

Machine Learning

Support Vector Machines (SVM)

Edgar F. Roman-Rangel.
`edgar.roman@itam.mx`

Digital Systems Department.
Instituto Tecnológico Autónomo de México, ITAM.

May 8th, 2021.

Outline

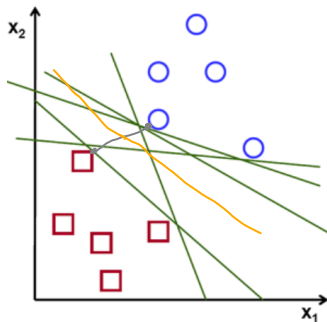
SVM

Kernels

Intro

Support Vector Machine: binary classifier.

Separate two sets (classes, $y = \{-1, +1\}$), by using a hyperplane.

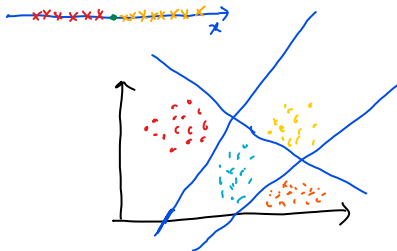
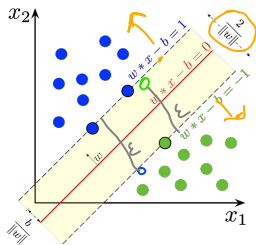


Challenge: find the optimal hyperplane.

Optimal hyperplane

$$W = [w_1, w_2 \dots w_n], \quad |W| = \sqrt{w_1^2 + w_2^2 \dots + w_n^2}$$

Largest margin with respect to closest point from each class.



- Distance → confidence on the decision.
- w : vector perpendicular to the hyperplane.
- Hyperplane of $(N - 1)$ -D, where N is the number of features.
- Number of hyperplanes: $|C| - 1$; $|C|$ is the number of classes.

Margin

Distance from positive to negative boundaries,

$$\begin{aligned} wx^+ + b &= 1 \\ -wx^- + b &= -1 \\ \hline w(x^+ - x^-) &= 2 \end{aligned}$$

Normalize distance to unit length:

$$\frac{w(x^+ - x^-)}{\|w\|} = \frac{2}{\|w\|}.$$

Maximizing $\frac{2}{\|w\|}$ is equivalent to minimizing $\|w\|$. Let's use:

$$\frac{1}{2}\|w\|^2.$$

Constraints

We must consider some constraints for the margin.

Hard margin classifier

cuando $y = 1$
 $(1)\langle x, w \rangle$
 si $\langle x, w \rangle = 1$
 $1 - (1)(1) = 0$
 si $\langle x, w \rangle = -1$
 $1 - (1)(-1) = 2$

$$\frac{1}{2}\|w\|^2 + \sum_{m=1}^M \left(1 - y^{(m)} \langle x^{(m)}, w \rangle\right).$$

Let's become robust to outliers by allowing a few mistakes.

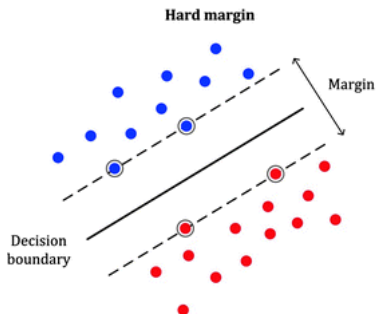
Soft margin classifier

distancia de un punto mal clasificado, al borde de su clase correcta

$$\frac{1}{2}\|w\|^2 + C \left(\sum_{k=1}^K \xi^{(k)} \right),$$

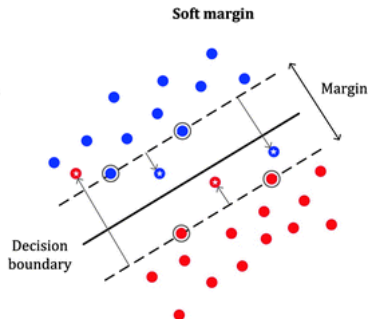
where, ξ is a slack variable that penalizes misclassifications, and C is a hyperparameter. K # puntos mal clasificados encontrados con la solución actual.

