



SVM – the original papers

Adrian Florea

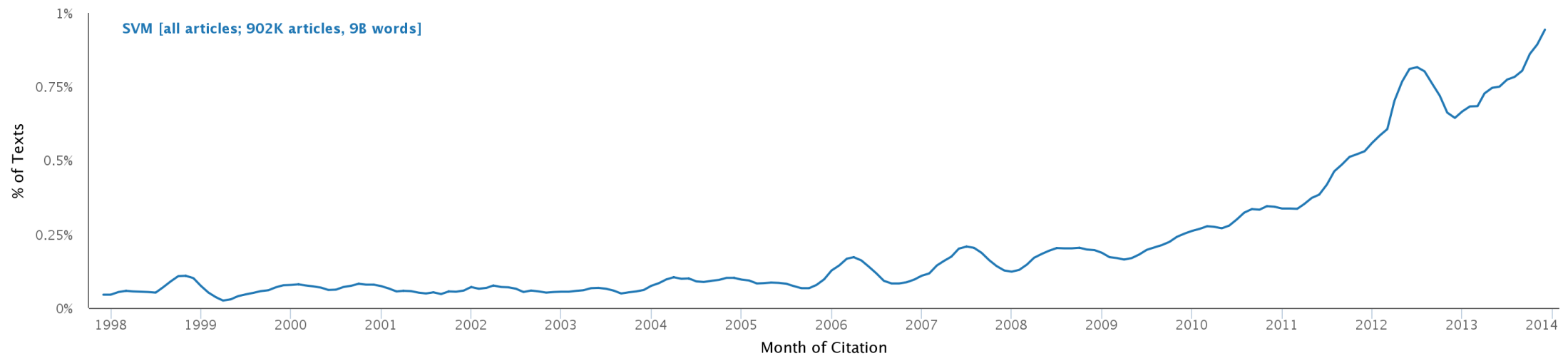
Papers We Love - Bucharest Chapter

November 3rd, 2016

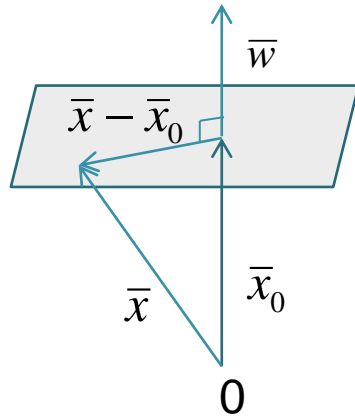
SVM-related papers trend on arXiv

data from Nov 1997 to Nov 2013

<http://bookworm.culturomics.org/arxiv/>



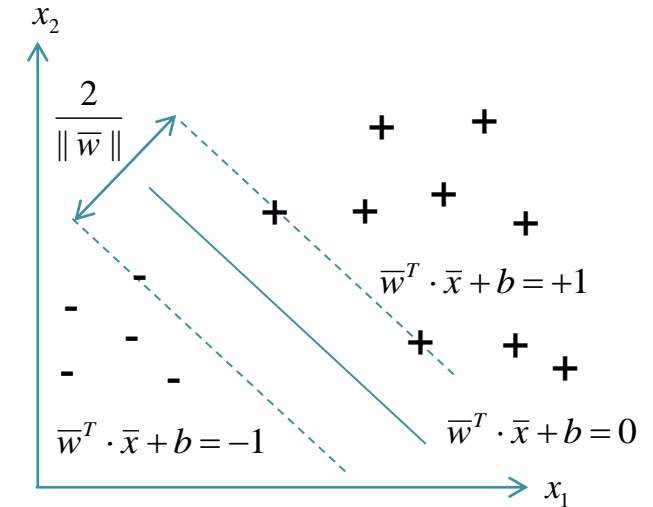
Introduction



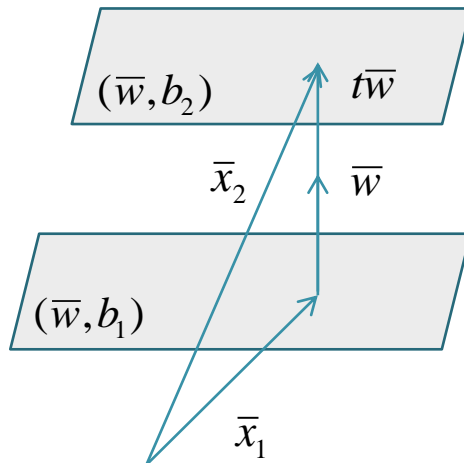
training set: $\left\{ (\bar{x}_i, y_i) \mid i = \overline{1, p}, \bar{x}_i \in R^N, y_i = \begin{cases} +1, \bar{x}_i \in A \\ -1, \bar{x}_i \in B \end{cases} \right\}$

$$d((\bar{w}, b+1), (\bar{w}, b-1)) = \frac{|b+1 - b-1|}{\|\bar{w}\|} = \frac{2}{\|\bar{w}\|}$$

$$\begin{cases} \bar{w} \cdot \bar{x}_i + b \geq +1, y_i = +1 \\ \bar{w} \cdot \bar{x}_i + b \leq -1, y_i = -1 \end{cases} \Leftrightarrow y_i(\bar{w} \cdot \bar{x}_i + b) \geq 1$$



$$\bar{w} \perp \bar{x} - \bar{x}_0 \Rightarrow \bar{w} \cdot (\bar{x} - \bar{x}_0) = 0 \Rightarrow \bar{w} \cdot \bar{x} + b = 0$$



$$d((\bar{w}, b_1), (\bar{w}, b_2)) = \|t\bar{w}\| = |t| \|\bar{w}\|$$

$$\bar{w} \cdot \bar{x}_2 + b_2 = 0 \Rightarrow \bar{w} \cdot (\bar{x}_1 + t\bar{w}) + b_2 = 0 \Rightarrow \bar{w} \cdot \bar{x}_1 + t \|\bar{w}\|^2 + b_2 = 0$$

$$(\bar{w} \cdot \bar{x}_1 + b_1) - b_1 + t \|\bar{w}\|^2 + b_2 = 0 \Rightarrow t = \frac{b_1 - b_2}{\|\bar{w}\|^2}$$

