodů lze rozbít orientovanými nadrovinami právě když zbylých $m-1$ bodů je lineárně nezávislých.

VC dimenze orientovaných nadrovin v \mathbf{R}^n je $n+1$.

2.4 VC dimenze orientované nadroviny, separabilní data

Algoritmus SVM je založen na struktuře množin funkcí rozdělujících nadrovin. Popišme si množiny v prostoru Z . Množina vektorů $\mathbf{x}_1, \dots, \mathbf{x}_r \in Z$.

Každá nadrovina $\mathbf{x} \in Z : (\mathbf{w} \cdot \mathbf{x} + b = 0)$ odpovídá dvojici $(\mathbf{w}, b) \in \mathbf{X} \times \mathbf{R}$.

Mějme podmínku na charakteristickou vzdálenost bodů $\mathbf{x}_1, \dots, \mathbf{x}_r$ od rozdělovací nadroviny

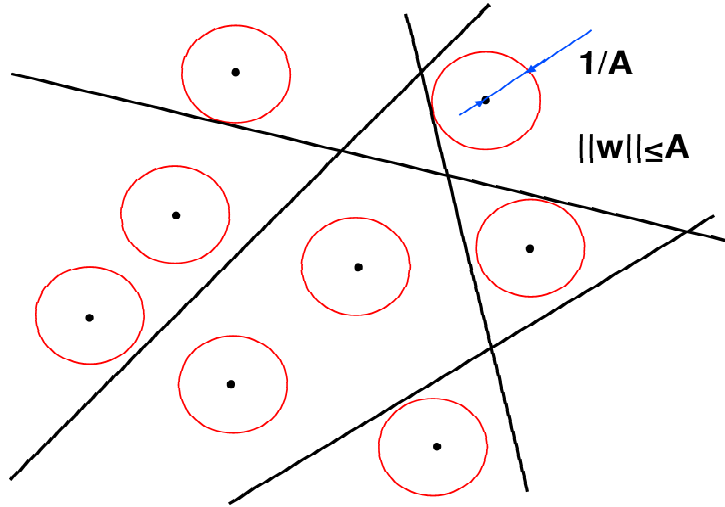
$$\max_{i=1, \dots, r} |(\mathbf{w} \cdot \mathbf{x}_i) + b| = 1 \quad (4)$$

Nejmenší koule obsahující body $\mathbf{x}_1, \dots, \mathbf{x}_r$, $B_{\mathbf{x}_1, \dots, \mathbf{x}_r} = \{\mathbf{x} \in Z : |(\mathbf{w} - \mathbf{a})| < R\}$ $\mathbf{a} \in Z$ a rozhodovací funkce definovaná na těchto bodech

$$f_{\mathbf{w}, b} = \text{sgn}((\mathbf{w} \cdot \mathbf{x}) + b) \quad (5)$$

množina $f_{\mathbf{w}, b} : |\mathbf{w}| \leq A$ má VC-dimenzi h (Vapnik, 1995) :

$$h \leq R^2 A^2. \quad (6)$$

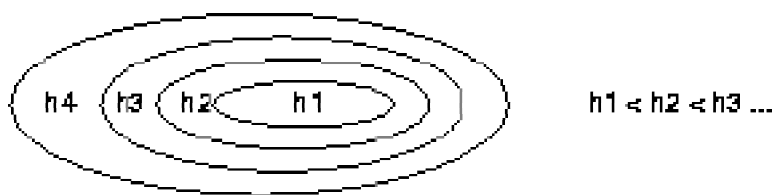


Obrázek 3: rozdělující nadroviny a koule obsahující body

2.5 Minimalizace strukturálního rizika

VC confidence v rovnici (3) závisí na volbě klasifikační funkce. Měli bychom vybrat takovou podmnožinu z množiny funkcí, pro kterou bude odhad chyby minimální, pro každou podmnožinu zjistíme její VC dimenzi. Minimalizace strukturálního rizika je hledání takové podmnožiny funkcí, aby hodnota empirického riziku a hodnota VC confidence byla minimální.

Na obrázku 4 jsou zobrazeny podmnožiny s rostoucí VC dimenzí.

