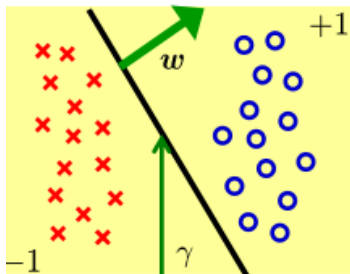


Support Vector Machines

March 28, 2019

- 1 Maximum margin classification
- 2 Nonlinearization by Kernel trick
- 3 Bibliografie

Hard margin classification



$$f(x) = w^T x + b \quad (1)$$

- w , b parametri - normala și intercept pentru hiperplanul de separație
- parametri învățați a.î. **marginea** pt. samples este pozitivă:

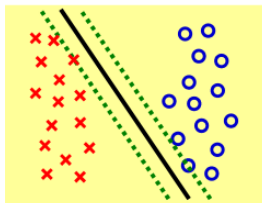
$$f(x_i)y_i = (w^T x_i + b)y_i > 0, \quad i = 1 \dots n \quad (2)$$

$$y_i \in \{-1, +1\}, \quad i = 1 \dots n$$

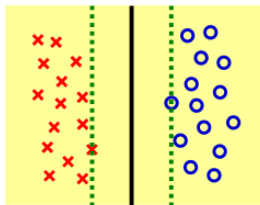
- dificil de rezolvat, constrângerea se rescrie - scalele pentru w și b sunt arbitrare:

$$(w^T x_i + b)y_i \geq 1, \quad i = 1 \dots n \quad (3)$$

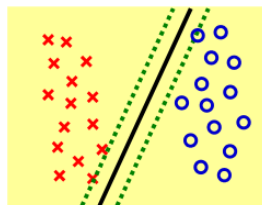
- dacă există $\{w, b\}$ setul $\{(x_i, y_i)\}_{i=1}^n$ este zis **liniar separabil**
- există un număr infinit de hiperplane de separație
- îl selectăm pe cel care are marginea maximă



(a) Small margin



(b) Large margin



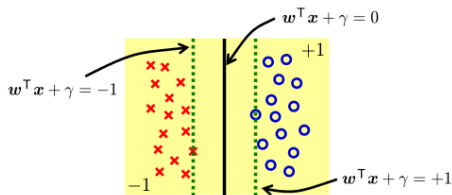
(c) Small margin

Marginea față de hiperplan

- distanța d de la un punct x la hiperplanul $g(x) = w^T x + b = 0$:

$$x = x_p + d \frac{w}{\|w\|} \quad (4)$$

$$\begin{aligned} g(x) &= g\left(x_p + d \frac{w}{\|w\|}\right) = \\ w^T \left(x_p + d \frac{w}{\|w\|}\right) + b &= (w^T x_p + b) + \frac{w^T w}{\|w\|} d = \frac{\|w\|^2}{\|w\|} d = d \|w\| \\ \Rightarrow d &= \frac{g(x)}{\|w\|} \quad (5) \end{aligned}$$



- pentru toate punctele x_i , marginile m_i vor fi:

$$m_i = (w^T x_i + b)y_i / \|w\| \quad (6)$$

$$\min_{i=1 \dots n} \frac{(w^T x_i + b)y_i}{\|w\|} = \frac{1}{\|w\|} \quad (7)$$

- geometric, marginea unui set (x sau o) este jumătate din distanța dintre cele două **hiperplane**
- denumit hard margin SVM:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \text{ s.t. } (w^T x_i + b)y_i \geq 1, \quad i = 1 \dots n \quad (8)$$

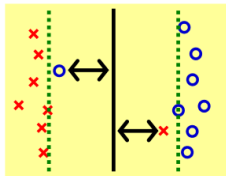
Soft margin SVM

- hard SVM cere linear separability - în practică rar satisfăcut
- soft margin SVM relaxează prin adăugarea erorii ξ_i pentru margini

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

$$s.t. (w^T x_i + b)y_i \geq 1 - \xi_i, \xi_i \geq 0, i = 1 \dots n \quad (9)$$

- C controlează erorile peste margine - C mare determină ξ_i mici (tinde la hard margin SVM)
- ξ_i sunt denumite **slack variables**