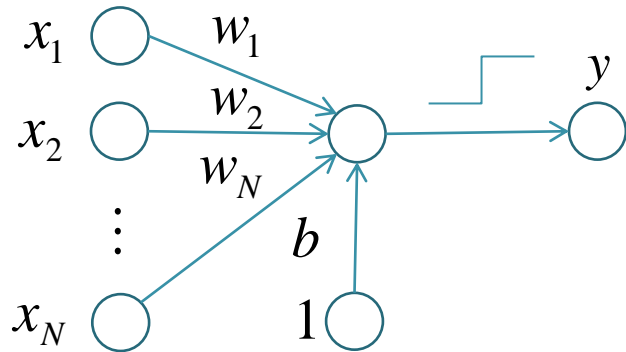
$$d((\bar{w}, b_1), (\bar{w}, b_2)) = |t| \|\bar{w}\| = \frac{|b_1 - b_2|}{\|\bar{w}\|}$$

$$\max_{\bar{w}, b} \frac{1}{\|\bar{w}\|}, \text{ s.t. } y_i(\bar{w} \cdot \bar{x}_i + b) \geq 1$$

$$\max \frac{1}{\|\bar{w}\|} \Rightarrow \min \|\bar{w}\| \Rightarrow \min \frac{\|\bar{w}\|^2}{2}$$

$$y_i(\bar{w} \cdot \bar{x}_i + b) \geq 1 \Rightarrow y_i = \text{sign}(\bar{w} \cdot \bar{x}_i + b)$$

Perceptron



decision function

$$\hat{f}(\bar{x}) = \text{sign}(\bar{w} \cdot \bar{x} + b)$$

$\bar{x} = (x_1, x_2, \dots, x_N)$ features of the customer \bar{x}

approve credit if $\sum_{i=1}^N w_i x_i > \text{threshold}$

deny credit if $\sum_{i=1}^N w_i x_i < \text{threshold}$

input: $T = \{(\bar{x}_i, y_i) \mid i = \overline{1, p}\} \subset \mathbf{R}^N \times \{+1, -1\}$

$\bar{w} \leftarrow \bar{0}, b \leftarrow 0$

repeat

for $i=1$ **to** p

if $\text{sign}(\bar{w} \cdot \bar{x}_i + b) \neq y_i$ **then**

$\bar{w} \leftarrow \bar{w} + y_i \bar{x}_i$

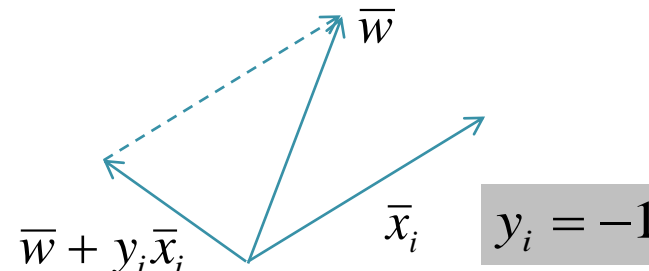
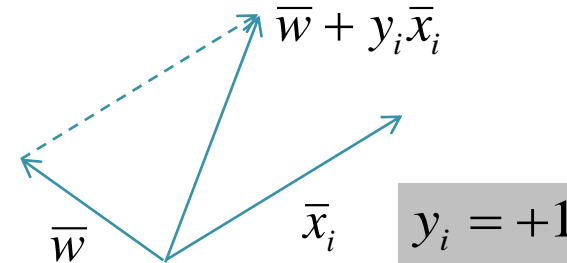
$b \leftarrow b + y_i$

end if

end for

until $\text{sign}(\bar{w} \cdot \bar{x}_j + b) = y_j, \forall j = \overline{1, p}$

return \bar{w}, b



Frank Rosenblatt

CORNELL AERONAUTICAL LABORATORY, INC.
BUFFALO, N. Y.

REPORT NO. 85-160-1

THE PERCEPTRON
A PERCEIVING AND RECOGNIZING AUTOMATON

(PROJECT PARA)

January, 1957

Prepared by: *Frank Rosenblatt*
Frank Rosenblatt,
Project Engineer

Novikoff's Theorem of Convergence

Assume $\exists \bar{w}^*, \|\bar{w}^*\| = 1, \rho > 0$ s.t. $y_i \bar{w}^* \cdot \bar{x}_i \geq \rho, \forall i = \overline{1, p}$

$$R \stackrel{\text{def}}{=} \max_{i=1, p} \|\bar{x}_i\|$$

Then the perceptron algorithm makes at most $\frac{R^2}{\rho^2}$ errors.

