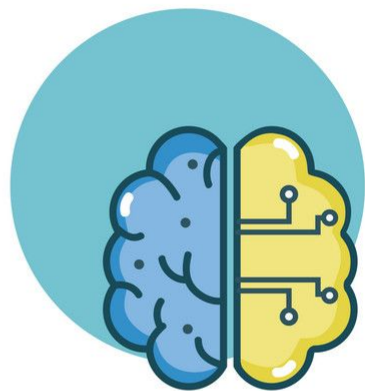


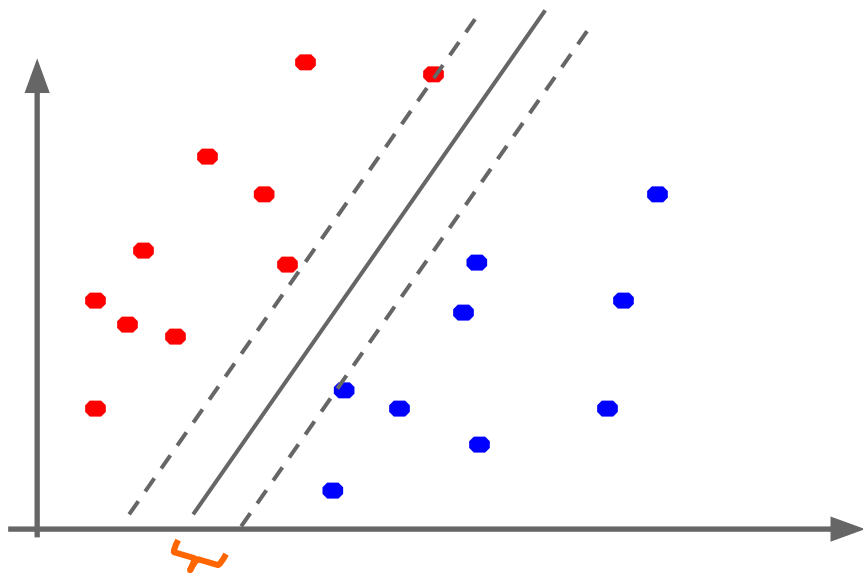
# INTRODUCTION TO MACHINE LEARNING

## SUPPORT VECTOR MACHINES



Elisa Ricci





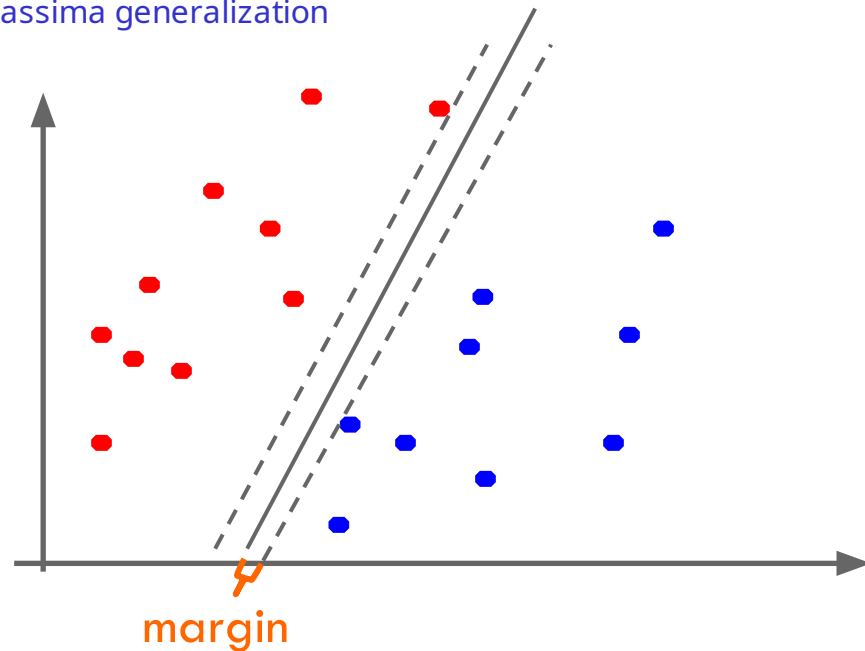
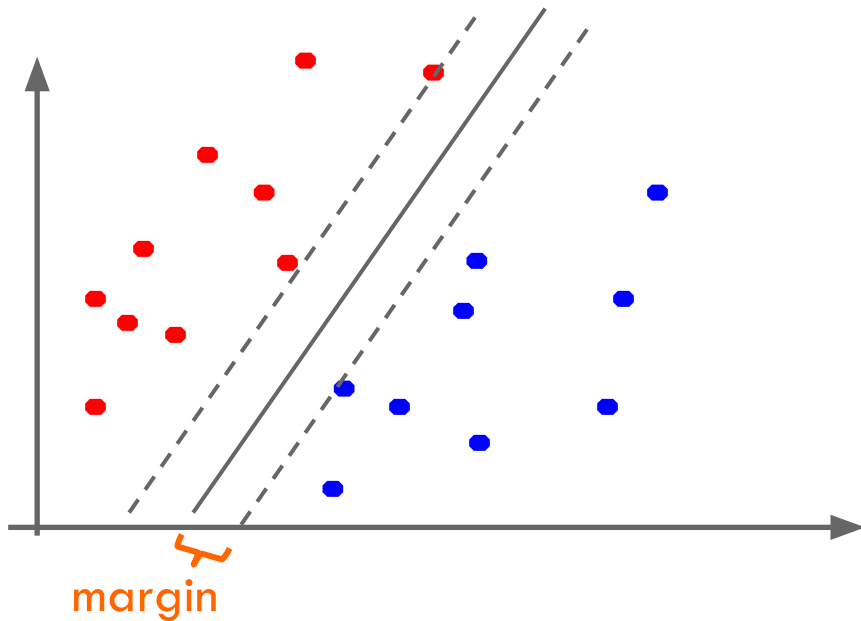
margin

non mi interessa imparare un linear model qualsiasi  
voglio imparare uno specifico linear model  
che massimizza il margine