Optimizarea duală a SV classification

- ecuația anterioară ia forma unei probleme quadratic programming cu constrângeri liniare
- numărul de parametri $d + n + 1$ nu scalează bine pentru dataseturi mari
- Lagrangianul¹ problemei de optimizare:

$$L(w, b, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i ((w^T x_i + b) y_i - 1 + \xi_i) - \sum_i \beta_i \xi_i \quad (10)$$

- duality principle: soluția maximizând problema duală este \leq soluția minimizând problema primal

$$\max_{\alpha, \beta} \min_{w, b, \xi} L(w, b, \xi, \alpha, \beta) \text{ s.t. } \alpha \geq 0, \beta \geq 0 \quad (11)$$

¹soluția la optimizarea lui $f(x, y)$ cu constrângerea $g(x, y)$ se bazează pe aceeași direcție a gradientilor, $\nabla f(x, y) = \lambda \nabla g(x, y)$

- condiția de optimalitate pentru $\min_{w,b,\xi} L(w, b, \xi, \alpha, \beta)$ determină:

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \quad (12)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \quad (13)$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \alpha_i + \beta_i = C \quad (14)$$

- w și β pot fi exprimate în funcție de α , iar ξ_i dispar din Lagrangian

- forma Lagrange duală se simplifică:

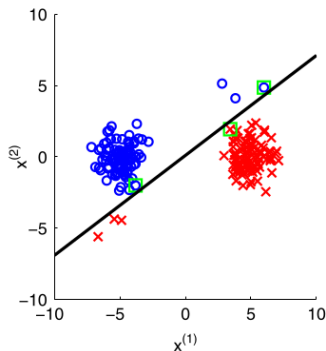
$$\hat{\alpha} = \arg \max_{\alpha} \left[\sum_i \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \right]$$

$$s.t. \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C, i = 1 \dots n \quad (15)$$

- quadratic problem cu doar n variabile
- sau rezolvare populară: John Platt, Sequential Minimal Optimization algorithm, construiește soluția pas cu pas (rapidă)
- soluția $\hat{\alpha}$ dă soluția SV classification:

$$\hat{w} = \sum_i \hat{\alpha}_i y_i x_i \quad \hat{b} = y_i - \sum_{j:\hat{\alpha}_j > 0} \hat{\alpha}_j y_j x_j^T x_j \quad (16)$$

Sparsitatea soluției duale



- condițiile de optim vin din minimizarea cu constrângeri de tip inegalitate - condițiile Karush-Kuhn-Tucker (KKT)
- variabilele duale α și constrângerile satisfac condițiile **complementary slackness**:

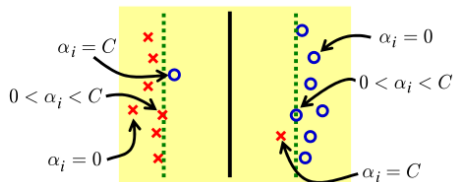
$$\alpha_i(m_i - 1 + \xi_i) = 0 \quad \text{și} \quad \beta_i \xi_i = 0 \quad (17)$$

- unde

$$m_i = (w^T x_i + b)y_i \quad (18)$$

variabilele duale α_i și marginile m_i satisfac relațiile:

- $\alpha_i = 0$ implică $m_i \geq 1$: sample-ul x_i e pe graniță sau înăuntrul ei, corect clasificat
- $0 < \alpha_i < C$ implică $m_i = 1$: x_i chiar pe graniță, corect clasificat (vector suport)
- $\alpha_i = C$ implică $m_i \leq 1$: x_i pe graniță sau în afara marginii; dacă $\xi_i > 1$, $m_i < 0$ și x_i e clasificat incorect
- $m_i > 1$ implică $\alpha_i = 0$: x_i este în margine, corect clasificat
- $m_i < 1$ implică $\alpha_i = C$: x_i este în afara marginii, incorect



$$\alpha_i((w^T x_i + b)y_i - 1 + \xi_i) = 0 \quad (19)$$

- 1 Maximum margin classification
- 2 Nonlinearization by Kernel trick
- 3 Bibliografie