A.B.J. Novikoff, "On convergence proofs on perceptrons", Proc. of the *Symposium on the Mathematical Theory of Automata*, **12**, pp 615–622 (1962)

Assume k^{th} error $\bar{w}^{(k)}$ is made on \bar{x}_t and $\bar{w}^{(1)} = \bar{0}$.

$$\bar{w}^{(k+1)} \cdot \bar{w}^* = (\bar{w}^{(k)} + y_t \bar{x}_t) \cdot \bar{w}^* = \bar{w}^{(k)} \cdot \bar{w}^* + y_t \bar{w}^* \cdot \bar{x}_t \geq \bar{w}^{(k)} \cdot \bar{w}^* + \rho$$

By induction on $k : \bar{w}^{(k+1)} \cdot \bar{w}^* \geq k\rho$

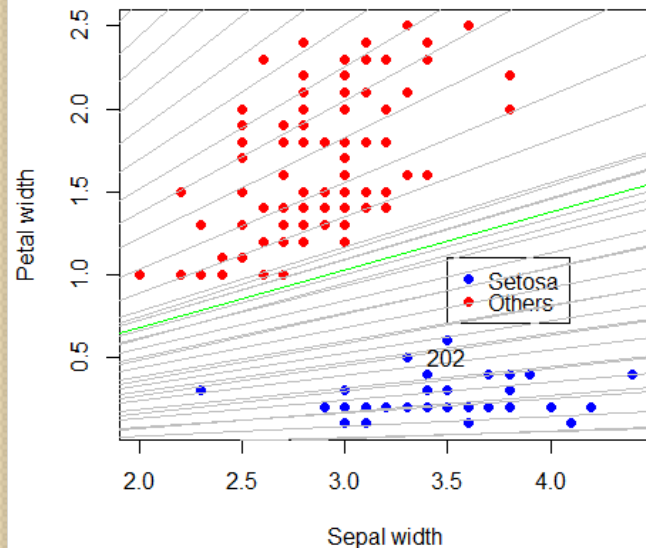
$$\text{Cauchy - Schwartz : } \|\bar{w}^{(k+1)}\| \cdot \|\bar{w}^*\| \geq \bar{w}^{(k+1)} \cdot \bar{w}^*$$

$$\|\bar{w}^*\| = 1 \Rightarrow \|\bar{w}^{(k+1)}\| \geq k\rho$$

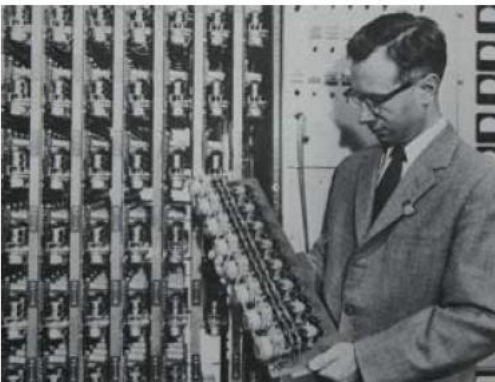
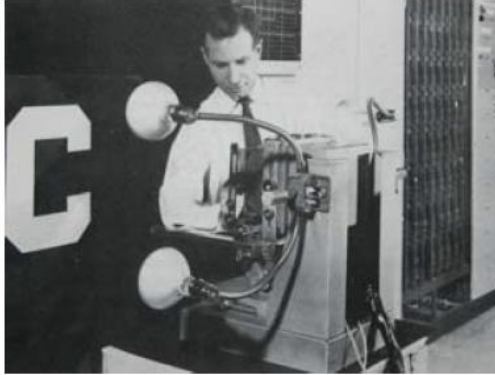
$$\|\bar{w}^{(k+1)}\|^2 = \|\bar{w}^{(k)} + y_t \bar{x}_t\|^2 = \|\bar{w}^{(k)}\|^2 + \underbrace{y_t^2}_{=1} \underbrace{\|\bar{x}_t\|^2}_{\leq R^2} + \underbrace{2 y_t \bar{w}^{(k)} \cdot \bar{x}_t}_{\leq 0 (k^{\text{th}} \text{ error})}$$

$$\text{By induction on } k : \|\bar{w}^{(k+1)}\|^2 \leq kR^2 \Rightarrow k^2 \rho^2 \leq \|\bar{w}^{(k+1)}\|^2 \leq kR^2 \Rightarrow k \leq \frac{R^2}{\rho^2}$$

Demo



Mark I Perceptron hardware



$N = 400$
 $p = 512$

- inputs captured by a 20×20 array of cadmium sulphide photocells

→ **input:** $T = \{(\bar{x}_i, y_i) \mid i = \overline{1, p}\} \subset \mathbf{R}^N \times \{+1, -1\}$

$w \leftarrow 0, b \leftarrow 0$

repeat

→ **for** $i=1$ **to** p

if $\text{sign}(\bar{w} \cdot \bar{x}_i + b) \neq y_i$ **then**

$\bar{w} \leftarrow \bar{w} + y_i \bar{x}_i$

$b \leftarrow b + y_i$

end if

end for

until $\text{sign}(\bar{w} \cdot \bar{x}_j + b) = y_j, \forall j = \overline{1, p}$

return \bar{w}, b

- a patch board allowed different configurations of input features to be tried

- each adaptive weight was implemented using a rotary variable resistor/potentiometer, driven by an electric motor

Linear separability as an LP problem

A necessary and sufficient condition for the **linear separability** of the pattern sets A and B is that:

$$\varphi(A, B) > 0$$

where $\varphi(A, B)$ is the solution of the linear programming problem:

$$\varphi(A, B) = \min_{\bar{u}, \bar{v}, \bar{p}} \{ \bar{1}_n^T \bar{p} \mid \bar{1}_m^T \bar{u} = 1, \bar{1}_k^T \bar{v} = 1, -A^T \bar{u} + B^T \bar{v} + \bar{p} \geq \bar{0}, A^T \bar{u} - B^T \bar{v} + \bar{p} \geq \bar{0}, \bar{u} \geq \bar{0}, \bar{v} \geq \bar{0} \}$$

where u , v , and p are m -, k -, and n -dimensional column vectors.

