

Rețele Mașini Vector Suport. Machine learning pentru clasificarea documentelor

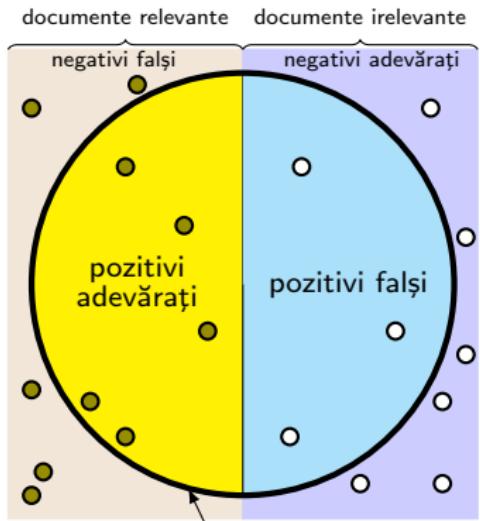
Mircea Marin

Departamentul of Informatică
Universitatea de Vest din Timișoara
mircea.marin@e-uvt.ro

noiembrie 2018

Machine learning pentru clasificarea documentelor

Calculul gradului acurateței



mulțimea documentelor selectate

$$\text{acuratețe} = \frac{\text{pozitivi adevărați} + \text{falși adevărați}}{\text{nr.total documente}} = \frac{11}{19} = 0.579$$

Acuratețea este o măsură a efectivității metodelor de clasificare bazate pe machine learning.

● : document relevant

○ : document irrelevant

$$\text{precizie} = \frac{\text{pozitivi adevărați}}{\text{nr.răspunsuri}} = \frac{5}{8} = 0.625$$

$$\text{reamintire} = \frac{\text{pozitivi adevărați}}{\text{nr.doc.relevante}} = \frac{5}{10} = 0.5$$

Precizia măsoară **exactitatea**
sau **calitatea** răspunsurilor.

Reamintirea măsoară **completitudinea**
sau **cantitatea** răspunsurilor.

Machine learning pentru clasificarea documentelor

Clasificatori cu grad mare de acuratețe

Eforturile de creștere a acurateții din ultimele decenii au produs metode performante de clasificare:

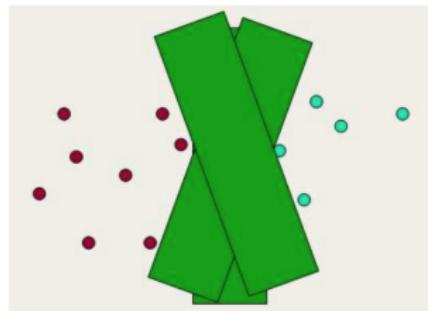
- **mașini vector suport (SVM)**: metodă de clasificare efectivă pt. documente reprezentate în spațiu vectorial.
- arbori amplificați de decizie (engl. boosted decision trees)
- regresie logistică regularizată
- rețele neuronale
- păduri aleatoare (engl. random forests)

Mașini vector suport (SVM)

Cazul liniar separabil



- SVM separă clasele calculând o suprafață de decizie aflată la distanță maximă de punctele clasificate.
- Exemplu ilustrat de suprafete (benzi) de decizie posibile, colorate verde:



Mașini vector suport (SVM)

Noțiuni algebrice (1)

Presupunem dată o mulțime de antrenare

$$\mathbb{D} = \{(\vec{d}_i, y_i) \mid 1 \leq i \leq N\}$$
 cu

- documente \vec{d}_i reprezentare în un spațiu vectorial \mathbb{R}^M .
- 2 clase: -1 și 1. Deci $y_i \in \{-1, 1\}$ pentru $1 \leq i \leq N$.

Un **hiperplan de decizie** $\langle \vec{w}, b \rangle$ este dat de ecuația

$$\vec{w} \cdot \vec{x} = -b$$

unde $\vec{w} \in \mathbb{R}^M$ și $b \in \mathbb{R}$. Hiperplanul de decizie definește
clasificatorul liniar

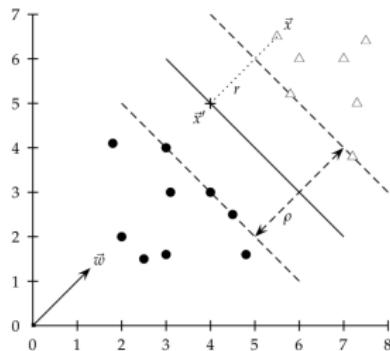
- $f : \mathbb{R}^M \rightarrow \mathbb{R}$, $f(\vec{d}) := \text{sign}(\vec{w} \cdot \vec{d} + b)$

Mașini vector suport

Margine funcțională

Distanța de la un punct \vec{d} într-o clasă $y \in \{-1, 1\}$ la un hiperplan $\langle \vec{w}, b \rangle$ dat de ecuația

$$\vec{w} \cdot \vec{x} = -b \quad \text{este} \quad r = y \cdot \frac{\vec{w} \cdot \vec{d} + b}{|\vec{w}|}$$



$y \cdot (\vec{w} \cdot \vec{d} + b)$ se numește marginea funcțională a lui $\langle \vec{d}, y \rangle$ în raport cu hiperplanul $\langle \vec{w}, b \rangle$.