# SUPPORT VECTOR MACHINES

# LARGE MARGIN CLASSIFIERS

massimo margine -> trovare hyperplane che garantisce la massima generalization



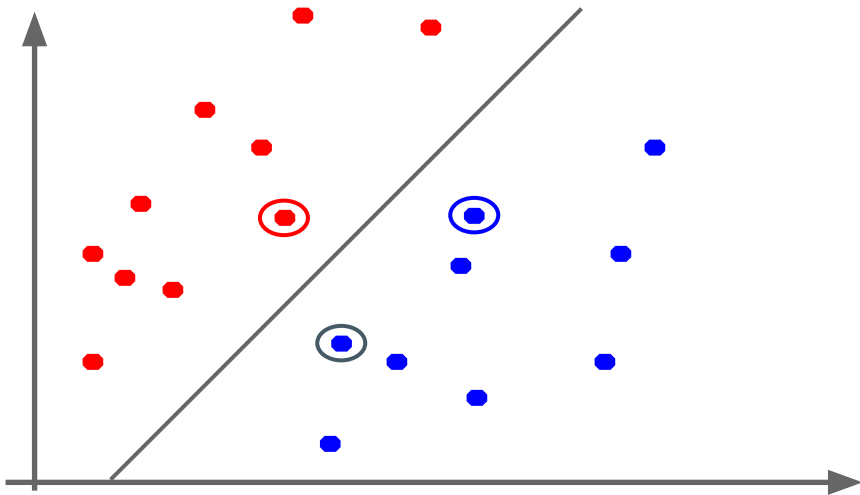
The **margin** of a classifier is the distance to the closest points of either class  
**Large margin** classifiers attempt to maximize this

# SUPPORT VECTORS

For any separating hyperplane, there exist some set of “closest points”

