

UCZENIE MASZYNOWE III - SVM

mgr inż. Adam Kupryjanow

Plan wykładu

- Wprowadzenie
- LSVM – dane separowalne liniowo
- SVM – dane nieseparowalne liniowo
- Nieliniowy SVM
- „Kernel trick”
- Przykłady zastosowań

Historia

- 1992 wprowadzony przez Boser, Guyon & Vapnik
- Algorytm z mocnymi podstawami teoretycznymi, wywodzący się ze statystyki
- Teoria uczenia (Vapnik & Chervonenkis) z lat 60tych
- Dobrze przebadany eksperymentalnie i zastosowany w wielu dziedzinach: bioinformatyka, rozpoznawanie tekstu/obrazu, ...

Dane liniowo separowalne

- Dane $\{\mathbf{x}_i\}$, $i = 1, \dots, l$, $\mathbf{x}_i \in R^d$ należące do dwóch klas określonych zmiennymi $\{y_i\}$, $i = 1, \dots, l$ są liniowo separowalne, jeśli istnieje hiperpłaszczyzna \mathcal{H} postaci $g(\mathbf{x})$:

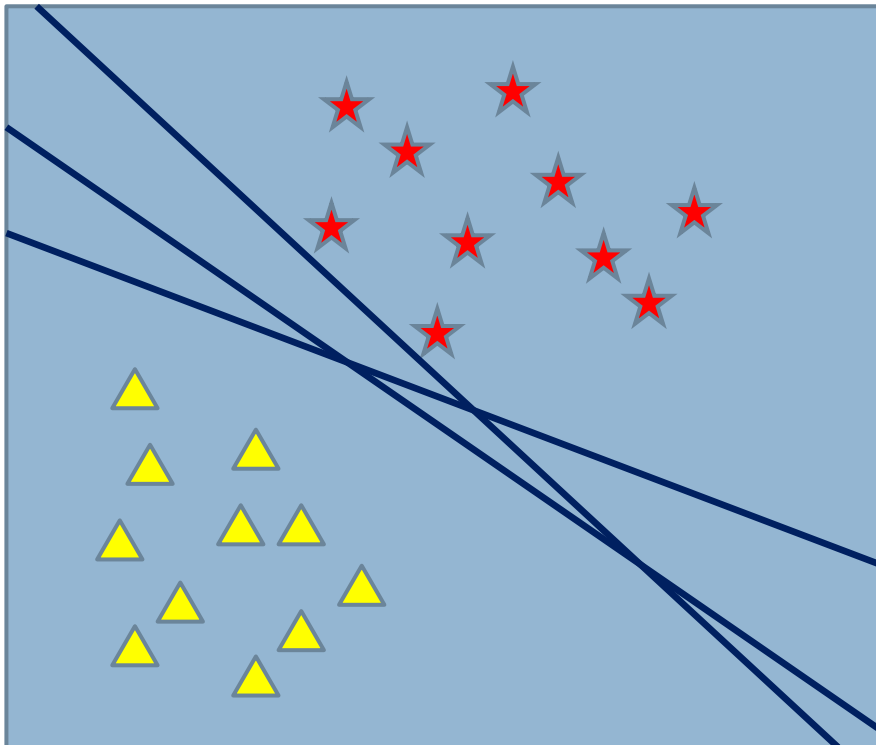
$$\mathcal{H}: g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + b$$

przyjmująca wartości ($i = 1, \dots, n$):

$$\begin{cases} g(\mathbf{x}_i) > 0 \text{ dla } \mathbf{x}_i \in w_1 \\ g(\mathbf{x}_i) < 0 \text{ dla } \mathbf{x}_i \in w_2 \end{cases}$$