Kernel model

- de regulă în modele precum regresia liniară, sau polinomială, funcțiile bază sunt fixate (polinoame, sinusoide):

$$f_{\theta}(x) = \sum_{j=1}^b \theta_j \phi_j(x) = \theta^T \phi(x) \quad (20)$$

- unde θ_j sunt parametri iar

$$\phi(x) = (1, x, x^2, \dots, x^{b-1}) \quad (21)$$

- modelul kernel utilizează sample-urile x_i pentru proiectarea funcțiilor bază

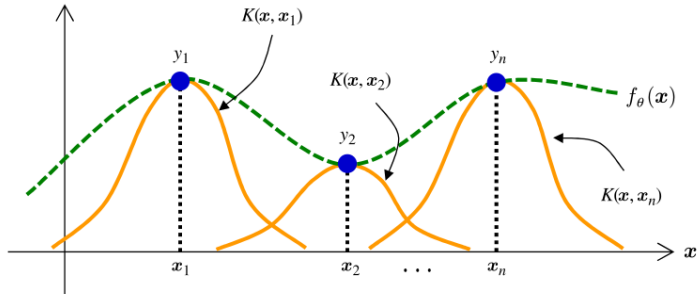
- fie o funcție de două variabile $K(\cdot, \cdot)$
- modelul kernel este definit ca o combinație liniară de funcții $\{K(x, x_j)\}_{j=1}^n$:

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j K(x, x_j) \quad (22)$$

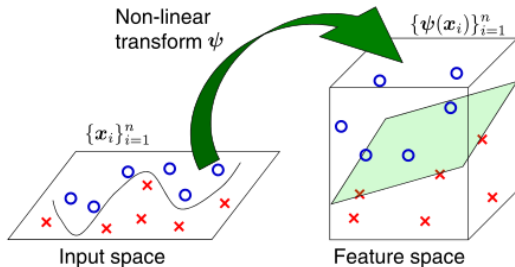
- ca funcție kernel, kernelul Gaussian este cel mai popular:

$$K(x, c) = \exp\left(-\frac{\|x - c\|^2}{2h^2}\right) \quad (23)$$

- h - Gaussian bandwidth, c - Gaussian center



- funcțiile Gaussiene localizate la x_i , este învățată înălțimea lor
- modelul aproximează funcția doar în vecinătatea x_i -urilor
- modelul multiplicativ încerca să aproximeze funcția pe tot spațiul
- nr. de parametri dat de n , independent de dimensionalitatea spațiului
- expresia lui x în sine nu contează, el apare doar în funcția kernel $K(x_i, x_j)$ - kernel trick



- sample-urile x_i sunt translatare într-un feature space folosind o funcție nonliniară Ψ
- SVM liniar este antrenat pe spațiul transformat $\{\Psi(x_i)\}_{i=1}^n$
- spațiul Ψ este de multe ori infinit dimensional - mai mari șanse să fie liniar separabile
- kernel trick: în optimizarea liniară SVM apar doar inner-product-urile:

$$x_i^T x_j = \langle x_i, x_j \rangle \quad (24)$$

- similar:

$$\langle \Psi(x_i), \Psi(x_j) \rangle \quad (25)$$

- trick-ul este specificarea directă a produsului scalar, chiar dacă nu se cunoaște forma analitică a lui $\Psi(x_i)$:

$$\langle \Psi(x_i), \Psi(x_j) \rangle = K(x_i, x_j) \quad (26)$$

- kernel trick se aplică și pentru alte probleme: clustering, dimensionality reduction

- 1 Maximum margin classification
- 2 Nonlinearization by Kernel trick
- 3 Bibliografie

Bibliografie

- M. Sugiyama, “Introduction to Machine Learning”, Elsevier 2016
- https://en.wikipedia.org/wiki/Sequential_minimal_optimization#Algorithm
- https://www.syncfusion.com/ebooks/support_vector_machines_succinctly/solving-the-optimization-problem
- <https://www.svm-tutorial.com/2016/09/duality-lagrange-multipliers/>
- N. Cristianini, J.S. Taylor, “An Introduction to SVM and Other Kernel-based Learning Methods”
- <http://www.csc.kth.se/utbildning/kth/kurser/DD3364/Lectures/KKT.pdf>