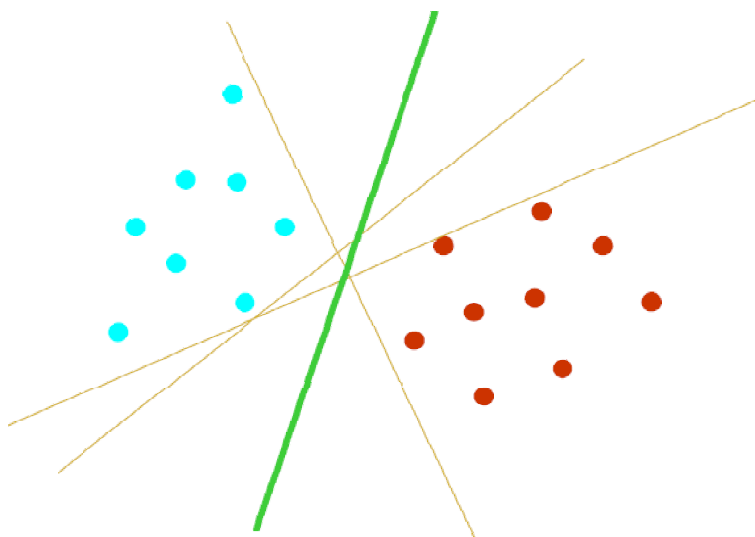


Obrázek 4: podmnožiny funkcí uspořádané podle VC dimenze

3 Separabilní data

Problém klasifikace je omezen na dvě třídy, jinak je obecný. Podívejme se na příklad na obr.5. Zde je mnoho schopných lineárních klasifikátorů, které umí rozdělit data, ale jen jeden maximalizuje „mezeru“ (tedy vzdálenost mezi nejbližšími body z obou tříd k rozdělující nadrovině). Tento lineární klasifikátor bývá nazýván optimální rozdělující nadrovinou.

Tedy hledáme lineární klasifikátor, jehož učení je založeno na minimalizaci strukturálního rizika, nalezením rozdělujícího pásu (mezery) maximální šířky.



Obrázek 5: Optimální rozdělující nadrovina

3.1 Optimální rozdělující nadrovina

Uvažujeme trénovací množinu ve tvaru $\{\mathbf{x}_i, y_i\}$, $i=1, \dots, l$, $y_i \in \{-1, 1\}$, $\mathbf{x}_i \in \mathbb{R}^d$, kde y_i je třída (pozitivní a negativní vzorky) a \mathbf{x}_i je vektor souřadnic d -dimenzionálního prostoru. Předpokládejme, že máme nějakou nadrovinu, která nám odděluje pozitivní vzorky od negativních (oddělující nadrovinu). Body (vzorky) \mathbf{x} ležící na oddělující nadrovině vyhovují rovnici $\mathbf{w} \cdot \mathbf{x} + b = 0$, kde \mathbf{w} je normálový vektor oddělující nadroviny (též si lze představit jako vektor vah) a $|b|/|\mathbf{w}|$ je vzdálenost od nadroviny k počátku souřadné soustavy, viz. obr.6.

Nechť $d_+(d_-)$ jsou nejkratší vzdálenosti od oddělující nadroviny k pozitivním (negativním) bodům.

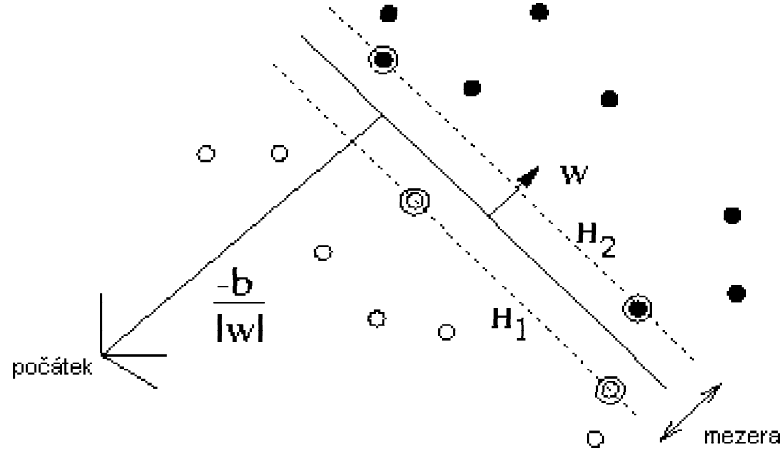
Definujme mezeru oddělující nadroviny: $\rho(\mathbf{w}, b) = d_+ + d_-$. Úlohu můžeme formulovat následovně: předpokládejme, že všechna trénovací data vyhovují následujícím podmínkám:

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \quad \text{pro } y_i = +1 \quad (7)$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \quad \text{pro } y_i = -1 \quad (8)$$

Tyto dvě podmínky lze též formulovat následovně:

$$y_i \cdot (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad \forall i \quad (9)$$



Obrázek 6: Lineární rozdělovací nadrovina pro dvoudimenzionální prostor

Tedy vyrobíme dva pásy o šířce 1, a změnou parametru \mathbf{w} budeme měnit šířku mezery, neboli měřítko námi zvolené soustavy. Je vhodné, aby všechna data z trénovací množiny spadala do tohoto měřítka. Nyní vezmeme v úvahu body, pro které platí podmínka (7).

Tyto body jsou odděleny nadrovinou H_1 :

$$\mathbf{x}_i \cdot \mathbf{w} + b = 1 \quad (10)$$

s normálou \mathbf{w} a vzdáleností od počátku $|1 - b| / |\mathbf{w}|$. Podobně pro body splňující podmínku (8) platí pro H_2 :

$$\mathbf{x}_i \cdot \mathbf{w} + b = -1 \quad (11)$$

Kde je vzdálenost od počátku rovna $|-1 - b| / |\mathbf{w}|$. Poznamenejme, že H_1 a H_2 jsou paralelní (mají stejné normály), a že mezi nimi neleží žádná trénovací data. Body \mathbf{x}_i splňující podmínky (9) resp. (10) leží na nadrovinách H_1 resp. H_2 se nazývají **support vectors** a jsou v obr.6 označeny kroužkem. Tedy hledáme body \mathbf{x}_i , které vyhovují podmínce

$$\min_{\mathbf{x}_i} |\mathbf{x} \cdot \mathbf{w} + b| = 1 \quad (12)$$

Vzdálenost $d(\mathbf{w}, b; \mathbf{x})$ od \mathbf{x} k nadrovině \mathbf{w}, b je

$$d(\mathbf{w}, b; \mathbf{x}) = \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{|\mathbf{w}|} \quad (13)$$

Tedy optimální oddělující nadrovina maximalizující mezeru nám dává velikost mezery

$$\begin{aligned} \rho(\mathbf{w}, b) &= d_+ + d_- = \min_{\{\mathbf{x}_i: y_i=1\}} d(\mathbf{w}, b; \mathbf{x}_i) + \min_{\{\mathbf{x}_i: y_i=-1\}} d(\mathbf{w}, b; \mathbf{x}_i) \\ &= \min_{\{\mathbf{x}_i: y_i=1\}} \frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{|\mathbf{w}|} + \min_{\{\mathbf{x}_i: y_i=-1\}} \frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{|\mathbf{w}|} \\ &= \frac{1}{|\mathbf{w}|} \left(\min_{\{\mathbf{x}_i: y_i=1\}} |\mathbf{w} \cdot \mathbf{x}_i + b| + \min_{\{\mathbf{x}_i: y_i=-1\}} |\mathbf{w} \cdot \mathbf{x}_i + b| \right) \\ &= \frac{2}{|\mathbf{w}|} \end{aligned} \quad (14)$$

Chceme-li maximalizovat mezeru musíme maximalizovat výraz $2/|\mathbf{w}|$ nebo-li minimalizovat $|\mathbf{w}|$. Jak je z výše uvedených vztahů vidět, je potřeba pro nalezení optimální rozdělující nadroviny nalézt vázaný extrém, který je funkcí \mathbf{w}, b , pro již určené \mathbf{x} jako support vectors. Použijeme minimalizačního kritéria Lagrangeových koeficientů pro nalezení optimální oddělující nadroviny:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} |\mathbf{w}|^2 - \sum_{i=1}^l \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^l \alpha_i \quad (15)$$

Kde α_i jsou Lagrangeovy koeficienty (multiplikátory), pro které platí $\alpha_i > 0 \Leftrightarrow \mathbf{x}_i$ je „support vector“.