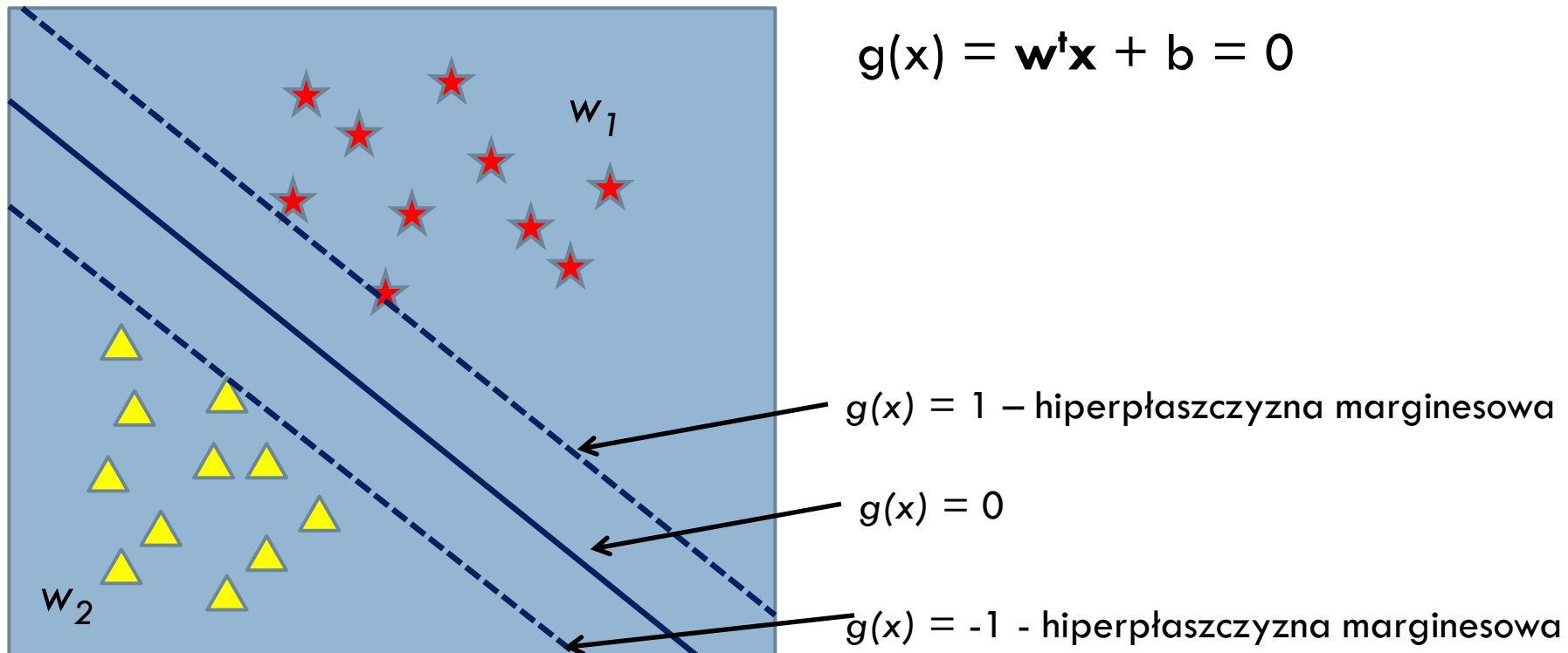
Jak wyznaczyć hiperpłaszczyznę?

Istnieje nieskończenie wiele funkcji rozdzielających dwie klasy



Liniowa maszyna wektorów nośnych LSVM

- Pomysł Vapnika metoda SVM (wektory nośne)
- Algorytm wyznaczanie hiperpłaszczyzny dąży do uzyskania jak największego marginesu



Wyznaczanie marginesu

- Można wykazać, że maksymalna odległość pomiędzy marginesami $\mathbf{w}^t \mathbf{x} + b = 1$ i $\mathbf{w}^t \mathbf{x} + b = -1$ wynosi $M = \frac{2}{\|\mathbf{w}\|}$, gdzie $\|\mathbf{w}\| = \sqrt{w_1^2 + w_2^2 + \dots + w_i^2}$
- Rozwiązanie powinno dążyć do uzyskania jak najkrótszego wektora \mathbf{w} , ponieważ wtedy uzyskany margines będzie największy
- Postulaty:
 - ▣ minimalizować wektor \mathbf{w} -> największy margines
 - ▣ próbki punktów uczących dla funkcji decyzyjnej mają dawać wartości ≥ 1 lub ≤ -1

LSVM - zagadnienie optymalizacji

- minimalizować po w wyrażenie $\frac{1}{2} \|w\|^2$
- przy warunku ograniczającym:

$$y_i \cdot ((w \cdot x_i) + b) \geq 1, \forall_i \in U$$

- Powyższe warunki prowadzą do uogólnionego równania Lagrange'a

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^l \alpha_i [y_i ((x_i \cdot w) + b) - 1]$$

gdzie α_i to mnożnik Lagrange'a i $\alpha_i > 0$

LSVM - zagadnienie optymalizacji

- Równanie Lagrange'a powinno różniczkować się po w i b
- Porównując pochodne $L(w,b,\alpha)$ względem w i b do zera otrzymujemy:

$$w = \sum_{i=1}^l \alpha_i y_i x_i \text{ oraz } \sum_{i=1}^l \alpha_i y_i = 0$$