These are called the support vectors

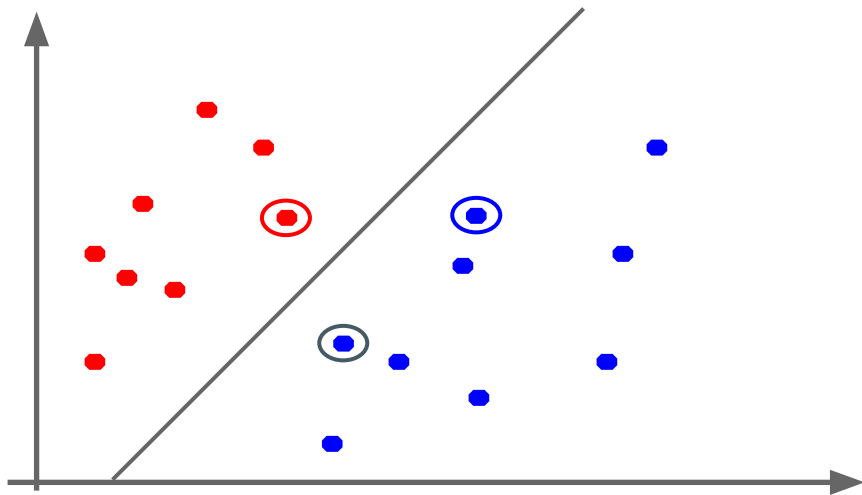
For  $n$  dimensions, there will be at least  $n+1$  support vectors



# LARGE MARGIN CLASSIFIERS

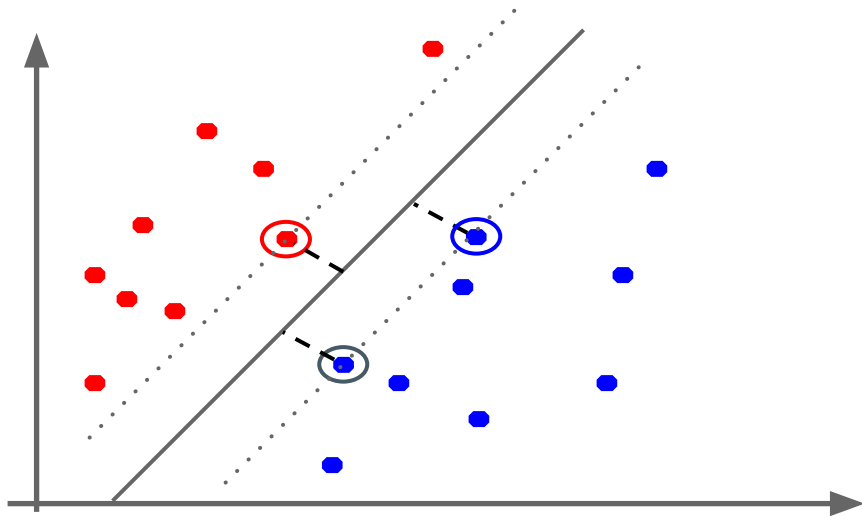
Maximizing the margin is good since it implies that only support vectors matter, other training examples are ignorable.

riduce notevolmente il training time



# MEASURING THE MARGIN

The margin is the distance to the support vectors, i.e. the “closest points”, on either side of the hyperplane