Teď si ještě naznačme řešení Lagrangeovy rovnice, předem je nutno podotknout, že se jedná o problém kvadratického programování. Hledáme minimum funkce $L(\mathbf{w}, b, \alpha)$ s ohledem ke koeficientům α_i . Tento problém bývá nazýván Wolfe dual (primárním problémem je označována rovnice (15) a sekundárním rovnice (16) viz. níže).

$$\max_{\alpha} W(\alpha) = \max_{\alpha} \left\{ \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) \right\} \quad (16)$$

Minimum $\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$ určíme pomocí

$$\frac{\partial L}{\partial b} = 0 \quad \Rightarrow \quad \sum_{i=1}^l \alpha_i y_i = 0 \quad (17)$$

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i \quad (18)$$

Ze vztahů (16), (17) a (18) plyne

$$\max_{\alpha} W(\alpha) = \max_{\alpha} -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^l \alpha_i \quad (19)$$

Řešením rovnice (19) je

$$\bar{\alpha} = \arg \min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \alpha_i \quad (20)$$

za podmínek (17) a $\alpha_i \geq 0, \quad i = 1, \dots, l$.

4 Neseparabilní data

Pro neseparabilní data výše popsany algoritmus nenajde vhodné řešení. Proto je třeba ho upravit tak, aby fungoval i pro neseparabilní data. Uvolníme proto omezení (7), (8) zavedením odchylky

$$\xi_i, \quad i = 1 \dots l$$

Omezení definující oddělovací nadrovinu potom přejdou v

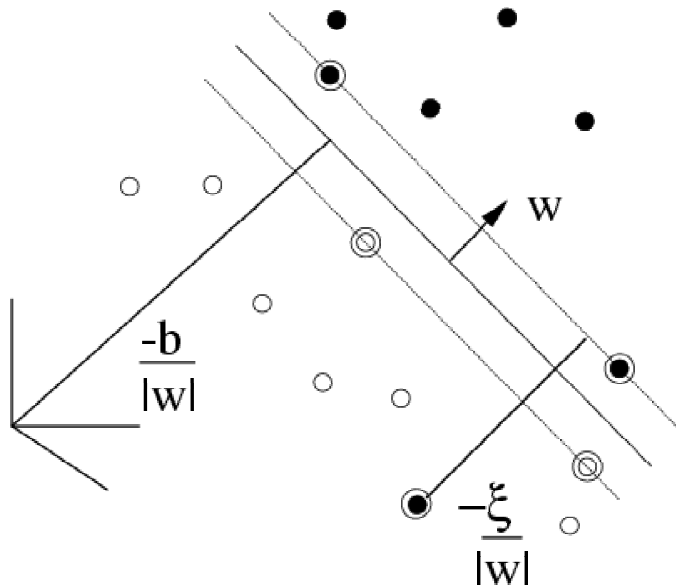
$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 - \xi_i \quad \text{pro } y_i = +1 \quad (21)$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i \quad \text{pro } y_i = -1 \quad (22)$$

$$\xi_i \geq 0 \quad \forall i. \quad (23)$$

Tím rozšíříme prostor příznaků a zároveň výraz $\sum_i \xi_i$ je maximum chyby na trénovacích datech. Potom místo hledání minima $\frac{\|\mathbf{w}\|^2}{2}$ pro separabilní data budeme hledat minimum, kde parametr C volí uživatel a odpovídá pokutě za danou chybu.

$$\frac{\|\mathbf{w}\|^2}{2} + C \left(\sum_i \xi_i \right)^k \quad (24)$$



Obrázek 7: Lineární rozděluující nadroviny pro neseparabilní data

Vyřešení nerovnic (21) až (24) je konvexní minimalizační problém matematického programování. Pro $k=1$, $k=2$ přechází v problém kvadratického programování. Pro $k=1$ je výhodné, že ani ξ stejně jako Langrangeovy multiplikátory se neobjeví ve Wolfově duálním problému maximalizace

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad (25)$$

$$0 \leq \alpha_i \leq C \quad (26)$$

$$\sum_i \alpha_i y_i = 0 \quad (27)$$

Řešením je:

$$\mathbf{w} = \sum_{i=1}^{N_s} \alpha_i y_i \mathbf{x}_i, \quad (28)$$

kde N_s je počet support vectors.