- Podstawiając otrzymane wartości do równania otrzymujemy funkcję decyzyjną $g(x)$:

$$g(x) = \operatorname{sgn} \left(\sum_{i=1}^l \alpha_i y_i \cdot (x_i \cdot x) + b \right)$$

LSVM - zagadnienie optymalizacji

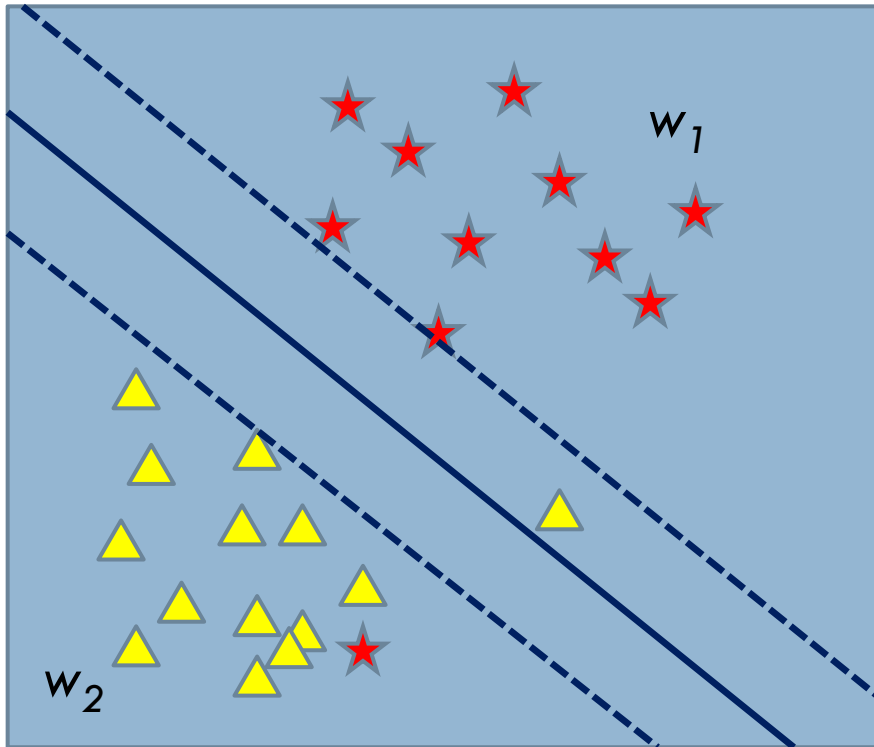
- Funkcja $g(x)$ zależy bezpośrednio od mnożników Lagrange'a (α_i)
- Mnożniki na podstawie twierdzenia KKT (Karush-Kuhn-Tucker) powinny spełniać warunek:

$$\sum_{i=1}^l \alpha_i [y_i((x_i \cdot w) + b) - 1] = 0, i = 1, \dots, l$$

Możliwe gdy:

- x_i leży na marginesie $\rightarrow \alpha_i$ dowolne
- x_i leży poza marginesami $\rightarrow \alpha_i$ równe zero