Mașini vector suport

Problema de învățare

Problema de învățare a valorilor \vec{w}, b, ρ a mașinilor vector suport (SVM) pentru 2 clase liniar separabile este:

- (1) Se impune condiția ca marginile funcționale ale tuturor documentelor \vec{d}_i din setul de antrenare să fie ≥ 1 , adică

$$y_i \cdot (\vec{w} \cdot \vec{d}_i + b) \geq 1$$

și 1 pentru cel puțin un document de antrenare $\langle \vec{d}_i, y_i \rangle \in \mathbb{D}$.

- ▶ vectorii \vec{d}_i ai documentelor de antrenare pt. care $y_i \cdot (\vec{w} \cdot \vec{d}_i + b) = 1$ se numesc **vectori suport**

- (2) Se maximizează valoarea lui $\rho = 2/|\vec{w}|$

$$\Leftrightarrow \text{minimizarea valorii lui } \frac{1}{2}|\vec{w}|^2$$

când au loc constrângerile impuse la (1).

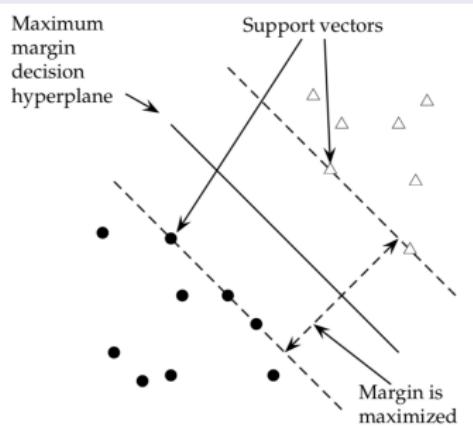
⇒ **problemă de minimizare pătratică** (vezi slide-ul următor)

Mașini vectori suport

Observații. Ilustrare diagramatică

- ρ este lățimea benzii de separare, și se numește **margine geometrică** a SVM.
- $\rho/2$ coincide cu distanța de la vectorii suport la hiperplanul de decizie.

Exemplu ilustrat cu 5 vectori suport



Mașini vector suport

O problemă de optimizare pătratică

Să se determine $\vec{w} \in \mathbb{R}^M$ și $b \in \mathbb{R}$ a.î.

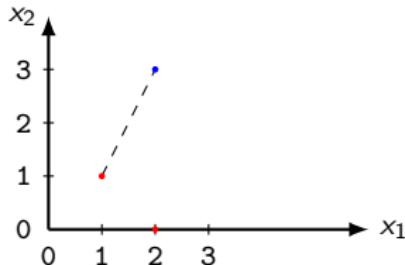
- $|\vec{w}|^2/2 = \frac{1}{2} \sum_{i=1}^M w_i^2$ are valoare minimă, și
- $y_i \cdot (\vec{w} \cdot \vec{d}_i + b) \geq 1$ pentru toți $\langle \vec{d}_i, y_i \rangle \in \mathbb{D}$

Metode de rezolvare a problemei de învățare pentru SVM:

- ① folosind biblioteci/algoritmi standard de rezolvare a problemelor de optimizare pătratică
- ② folosind algoritmi optimizați (mai rapizi și mai scalabili) pt. problema de optimizare pătratică specifică mașinilor vector suport.

Mașini vector suport

Rezolvare geometrică a unui exemplu concret



$$\mathbb{D} = \{(1, 1), -1\}, \{(2, 0), -1\}, \{(2, 3), 1\}\}$$

- ① În acest exemplu, **hiperplanul optim de decizie** $\langle \vec{w}, b \rangle$ este paralel cu segmentul cel mai scurt dintre vectori din clase diferite, care este $(1, 1) - (2, 3)$:

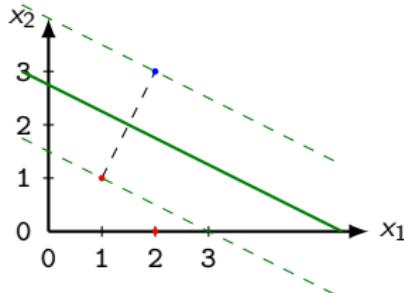
- $(1, 1)$ și $(2, 3)$ devin vectori suport.
- $\vec{w} = \langle 1, 2 \rangle$ fiindcă \vec{w} este paralel cu vectorul de la $(1, 1)$ la $(2, 3)$,
- $b = 5.5$ fiindcă hiperplanul trece prin mijlocul segmentului $(1, 1) - (2, 3)$, adică prin $(1.5, 2)$.