LINEAR AND NONLINEAR SEPARATION OF PATTERNS BY LINEAR PROGRAMMING

O. L. Mangasarian

Shell Development Company, Emeryville, California

(Received September, 1964)



A necessary and sufficient condition that the sets of patterns A and B be **linearly inseparable** is that the system has a solution:

$$\begin{cases} A^T \bar{u} - B^T \bar{v} = \bar{0} \\ \bar{1}^T \bar{u} = 1 \\ \bar{1}^T \bar{v} = 1 \\ \bar{u} \geq \bar{0} \\ \bar{v} \geq \bar{0} \end{cases}$$

$\text{card}(A)=m, \text{card}(B)=k, R^n$

$(\underline{m + k + n + 2} + n) \times (\underline{m + k} + n)$

Linear separability as an LP problem

$$\begin{bmatrix}
 1 & \dots & 1 & 0 & \dots & 0 & 0 & \dots & 0 \\
 0 & \dots & 0 & 1 & \dots & 1 & 0 & \dots & 0 \\
 -a_{11} & \dots & -a_{m1} & b_{11} & \dots & b_{k1} & 1 & \dots & 0 \\
 \vdots & \dots & \vdots & \vdots & \dots & \vdots & \vdots & \dots & \vdots \\
 -a_{1n} & \dots & -a_{mn} & b_{1n} & \dots & b_{kn} & 0 & \dots & 1 \\
 a_{11} & \dots & a_{m1} & -b_{11} & \dots & -b_{k1} & 1 & \dots & 0 \\
 \vdots & \dots & \vdots & \vdots & \dots & \vdots & \vdots & \dots & \vdots \\
 a_{1n} & \dots & a_{mn} & -b_{1n} & \dots & -b_{kn} & 0 & \dots & 1 \\
 1 & \dots & 0 & 0 & \dots & 0 & 0 & \dots & 0 \\
 \vdots & \dots & \vdots & \vdots & \dots & \vdots & \vdots & \dots & \vdots \\
 0 & \dots & 1 & 0 & \dots & 0 & 0 & \dots & 0 \\
 0 & \dots & 0 & 1 & \dots & 0 & 0 & \dots & 0 \\
 \vdots & \dots & \vdots & \vdots & \dots & \vdots & \vdots & \dots & \vdots \\
 0 & \dots & 0 & 0 & \dots & 1 & 0 & \dots & 0
 \end{bmatrix}
 \begin{Bmatrix}
 u_1 \\
 \vdots \\
 u_m \\
 v_1 \\
 \vdots \\
 v_k \\
 p_1 \\
 \vdots \\
 p_n
 \end{Bmatrix}
 \geq
 \begin{Bmatrix}
 1 \\
 1 \\
 0 \\
 \vdots \\
 0 \\
 0 \\
 0 \\
 \vdots \\
 0 \\
 \vdots \\
 0
 \end{Bmatrix}$$

A and B are **linear separable** iff:

$$\min_{\bar{u}, \bar{v}, \bar{p}} \sum_{i=1}^n p_i > 0$$

A and B are **linear inseparable** iff
this system has a solution

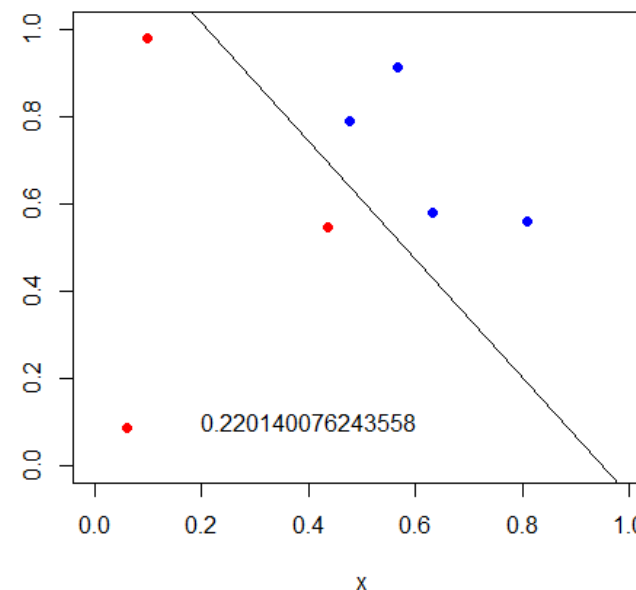
$$\begin{bmatrix}
 a_{11} & \dots & a_{m1} & -b_{11} & \dots & -b_{k1} \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 a_{1n} & \dots & a_{mn} & -b_{1n} & \dots & -b_{kn} \\
 1 & \dots & 1 & 0 & \dots & 0 \\
 0 & \dots & 0 & 1 & \dots & 1 \\
 1 & \dots & 0 & 0 & \dots & 0 \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 0 & \dots & 1 & 0 & \dots & 0 \\
 0 & \dots & 0 & 1 & \dots & 0 \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 0 & \dots & 0 & 0 & \dots & 1
 \end{bmatrix}
 \begin{Bmatrix}
 u_1 \\
 \vdots \\
 u_m \\
 v_1 \\
 \vdots \\
 v_k
 \end{Bmatrix}
 \geq
 \begin{Bmatrix}
 0 \\
 \dots \\
 0 \\
 1 \\
 1 \\
 0 \\
 \dots \\
 0 \\
 0 \\
 \dots \\
 0
 \end{Bmatrix}$$