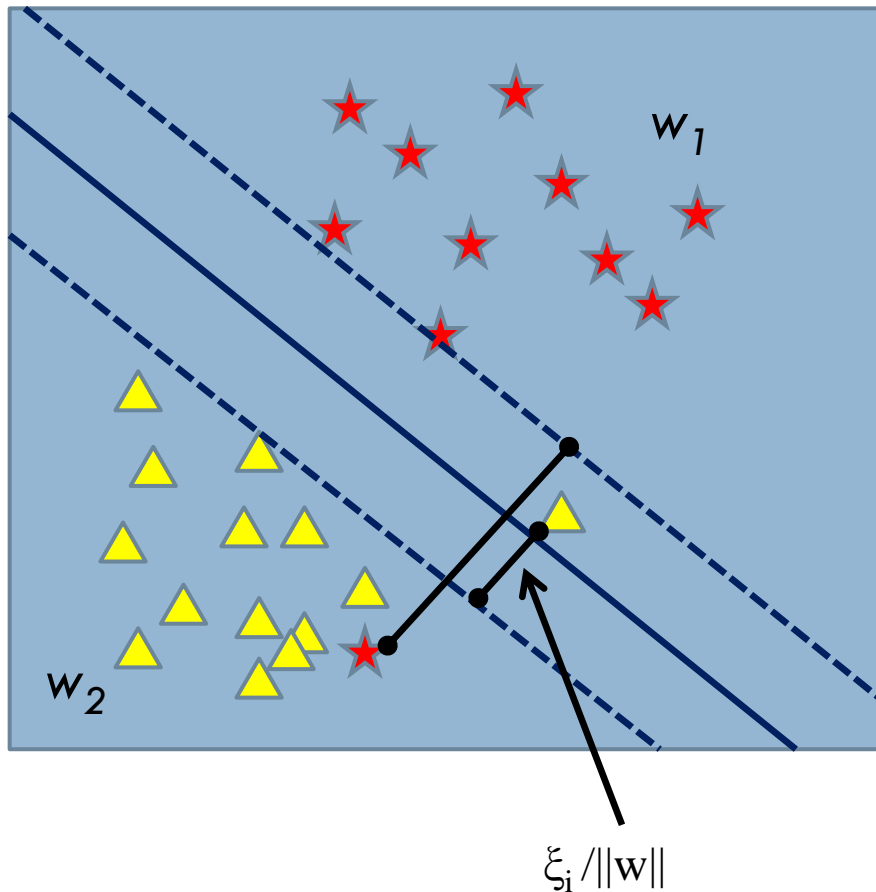
Dane liniowo nieseparowalne



- Cortes i Vapnik wykazali, że możliwe jest zastosowanie wektorów nośnych stosując pewne rozszerzenie założeń

$$y_i \cdot ((w \cdot x_i) + b) \geq 1, \forall_i \in U \longrightarrow y_i \cdot ((w \cdot x_i) + b) \geq 1 - \xi_i, \forall_i \in U$$

Dane liniowo nieseparowalne



- Nowa zmienna ξ_i nazywana „zwisem” (slack variable)
- Wartość ξ_i powinna być „mała”, aby ją określić rozpatrujemy:

$$C \sum_i \xi_i$$

Gdzie C to parametr generalizujący deklarowany przez użytkownika

- Jeżeli $0 \leq \xi_i \leq 1$ to punkt danych leży wewnątrz strefy separującej, po właściwej stronie
- Jeżeli $\xi_i > 1$, punkt po niewłaściwej stronie hiperpłaszczyzny = błąd klasyfikacji

Zagadnienie optymalizacji -SVM

- minimalizuj wyrażenie: $\frac{1}{2} \|w\|^2 + C \sum_i^l \xi_i$
- przy warunkach:

$$y_i \cdot ((w \cdot x_i) + b) \geq 1 - \xi_i, \forall_i \in U \quad \xi_i \geq 0, \forall_i$$

- Otrzymujemy Lagrangian:

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i \{ [y_i ((x_i \cdot w) + b)] - 1 + \xi_i \} - \sum_{i=1}^l \beta_i \xi_i$$

- Wyznaczamy pochodne cząstkowe względem w , b i ξ , i podstawiamy otrzymane wartości do Lagrangianu. Z warunku zerowania pochodnej Lagrangianu względem ξ otrzymujemy:

$$\alpha_i + \beta_i = C, \forall_i$$

Zagadnienie optymalizacji -SVM

- W tej sytuacji można wyeliminować β i zastąpić przez α
- Do rozwiązania pozostaje problem dualny:

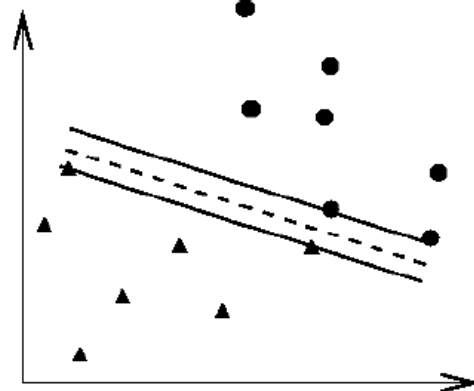
$$\max_{\alpha} W(\alpha) = \max_{\alpha} \left(\alpha \cdot 1 - \frac{1}{2} \alpha \cdot D\alpha \right)$$