Další možná řešení lze najít pomocí:

- constrained confugate gradient descent

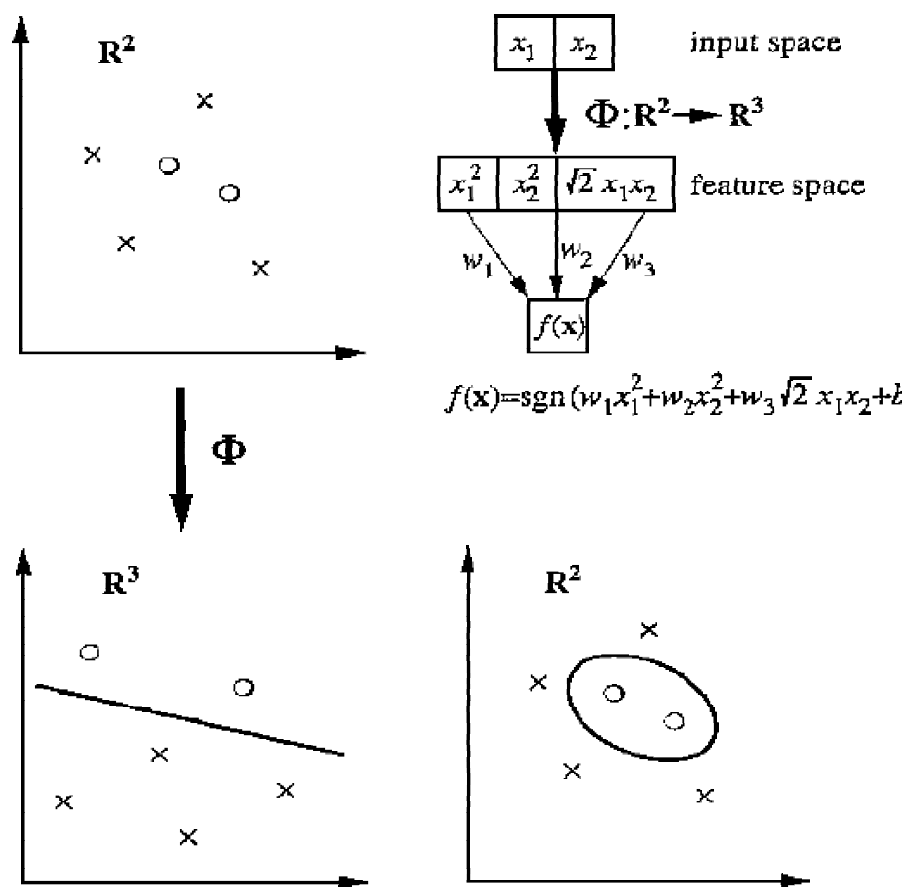
- interior point methods
- projection methods
- Bunhc Kaufman decomposition

5 Nelineární SVM

Nelineární SVM musíme použít v případě, když rozhodovací funkce není lineární. Boser, Guyon a Vapnik ukázali, že cesta vede přes vlastnosti skalárního součinu $\mathbf{x}_i \cdot \mathbf{y}_i$. Trénovací data jednoduše přetransformujeme z prostoru R^d do \mathcal{H} pomocí zobrazení

$$\Phi : R^d \rightarrow \mathcal{H} \quad (29)$$

Řešení tedy vede na hledání v prostoru vyšší dimenze, kde jsou data separabilní a kde tedy můžeme použít postup z předchozího odstavce.



Obrázek 8: Transformace trénovacích dat (vlevo nahoře) nelineárně přes Φ do vícedimenziálního prostoru \mathbb{R}^3 a konstrukce oddělující nadroviny (dole vlevo), SV klasifikátor (vpravo nahoře) odpovídající nelineární rozhodovací oblasti ve vstupním prostoru (vpravo dole).