Linear separability as an LP problem

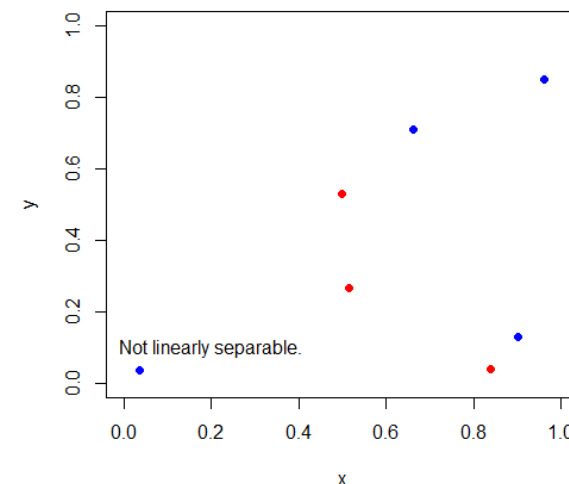
$$\exists \beta \in R, \bar{h} \in R^n, s.t. \forall \bar{a} \in A, \bar{b} \in B:$$

$$\begin{cases} \bar{h}^T \bar{a} > \beta \\ \bar{h}^T \bar{b} < \beta \end{cases} \Rightarrow \begin{cases} \bar{h}^T \bar{a} \geq \beta + \varepsilon \\ \bar{h}^T \bar{b} \leq \beta - \varepsilon \end{cases} \left(\frac{1}{\varepsilon} \right) \Rightarrow \begin{cases} -\bar{h}^T \bar{a} + \beta \leq -1 \\ \bar{h}^T \bar{b} - \beta \leq -1 \end{cases} \quad y$$

$$\begin{bmatrix} -a_{11} & \dots & -a_{1n} & 1 \\ \dots & \dots & \dots & \dots \\ -a_{m1} & \dots & -a_{mn} & 1 \\ b_{11} & \dots & b_{1n} & -1 \\ \dots & \dots & \dots & \dots \\ b_{k1} & \dots & b_{kn} & -1 \end{bmatrix} \begin{Bmatrix} h_1 \\ \dots \\ h_n \\ \beta \end{Bmatrix} \leq \begin{Bmatrix} -1 \\ \dots \\ -1 \\ -1 \end{Bmatrix}$$



Demo

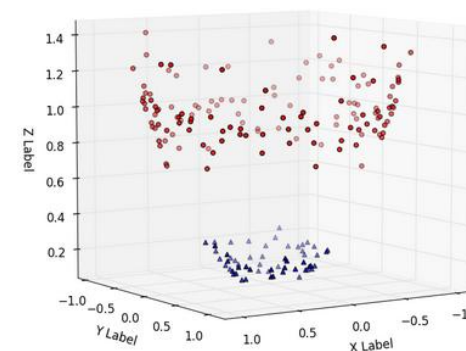
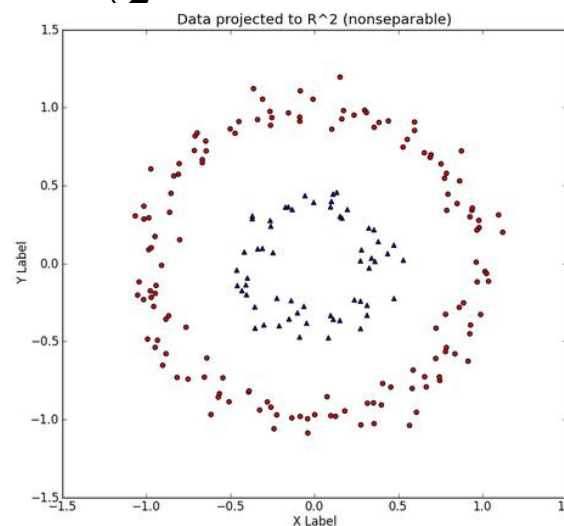
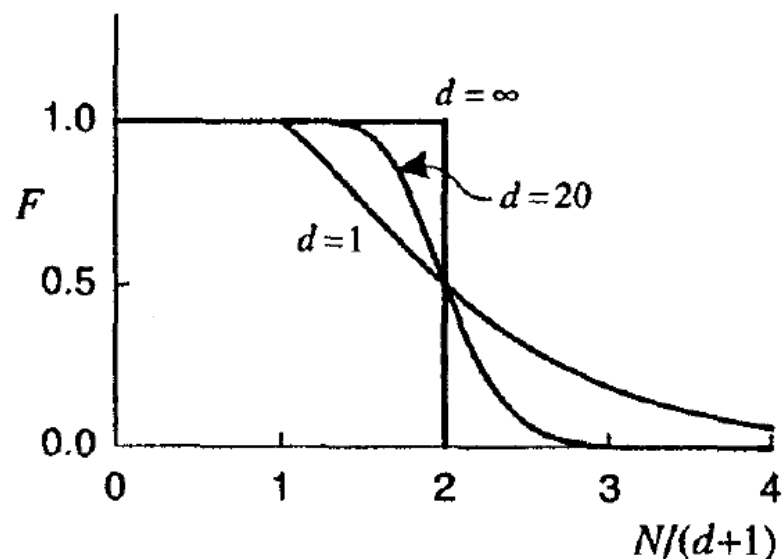


Cover's Theorem

dichotomy : $\{(\bar{x}_i, y_i) \mid i = \overline{1, N}\} \subset \mathbb{R}^d \times \{-1, 1\}$

Assume there is no subset of d or fewer points linearly dependent.