gdzie, $D_{ij} = y_i y_j \cdot (x_i \cdot x_j)$

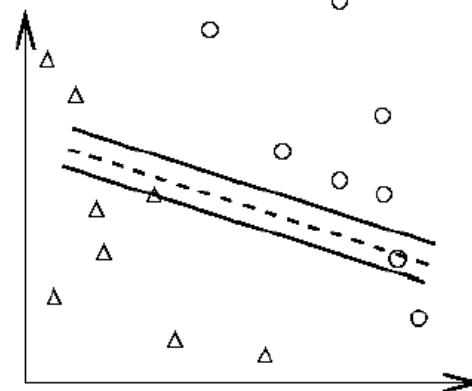
z ograniczeniami: $\alpha \cdot \beta = 0$, oraz $0 \leq \alpha_i \leq C$

Jakiego C używać?

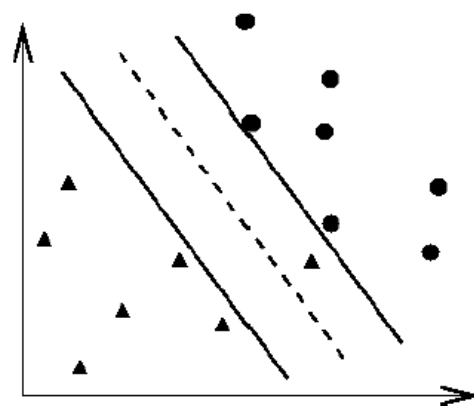
- Blanz i Vapnik zalecali stosowania $C = 5$, ale tak naprawdę C należy dobierać w zależności od danych uczących.



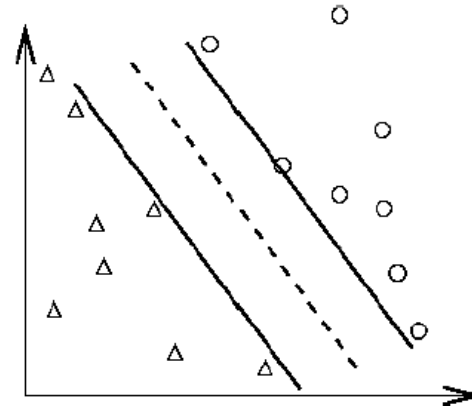
(a) Training data and an overfitting classifier



(b) Applying an overfitting classifier on testing data



(c) Training data and a better classifier



(d) Applying a better classifier on testing data

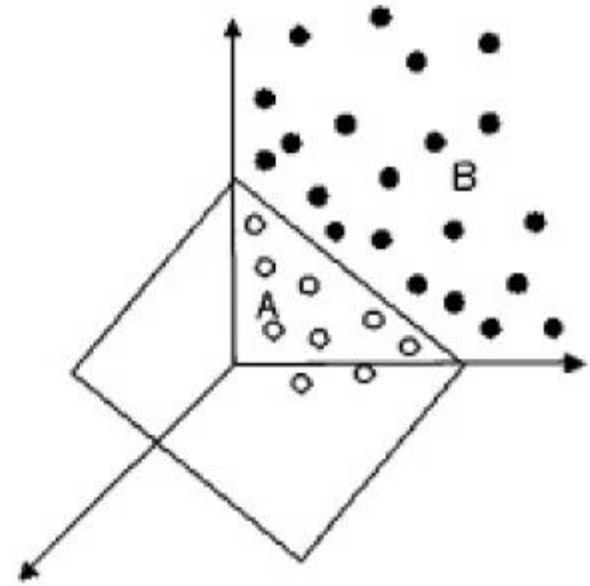
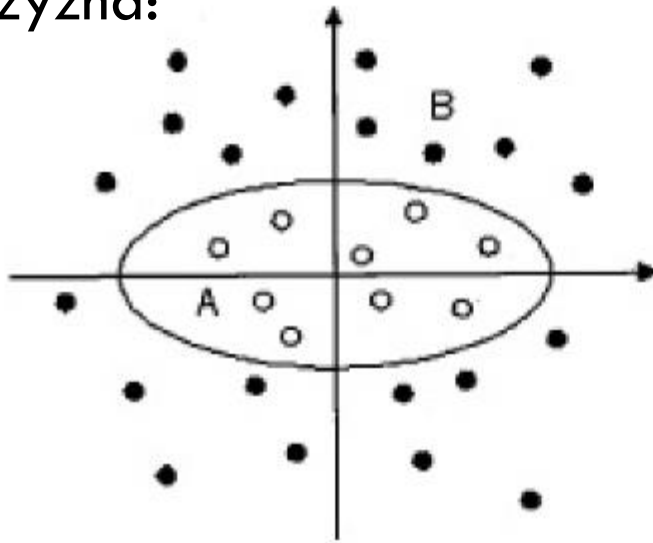
Nieliniowy SVM