Hard vs soft margins



Caso común:

Train \rightarrow acc: 100%
val \rightarrow acc: 80%



Train \rightarrow acc: 96%
val \rightarrow acc: 95%

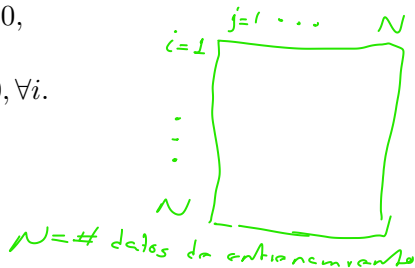
Lagrangian dual

After some math, and optimization rules, we end up with:

$$\max_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underline{x_i^T x_j}$$

$$\text{s.t.} \quad \sum_i \alpha_i y_i = 0,$$

$$C \geq \alpha_i \geq 0, \forall i.$$



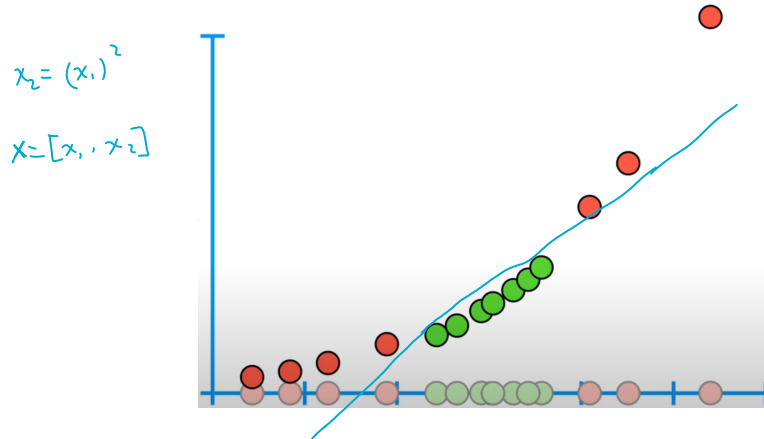
Outline

SVM

Kernels

Non-linearly separable set

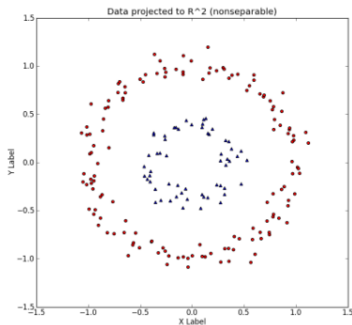
Often, real-world problems are not really linearly separable.



We must find a projection of data, where it becomes linearly separable.

Projection

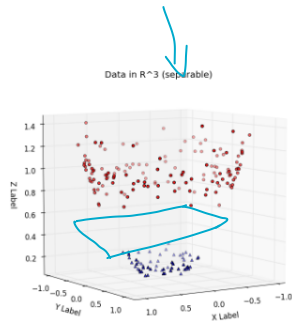
2-D data projected onto a 3-D plane.



$$x = [x_1, x_2]$$

$$x_3 = x_1^2 \quad x_5 = x_1 x_2$$

$$x_4 = x_2^2$$



Kernel functions

$$a = x^{(i)}$$
$$b = x^{(j)}$$

Linear:

$$\kappa(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b}.$$

Polynomial:

$$\kappa(\mathbf{a}, \mathbf{b}) = \left(1 + \sum_{j=1} a_j b_j \right)^d.$$

Radial Basis Function: (Gaussian)

$$\kappa(\mathbf{a}, \mathbf{b}) = \exp \left(\frac{-\gamma \|\mathbf{a} - \mathbf{b}\|^2}{2\sigma^2} \right).$$

Sigmoid:

$$\kappa(\mathbf{a}, \mathbf{b}) = \tanh (c \mathbf{a}^T \mathbf{b} + h).$$

Kernel trick

Compute the relationship between pairs of point, as if they were in the higher dimensional space (do not perform the actual mapping).

Somehow, it is equivalent to map each point to a 1-D space, where its position is the distance with respect to a reference point.

$$\begin{aligned}
 E_i) \text{ Polynomial} \\
 \left(ab + \frac{1}{2}\right)^2 &= \left(ab + \frac{1}{2}\right) \left(ab + \frac{1}{2}\right) \\
 &= ab + a^2b^2 + \frac{1}{2} \quad \leftarrow \text{red arrow} \\
 &= (a, a^2, \frac{1}{2}) (b, b^2, \frac{1}{2})
 \end{aligned}$$

Linear projection y inspires
la matriz de kernel

$$\begin{aligned}
 E_i) \quad a=9, \quad b=14 \\
 \left(ab + \frac{1}{2}\right)^2 &= \left(9 \cdot 14 + \frac{1}{2}\right)^2 \\
 &= \left(26 + \frac{1}{2}\right)^2 \\
 &= 26.5^2 \\
 &= 16.002
 \end{aligned}$$

Linear el truco del
kernel

Q&A

Thank you!

`edgar.roman@itam.mx`