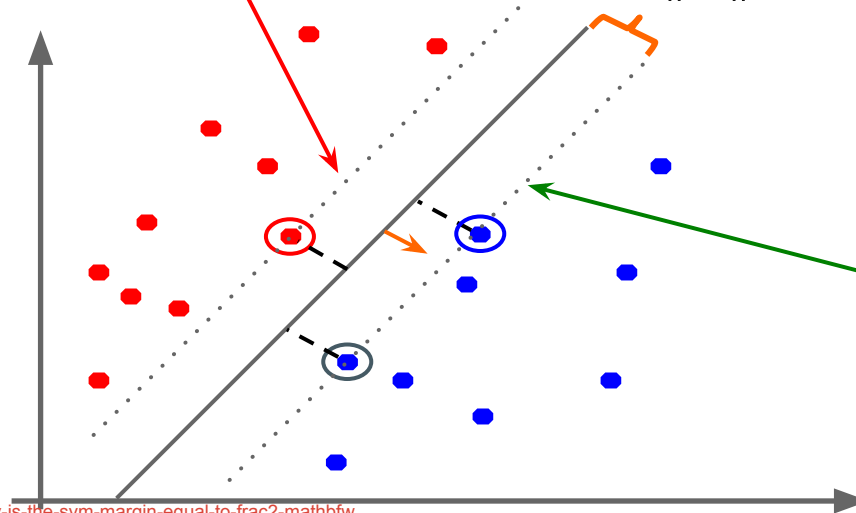
# MEASURING THE MARGIN

$$w \cdot x_i + b = -1$$

$$\frac{w \cdot x_i + b}{\|w\|} = \frac{1}{\|w\|}$$

margine

cerco  $w$   
che massimizza il margine



$$w \cdot x_i + b = 1$$

# MAXIMIZING THE MARGIN

Select the hyperplane with the largest margin where the points are classified correctly and outside the margin!

Setup as a **constrained optimization problem**:

$$\begin{aligned} & \max_{w,b} \quad \text{margin}(w,b) \\ & \text{subject to:} \\ & \quad \underline{y_i(w \cdot x_i + b) \geq 1 \quad \forall i} \end{aligned}$$

*what does  
this mean?*

questo constrain ci assicura che tutti i punti sono fuori dal margine



# MAXIMIZING THE MARGIN

Select the hyperplane with the largest margin where the points are classified correctly and outside the margin!

Setup as a **constrained optimization problem**:

$$\max_{w,b} \frac{1}{\|w\|}$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

# MAXIMIZING THE MARGIN

$$\begin{aligned} & \min_{w,b} \quad \|w\| \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

Maximizing the margin is equivalent to minimize the norm of the weights (subject to separating constraints).

# MAXIMIZING THE MARGIN

The minimization criterion wants  $w$  to be as small as possible

$$\min_{w,b} \|w\|$$

subject to:

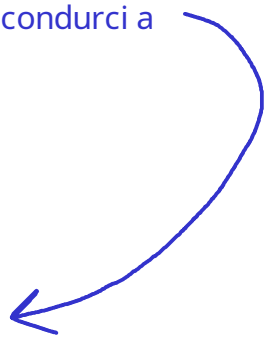
$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

The constraints make sure that the data is separable

# SUPPORT VECTOR MACHINE PROBLEM

$$\begin{array}{ll} \min_{w,b} & \|w\|^2 \\ \text{subject to:} & \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{array}$$

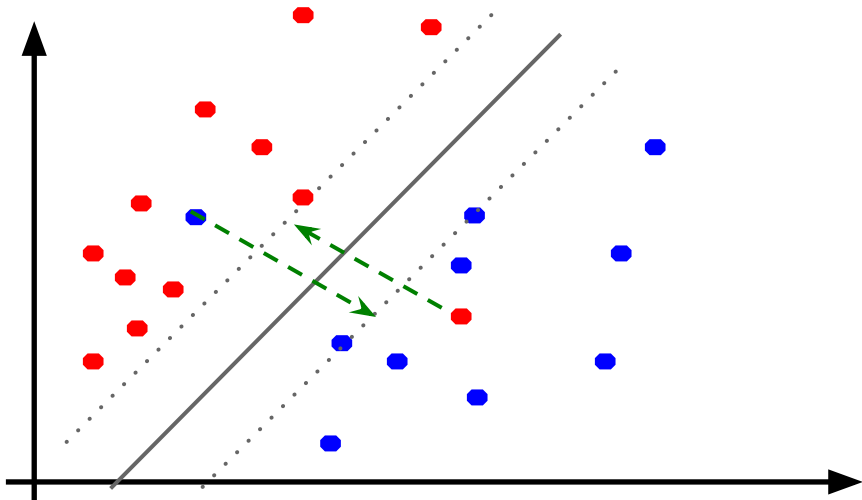
prendiamo la norma al quadrato per ricondurci a