\Rightarrow hiperplanul de decizie $a \cdot (x_1 + 2x_2 - 5.5) = 0$ cu $a > 0$.

- ② **Marginile funcționale** ale vectorilor suport $(1, 1)$ și $(2, 3)$ trebuie să fie 1, deci $-a \cdot (1 + 2 - 5.5) = 1$ și $a \cdot (2 + 6 - 5.5) = 1 \Rightarrow a = \frac{2}{5}$, deci vom considera $\vec{w} = (2/5, 4/5)$ și $b = 11/5$.

Mașini vector suport

Rezolvare geometrică a unui exemplu concret

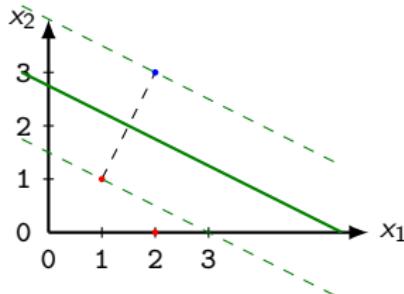


$$\mathbb{D} = \{(1, 1), -1\}, \{(2, 0), -1\}, \{(2, 3), 1\}\}$$

- ① În acest exemplu, **hiperplanul optim de decizie** $\langle \vec{w}, b \rangle$ este paralel cu segmentul cel mai scurt dintre vectori din clase diferite, care este $(1, 1) - (2, 3)$:
 - $(1, 1)$ și $(2, 3)$ devin vectori suport.
 - $\vec{w} = \langle 1, 2 \rangle$ fiindcă \vec{w} este paralel cu vectorul de la $(1, 1)$ la $(2, 3)$,
 - $b = 5.5$ fiindcă hiperplanul trece prin mijlocul segmentului $(1, 1) - (2, 3)$, adică prin $(1.5, 2)$. \Rightarrow hiperplanul de decizie $a \cdot (x_1 + 2x_2 - 5.5) = 0$ cu $a > 0$.
- ② **Marginile funcționale** ale vectorilor suport $(1, 1)$ și $(2, 3)$ trebuie să fie 1, deci $-a \cdot (1 + 2 - 5.5) = 1$ și $a \cdot (2 + 6 - 5.5) = 1 \Rightarrow a = \frac{2}{5}$, deci vom considera $\vec{w} = (2/5, 4/5)$ și $b = 11/5$.

Mașini vector suport

Rezolvare geometrică a unui exemplu concret



$$\mathbb{D} = \{(1, 1), -1\}, \{(2, 0), -1\}, \{(2, 3), 1\}\}$$

- ① În acest exemplu, **hiperplanul optim de decizie** $\langle \vec{w}, b \rangle$ este paralel cu segmentul cel mai scurt dintre vectori din clase diferite, care este $(1, 1) - (2, 3)$:
 - $(1, 1)$ și $(2, 3)$ devin vectori suport.
 - $\vec{w} = \langle 1, 2 \rangle$ fiindcă \vec{w} este paralel cu vectorul de la $(1, 1)$ la $(2, 3)$,
 - $b = 5.5$ fiindcă hiperplanul trece prin mijlocul segmentului $(1, 1) - (2, 3)$, adică prin $(1.5, 2)$.

\Rightarrow hiperplanul de decizie $a \cdot (x_1 + 2x_2 - 5.5) = 0$ cu $a > 0$.
- ② **Marginile funcționale** ale vectorilor suport $(1, 1)$ și $(2, 3)$ trebuie să fie 1, deci $-a \cdot (1 + 2 - 5.5) = 1$ și $a \cdot (2 + 6 - 5.5) = 1 \Rightarrow a = \frac{2}{5}$, deci vom considera $\vec{w} = (2/5, 4/5)$ și $b = 11/5$.
- ③ **Marginea geometrică** a mașinii vector suport este $\rho = 2/|\vec{w}| = 2/\sqrt{20/25} = \sqrt{5}$.

Extensiile ale modelului SVM

Clasificarea cu margini soft (engl. soft margin classification)

Permite definirea unei mașini vector suport pt. colecții de date care nu sunt separabile liniar.

- Dacă mulțimea de antrenare \mathbb{D} nu este liniar separabilă, SVM permite clasificarea eronată a unui nr. mic de documente, numite *zgomot* (engl. noise)
 - Fiecare clasificare greșită a unui exemplu \vec{d}_i este penalizată cu un cost proporțional cu val. unei variabile ζ_i

Definiție formală a problemei de învățare SVM cu margini soft:
Să se determine \vec{w} , b și $\zeta_i \geq 0$ astfel încât:

- $\frac{1}{2}|\vec{w}|^2 + C \sum_i \zeta_i$ este minimizat și
- $y_i \cdot (\vec{w} \cdot \vec{d}_i + b) \geq 1 - \zeta_i$ pentru toți $(\vec{d}_i, y_i) \in \mathbb{D}$.

Parametrul C este un termen de regularizare