- Transformacja do przestrzeni o wyższym wymiarze
- Projekcja danych oryginalnych $x \in \mathbb{R}^d$ do przestrzeni wymiarowej $n > d$ w której dane z dużym prawdopodobieństwem będą separowane liniowo

Przykład

- Mamy nieliniową funkcję mapującą $\varphi: I=\mathbb{R}^2 \rightarrow F=\mathbb{R}^3$ 2-wymiarową przestrzeń wejściową (input space) do 3-wymiarowej przestrzeni zmiennych przekształconych (feature space)
- $(x_1, x_2) \rightarrow (z_1, z_2, z_3) := (x_1^2, 2^{0.5} x_1 x_2, x_2)$
- hipertaszczyzna:

$$\langle w \cdot x \rangle = 0$$



Model nieliniowy SVM

- funkcja decyzyjna $g(\mathbf{x}) = \mathbf{w}\phi(\mathbf{x}) + b$
- problem optymalizacji

minimalizuj wyrażenie: $\frac{1}{2} \|\mathbf{w}\|^2$

Przy warunkach ograniczających:

$$y_i \cdot ((\mathbf{w} \cdot \Phi(x_i)) + b) \geq 1, \forall_i \in U$$

- Funkcja z mnożnikiem Lagrange'a:

$$\max_{\alpha} W(\alpha) = \max_{\alpha} \left(\alpha \cdot 1 - \frac{1}{2} \alpha \cdot D\alpha \right)$$

gdzie $D_{ij} = y_i y_j \cdot (\Phi(x_i) \cdot \Phi(x_j))$

- Funkcja klasyfikująca: $\text{sgn} \left(\sum_{i=1}^l \alpha_i y_i \cdot (\Phi(x_i) \cdot \Phi(x_j)) + b \right)$

Kernel trick

- Jak obliczyć $\Phi(x_i) \cdot \Phi(x_j)$
 - $K(x,z) = (x \cdot z)^2$, $x = (x_1, x_2)$, $z = (z_1, z_2)$
 - $K(x,z) = (x \cdot z)^2 = (x_1 z_1 + x_2 z_2)^2 = (x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2) \cdot (z_1^2, \sqrt{2}z_1 z_2, z_2^2) = \phi(x) \cdot \phi(z)$
 - Więc D_{ij} można zapisać: $D_{ij} = y_i y_j \cdot K(x_i, x_j)$
- nie trzeba znać funkcji $\phi(x)$, do operacji w wyższej przestrzeni wystarczy znajomość jądra (kernel)

Funkcje jądra (kernel functions)

- wielomianowe (polynomial):

$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z} + d)^p, p > 0$$

- gaussowskie (radial basis function):

$$K(\mathbf{x}, \mathbf{z}) = e^{-\frac{\|\mathbf{x}-\mathbf{z}\|}{2\sigma^2}}$$

- sigmoidalne:

$$K(\mathbf{x}, \mathbf{z}) = \operatorname{tgh}(\eta \mathbf{x} \cdot \mathbf{z} - \delta)$$

Kilka uwag praktycznych

- Normalizuj danej wejściowe
- Rozpoczynaj od zastosowania jądra RBF
- Znajdź optymalne wartości C i σ . Jak? np. grid search
- W klasyfikatorze wykorzystaj parametry znalezione w grid-search