This is a version of a quadratic optimization problem

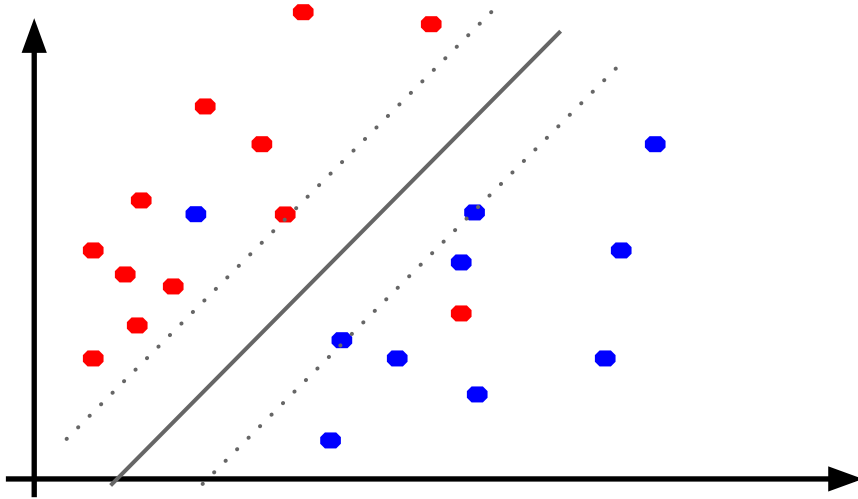
problema noto per cui esistono già solver ottimizzati

Maximize/minimize a quadratic function subject to a set of linear constraints



SOFT MARGIN  
CLASSIFICATION

# SOFT MARGIN CLASSIFICATION

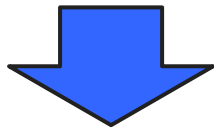


$$\begin{aligned} & \min_{w,b} \|w\|^2 \\ & \text{subject to:} \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

We would like to learn something like this, but our constraints do not allow it...

# SLACK VARIABLES

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

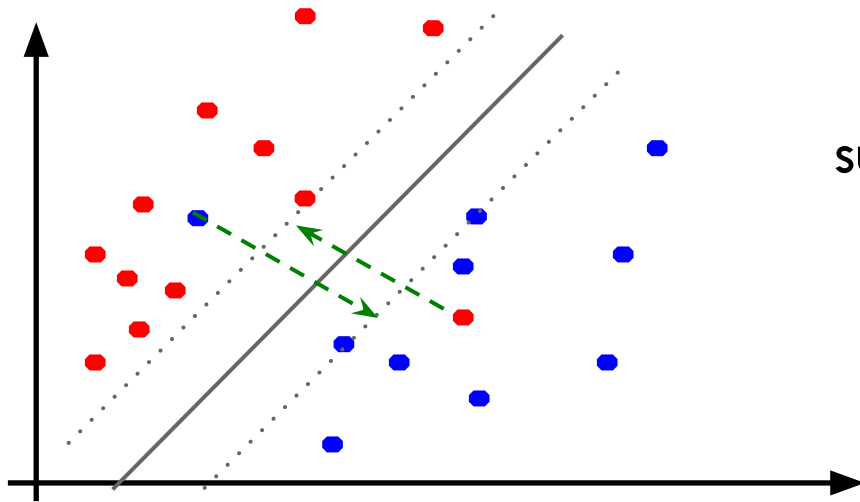


$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \quad \zeta_i \geq 0 \end{aligned}$$

**slack variables**  
(one for each example)

What effect do they have?

# SLACK VARIABLES



slack penalties

$$\min_{w,b} \|w\|^2 + C \sum_i \varsigma_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$

$$\varsigma_i \geq 0$$



# SLACK VARIABLES

margin

trade-off between margin maximization and penalization

$$\min_{w,b} \|w\|^2 + C \sum_i \varsigma_i$$

penalized by how far from “correct”

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$

allowed to make a mistake

$$\varsigma_i \geq 0$$

Still a quadratic optimization problem

# SOFT MARGIN SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

Parameter  $C$  can be viewed as a way to control **overfitting**:  
it “trades off” the relative importance of maximizing the  
margin and fitting the training data.