Number of dichotomies linearly separable : $F(N, d) = \begin{cases} 1, & N \leq d + 1 \\ \frac{1}{2^{N-1}} \sum_{i=0}^d C_{N-1}^i, & N \geq d + 1 \end{cases}$

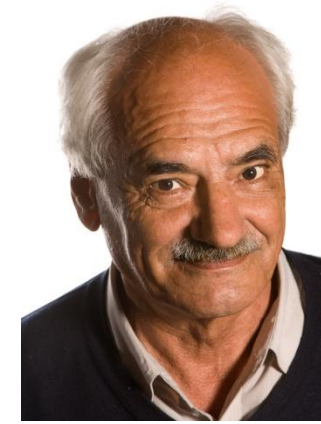


Thomas Cover
Data in \mathbb{R}^3 (separable)

$$F(N, d + 1) = F(N, d) + \frac{C_{N-1}^{d+1}}{2^{N-1}}$$

A complex pattern-classification problem, cast in a **high-dimensional space** nonlinearly, is more likely to be linearly separable than in a **low-dimensional space**

V. Vapnik, A. Chervonenkis,
"A note on one class of perceptrons",
Automation and remote control, **25**, 1 (1964)



Alexey Chervonenkis Vladimir Vapnik

A Training Algorithm for Optimal Margin Classifiers



Bernhard E. Boser*
EECS Department
University of California
Berkeley, CA 94720
boser@eecs.berkeley.edu

Isabelle M. Guyon
AT&T Bell Laboratories
50 Fremont Street, 6th Floor
San Francisco, CA 94105
isabelle@neural.att.com

Vladimir N. Vapnik
AT&T Bell Laboratories
Crawford Corner Road
Holmdel, NJ 07733
vlad@neural.att.com

Bernhard Boser Isabelle Guyon Vladimir Vapnik

Machine Learning, 20, 273–297 (1995)

© 1995 Kluwer Academic Publishers, Boston. Manufactured in The Netherlands.

Support-Vector Networks

CORINNA CORTES
VLADIMIR VAPNIK
AT&T Bell Labs., Holmdel, NJ 07733, USA

corinna@neural.att.com
vlad@neural.att.com



Corinna Cortes Vladimir Vapnik

Lagrangian optimization problem

Primal formulation:

$$\min_{\bar{x}} \Phi(\bar{x}) \text{ s.t. } g_i(\bar{x}) \geq 0 \quad i = \overline{1, p}, \bar{x} \in \mathbb{R}^N, \Phi \text{ convex, } g_i \text{ linear}$$

Dual formulation: $\max_{\bar{\alpha}} \min_{\bar{x}} L(\bar{\alpha}, \bar{x}) = \max_{\bar{\alpha}} \min_{\bar{x}} (\Phi(\bar{x}) - \sum_{i=1}^p \bar{\alpha}_i g_i(\bar{x})) \text{ s.t. } \bar{\alpha}_i \geq \bar{0} \quad i = \overline{1, p}, \bar{x} \in \mathbb{R}^N$

$$\max_{\bar{\alpha}} L(\bar{\alpha}, \bar{x}^*) = \max_{\bar{\alpha}} (\Phi(\bar{x}^*) - \sum_{i=1}^p \bar{\alpha}_i g_i(\bar{x}^*)) \quad \min_{\bar{x}} L(\bar{\alpha}^*, \bar{x}) = \min_{\bar{x}} (\Phi(\bar{x}) - \sum_{i=1}^p \bar{\alpha}_i^* g_i(\bar{x}))$$

$$\max_{\bar{\alpha}} \min_{\bar{x}} L(\bar{\alpha}, \bar{x}) = L(\bar{\alpha}^*, \bar{x}^*) = \Phi(\bar{x}^*) - \sum_{i=1}^p \bar{\alpha}_i^* g_i(\bar{x}^*) \quad \exists! (\bar{\alpha}^*, \bar{x}^*) \in \mathbb{R}^p \times \mathbb{R}^N \text{ saddle point}$$

Karush-Kuhn-Tucker
KKT conditions

$$\frac{\partial L}{\partial \bar{x}}(\bar{\alpha}^*, \bar{x}^*) = \bar{0}$$

$$\bar{\alpha}_i^* g_i(\bar{x}_i^*) = 0, i = \overline{1, p}$$

KKT complementarity condition

$$g_i(\bar{x}_i^*) \geq 0, i = \overline{1, p}$$

$$\bar{\alpha}_i^* \geq 0, i = \overline{1, p}$$

Generalized Portrait Method

Primal formulation for linear SVM

training set: $\{(\bar{x}_i, y_i) \mid i = \overline{1, p}\} \subset R^N \times \{1, -1\}$

$$y_i = \begin{cases} +1, \bar{x}_i \in A \\ -1, \bar{x}_i \in B \end{cases}$$

Classifier for new instances: $f(\bar{x}) = \text{sign}(\bar{w} \cdot \bar{x} + b)$

$$\min \frac{\|\bar{w}\|^2}{2}$$

$$y_i(\bar{w} \cdot \bar{x}_i + b) \geq 1, i = \overline{1, p}$$

convex QP with **N** variables $\bar{w}_i, i = \overline{1, N}$

Dual formulation for linear SVM

Apply the method of Lagrange multipliers:

$$L(\bar{w}, b, \bar{\alpha}) = \frac{1}{2} \sum_{i=1}^N w_i^2 - \sum_{i=1}^p \alpha_i [y_i(\bar{w} \cdot \bar{x}_i + b) - 1]$$