Grid-search – znalezienie maksimum

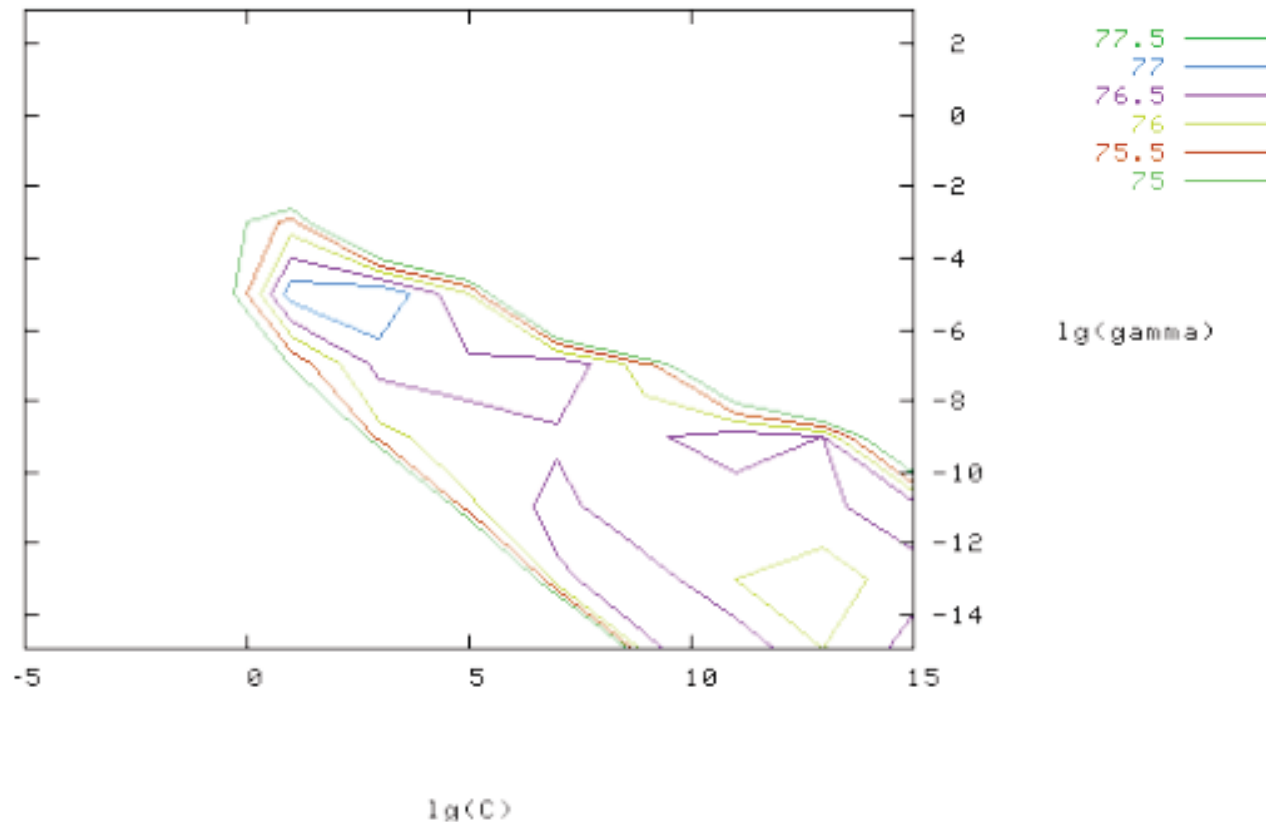


Figure 2: Loose grid search on $C = 2^{-5}, 2^{-3}, \dots, 2^{15}$ and $\gamma = 2^{-15}, 2^{-13}, \dots, 2^3$.