# SOFT MARGIN SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i$$

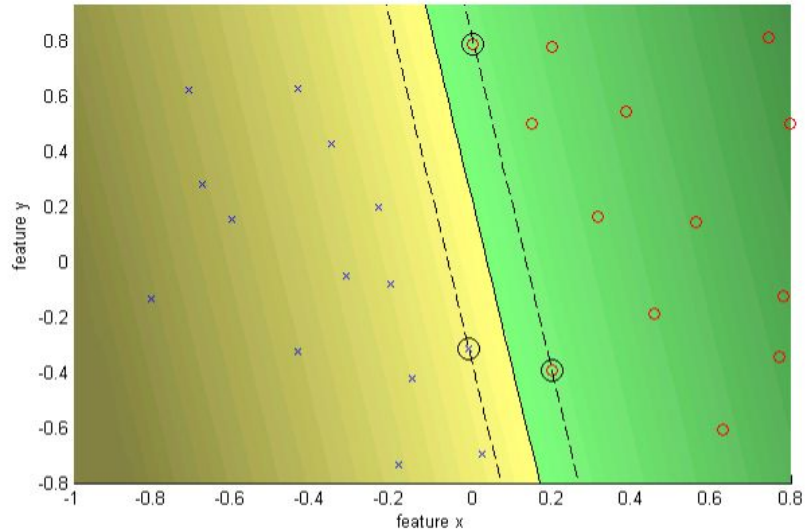
subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

$C$  is a regularization parameter:

- small  $C$  allows constraints to be easily ignored  $\rightarrow$  large margin
- large  $C$  makes constraints hard to ignore  $\rightarrow$  narrow margin
- $C = \infty$  enforces all constraints: hard margin

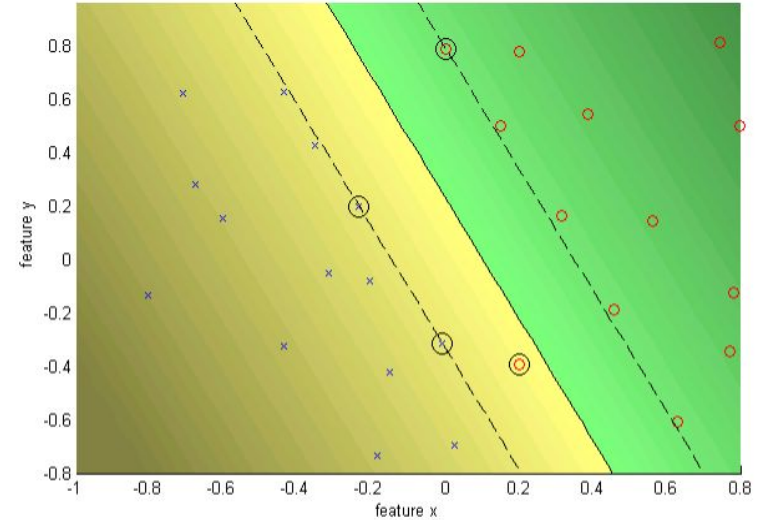
$C = \text{Infinity}$  hard margin



Comment Window

SVM (L1) by Sequential Minimal Optimizer  
Kernel: linear (-), C: Inf  
Kernel evaluations: 971  
Number of Support Vectors: 3  
Margin: 0.0966  
Training error: 0.00%