The saddle point is $\min_{\bar{w}} L$ (\bar{w} solves primal)

and $\max_{\bar{\alpha}} L$ ($\bar{\alpha}$ solves dual), $\bar{\alpha} \geq \bar{0}$

$$\frac{\partial L(\bar{w}, b, \bar{\alpha})}{\partial \bar{w}} = 0 \Rightarrow \bar{w} = \sum_{i=1}^p \alpha_i y_i \bar{x}_i$$

$$\frac{\partial L(\bar{w}, b, \bar{\alpha})}{\partial b} = 0 \Rightarrow \sum_{i=1}^p \alpha_i y_i = 0$$

$$L(\bar{\alpha}) = \frac{1}{2} \left(\sum_{i=1}^p \alpha_i y_i \bar{x}_i \right) \left(\sum_{i=1}^p \alpha_i y_i \bar{x}_i \right) - \sum_{i=1}^p \alpha_i y_i \bar{x}_i \left(\sum_{j=1}^p \alpha_j y_j \bar{x}_j \right) - \underbrace{\sum_{i=1}^p \alpha_i y_i}_{=0} b + \sum_{i=1}^p \alpha_i$$

convex QP with **p** variables $\alpha_i, i = \overline{1, p}$

$$L(\bar{\alpha}) = \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \alpha_i \alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j$$

$$\max_{\bar{\alpha}} L(\bar{\alpha})$$

$$\alpha_i \geq 0, i = \overline{1, p}, \sum_{i=1}^p \alpha_i y_i = 0$$

Support Vectors

KKT complementarity condition: $\bar{\alpha}_i^* [y_i^* (\bar{w} \cdot \bar{x}_i^* + b) - 1] = 0, i = \overline{1, p}$

Primal formulation conditions for linear SVM: $y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1, i = \overline{1, p}$

$$\alpha_i \geq 0, i = \overline{1, p}$$

training set linear separability

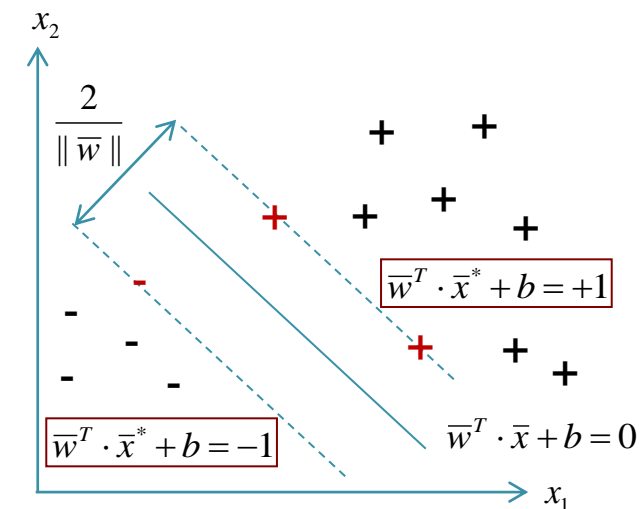
$$\exists S \subset \{1, 2, \dots, p\}, \bar{\alpha}_s^* > 0, \forall s \in S$$

$$y_s^* (\bar{w} \cdot \bar{x}_s^* + b) = 1 \\ \forall s \in S$$

Support vectors = points with non-zero Lagrangian multipliers

$$\text{card}(S) \ll p$$

Primal	Supporting hyperplanes
Dual	Support vectors



Kernel functions

$$\Phi: R^2 \rightarrow R^3, \Phi(\bar{x}) = \Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$\Phi(\bar{x}) \cdot \Phi(\bar{y}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (y_1^2, y_2^2, \sqrt{2}y_1y_2) =$$

$$= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2 =$$

$$= (x_1 y_1 + x_2 y_2)(x_1 y_1 + x_2 y_2) =$$

$$= (\bar{x} \cdot \bar{y})(\bar{x} \cdot \bar{y}) = (\bar{x} \cdot \bar{y})^2$$

$$K(\bar{x}, \bar{y}) = \Phi(\bar{x}) \cdot \Phi(\bar{y}) \quad \text{kernel function}$$

$$K(\bar{x}, \bar{y}) = \bar{x} \cdot \bar{y} \quad \text{Examples}$$

$$K(\bar{x}, \bar{y}) = (p + \bar{x} \cdot \bar{y})^q$$

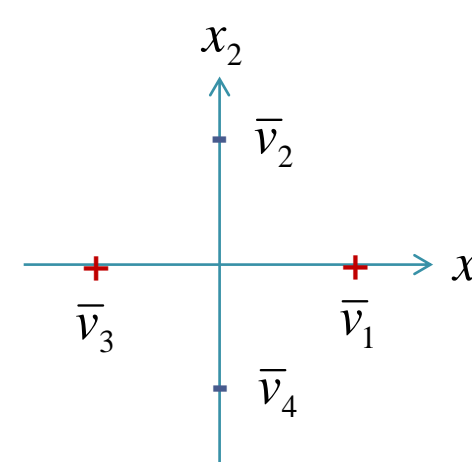
$$K(\bar{x}, \bar{y}) = \exp(-\gamma \|\bar{x} - \bar{y}\|)$$