Grid-search – szukanie w okolicy maksimum

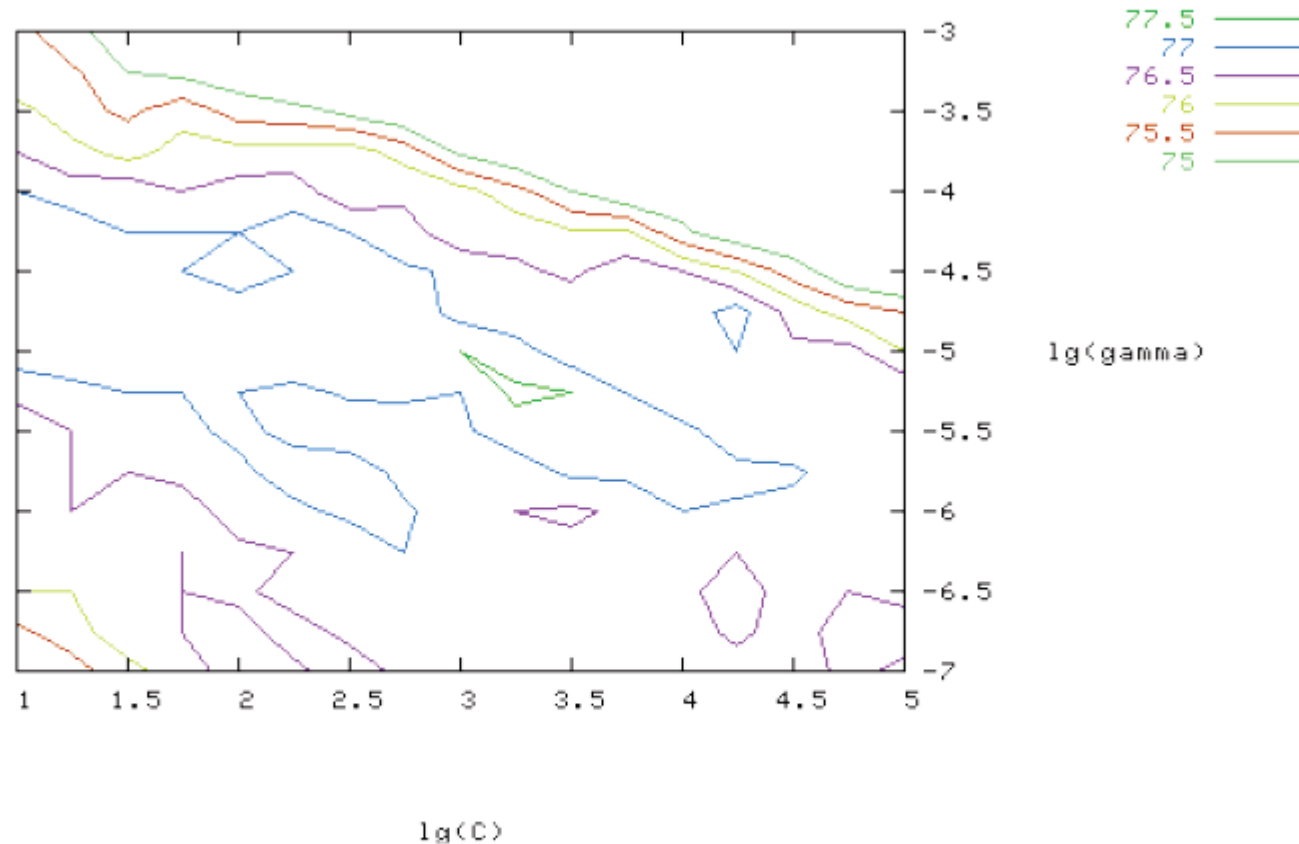


Figure 3: Fine grid-search on $C = 2^1, 2^{1.25}, \dots, 2^5$ and $\gamma = 2^{-7}, 2^{-6.75}, \dots, 2^{-3}$.

Klasyfikacja wieloklasowa- One-versus-all

- One-versus-all – wiele klasyfikatorów dwuklasowych. Każdy klasyfikator dzieli dane wejściowe na klasę zainteresowania i na „resztę”.
- Decyzja o przynależności do klasy podejmowana może być w różny sposób np. głosowanie większościowe, pewność decyzji ...
- Konieczność wytrenowania tylu klasyfikatorów ile klas.

Klasyfikacja wieloklasowa- One-against-one

- One-against-one – wiele klasyfikatorów dwuklasowych. Klasyfikatory dla każdej pary klas
- Decyzja podejmowana podobnie jak dla One-versus-all
- Konieczność wytrenowania $k(k-1)/2$ klasyfikatorów, gdzie k to liczba klas

Implementacje

- C++
 - libSVM
 - SVM light
- Java
 - Weka
- Matlab:
 - libSVM
 - Spider

Bibliografia

- A. Bartkowiak: Wykłady nt. Sieci Neuronowych: w1 1 Kernele, siecie SVM i sieci GDA.
<http://www.ii.uni.wroc.pl/~aba/>
- J. STEFANOWSKI, SVM – Support Vector Machines Metoda wektorów nośnych,
<http://www.cs.put.poznan.pl/jstefanowski/ml/SVM.pdf>
- J. Weston: Support Vector Machine (and Statistical Learning Theory) Tutorial,
http://www.cs.columbia.edu/~kathy/cs4701/documents/jason_svm_tutorial.pdf
- <http://www.kernel-machines.org/>

Bibliografia

- C. Hsu, C. Chang, C. Lin, “A practical guide to Support Vector Classification”, April 2010, Taiwan. (<http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>)
- M. Hoffman, Support Vector Machines — Kernels and the Kernel Trick, http://www.cogsys.wiai.uni-bamberg.de/teaching/ss06/hs_svm/slides/SVM_Seminarbericht_Hofmann.pdf



Dziękuję za uwagę !!!