Trénovací algoritmus závisí pouze na datech ze skalárního součinu v \mathcal{H} (tj. $\Phi(x_i) \cdot \Phi(x_j)$). Pokud tedy existuje Kernel funkce:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \quad (30)$$

nemusíme funkci Φ počítat ani ji nemusíme znát. Ne vždy lze jednoduše určit, zda existuje pár $\{\mathcal{H}, \Phi\}$. Odpověď lze určit pomocí Mercerovy podmínky:

existuje

$$K(\mathbf{x}_i, \mathbf{y}_i) = \sum_i \Phi(\mathbf{x})_i \Phi(\mathbf{y})_i \quad (31)$$

a jakákoliv funkce $g(x)$ taková, že integrál 32 je konečný.

$$\int g(\mathbf{x})^2 d\mathbf{x} \quad (32)$$

potom pár $\{\mathcal{H}, \Phi\}$ existuje pouze když platí

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x})^2 g(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0 \quad (33)$$

Příklad: (ukázka povoleného kernelu, pro který můžeme určit Φ)
Trénovací data jsou v dimenzi R^2 a jako kernel vyberem $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i, \mathbf{x}_j)^2$.
Potom není těžké najít prostor $\mathcal{H} = \mathcal{R}^3$ a transformaci Φ z R^2 do \mathcal{H} takovou, že

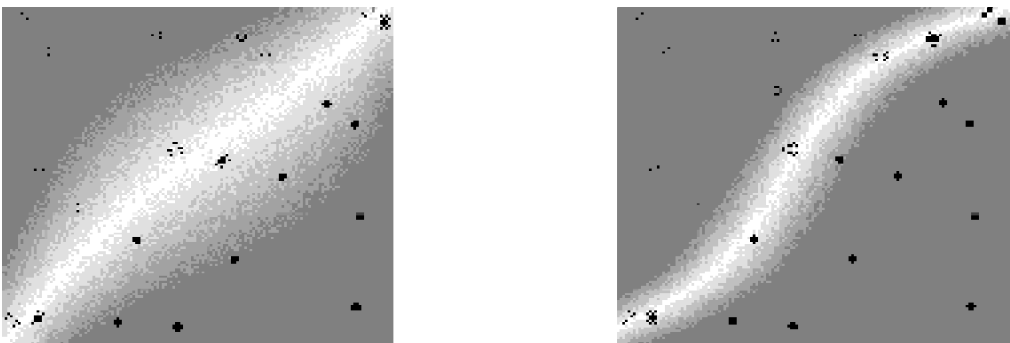
$$(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y}) \text{ a}$$

$$\Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

Rozhodovací funkce klasifikátoru pro nelineární SVM se počítá dle

$$f(\mathbf{x}) = \sum_{i=1}^{N_s} \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}_i) + b, \quad (34)$$

kde \mathbf{s}_i jsou support vectors.



Obrázek 9: Příklad nelineární SVM, kde kernel je kubický polynom 3. stupně

6 Generalizace SVM

- VC dimenze množiny funkcí rozdělujících nadrovin je

$$n + 1$$

- podmnožina orientovaných nadrovin

$$f(x, \omega, b) = \text{sign}(\omega \cdot x) + b$$

definovaná na $\mathbf{X}^r = x_1, x_2, \dots, x_r; |x_i - a| \leq R$ a splňující $|\omega| \leq A$ má VC dimenzi omezenou nerovností

$$h \leq \min(R^2 A^2)$$

výše dva zmíněné případy lze shrnout a omezit VC dimenzi:

$$h \leq \min(R^2 A^2, n) + 1 \quad (35)$$

- pro separabilní data (Vapnik, 1995) přináší druhou možnost jak odhadnout skutečnou chybu klasifikátoru SVM

$$E[P(\text{error})] \leq \frac{E[\text{number of SUPPORT VECTORS}]}{\text{Number of training vectors}} \quad (36)$$

kde $P(\text{error})$ je hodnota skutečného rizika, klasifikátoru naučeného z $l - 1$ vzoru. $E[P(\text{error})]$ očekávání skutečného rizika přes všechny možnosti trénovací množiny $l - 1$. $E[\text{Number of support vectors}]$ je očekávání z počtu support vectors přes všechny možnosti trénovacích množin velikosti l .

7 Závěr

Metoda SVM má velmi jednoduchou geometrickou interpretaci. Učení SVM klasifikátoru vždy nalezne globální minimum.