$$K(\bar{x}, \bar{y}) = \exp(-\gamma \|\bar{x} - \bar{y}\|^2)$$

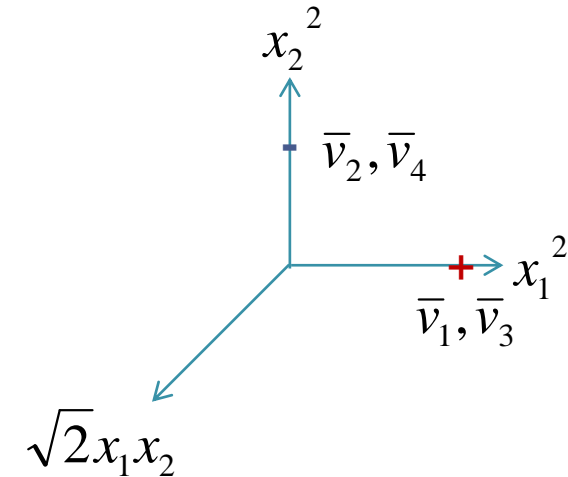
$$K(\bar{x}, \bar{y}) = (p + \bar{x} \cdot \bar{y})^q \exp(-\gamma \|\bar{x} - \bar{y}\|^2)$$

$$K(\bar{x}, \bar{y}) = \tanh(k\bar{x} \cdot \bar{y} - \delta)$$

a kernel expresses a measure of “similarity” between vectors



linear inseparable



linear separable

From Cover's theorem:

A complex pattern-classification problem, cast in a **high-dimensional space** nonlinearly, is more likely to be linearly separable than in a **low-dimensional space**

Kernel trick

input space	image space
$\hat{f}(\bar{x}) = \text{sign}(\bar{w}^* \cdot \bar{x} + b^*) \Leftrightarrow$ $\bar{w}^* = \sum_{i=1}^p \alpha_i^* y_i \bar{x}_i$	$\hat{f}(\bar{x}) = \text{sign}(\bar{w}^* \cdot \Phi(\bar{x}) + b^*) =$ $\bar{w}^* = \sum_{i=1}^p \alpha_i^* y_i \Phi(\bar{x}_i)$
$\hat{f}(\bar{x}) = \text{sign}(\sum_{i=1}^p \alpha_i^* y_i \Phi(\bar{x}_i) \cdot \Phi(\bar{x}) + b) = \text{sign}(\sum_{i=1}^p \alpha_i^* y_i K(\bar{x}_i, \bar{x}) + b)$	
<u>No need to know Φ explicitly!</u>	

$\alpha_i^* = 0, i \notin S \Rightarrow$
{
 instead of obtaining a function with complexity proportional to the image space dimension, we obtained an expression with complexity proportional to the number of support vectors

$$\max_{\bar{\alpha}} \left(\sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i,j=1}^p \alpha_i \alpha_j y_i y_j K(\bar{x}_i, \bar{x}_j) \right)$$

$$\text{s.t. } \sum_{i=1}^p \alpha_i y_i = 0, \alpha_i \geq 0, i = \overline{1, p}$$

Long history of development

- "You look at stuff like this in a book and you think, well, Vladimir Vapnik just figured this out one Saturday afternoon when the weather was too bad to go outside. That's not how it happened." - Patrick Winston
- "The invention of SVMs happened when Bernhard decided to implement Vladimir's algorithm in the three months we had left before we moved to Berkeley. After some initial success of the linear algorithm, Vladimir suggested introducing products of features. I proposed to rather use the **kernel trick** of the 'potential function' algorithm. Vladimir initially resisted the idea because the inventors of the 'potential functions' algorithm (Aizerman, Braverman, and Rozonoer) were from a competing team of his institute back in the 1960's in Russia! But Bernhard tried it anyways, and the SVMs were born!" – Isabel Guyon

Bibliography

- C. Bishop, "Pattern Recognition and Machine Learning", Springer (2006)
- B. Boser, I. Guyon, V. Vapnik, "A Training Algorithm for Optimal Margin Classifiers", Proc. of the 5th annual workshop on Computational Learning Theory (COLT '92), pp. 144-152 (1992)
- T. Cover, "Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition", IEEE Transactions on Electronic Computers, EC-14, No. 3, pp. 326-334 (1965)
- I. Guyon - <https://www.facebook.com/WomenInMachineLearning/posts/1314864275195877>
- L. Hamel, "Knowledge Discovery with Support Vector Machines", Wiley (2009)
- E. Kim, "Everything You Wanted to Know about the Kernel Trick (But Were Too Afraid to Ask)" - http://www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html
- O. Mangasarian, "Linear and Nonlinear Separation of Patterns by Linear Programming", *Operations Research*, **13**, pp. 444-452 (1965)
- H. Murrell, "Data Mining with R" - <http://www.cs.ukzn.ac.za/~hughm/dm/>
- F. Rosenblatt, "The Perceptron. A Perceiving and Recognizing Automaton (Project PARA)", Cornell Aeronautical Laboratory, Report No. 85-460-1 (1957)
- A. Smola, "Introduction to Machine Learning, Class 10-701/15-781", Carnegie Mellon University - <http://alex.smola.org/teaching/10-701-15/>
- A. Statnikov, D. Hardin, I. Guyon, C. Aliferis, "A Gentle Introduction to Support Vector Machines in Biomedicine", The American Medical Informatics Association (AMIA) 2009 Annual Symposium, San Francisco (2009)
- V. Vapnik, "The Nature of Statistical Learning Theory", 2nd Ed., Springer (2000)
- V. Vapnik, "Estimation of Dependences Based on Empirical Data", 2nd ed., Springer (2006)
- R. Vogler, "Testing for Linear Separability with Linear Programming in R" - <http://www.joyofdata.de/blog/testing-linear-separability-linear-programming-r-glpk/>
- P. Winston, "6.034 Artificial Intelligence. Lecture 16: Learning: Support Vector Machines" - <https://www.youtube.com/watch?v=PwhiWxHK8o>