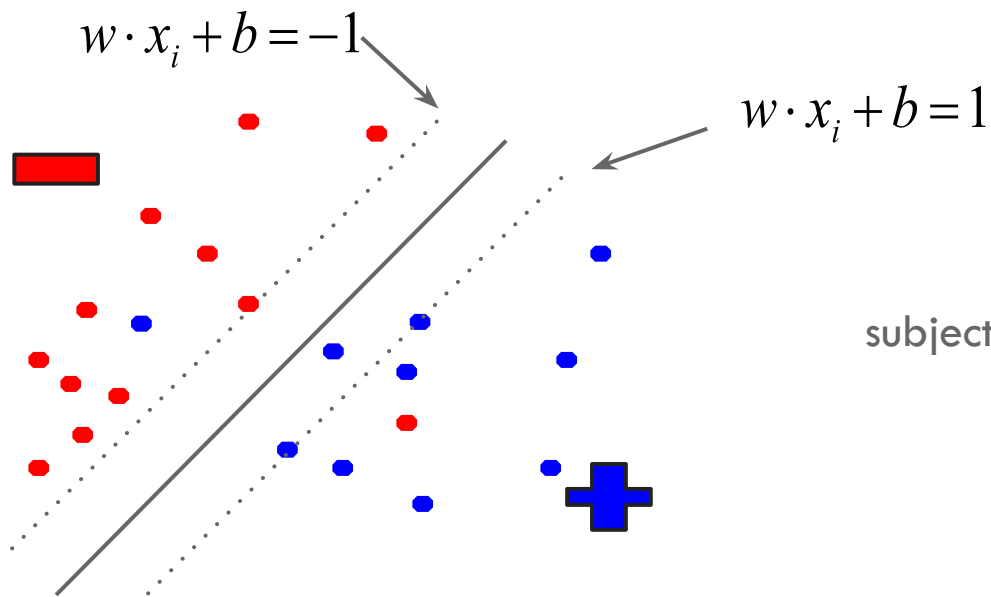
$C = 10$  soft margin



Comment Window

SVM (L1) by Sequential Minimal Optimizer  
Kernel: linear (-), C: 10.0000  
Kernel evaluations: 2645  
Number of Support Vectors: 4  
Margin: 0.2265  
Training error: 3.70%

# UNDERSTANDING THE SOFT MARGIN SVM



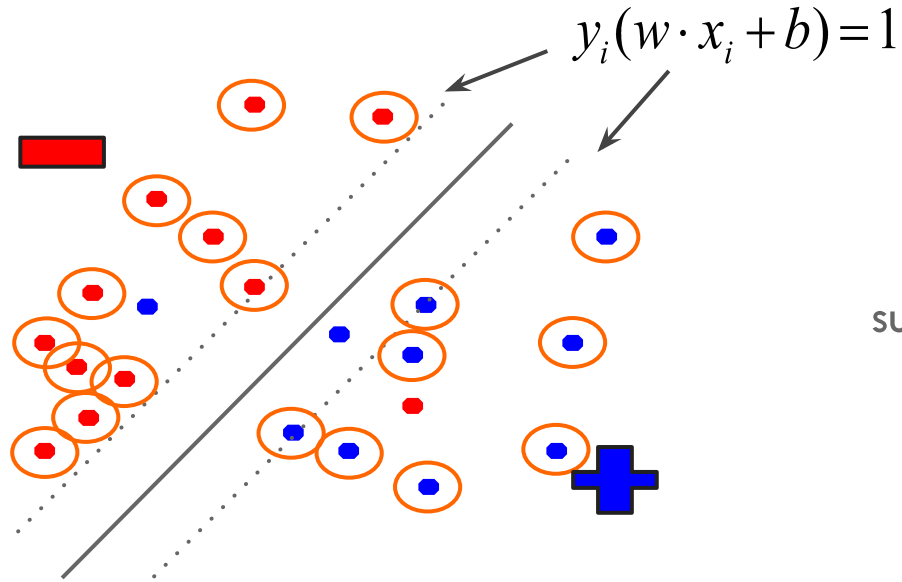
subject to:

$$\min_{w,b} \|w\|^2 + C \sum_i \zeta_i$$
$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

or:

$$y_i(w \cdot x_i + b) = 1$$

# UNDERSTANDING THE SOFT MARGIN SVM

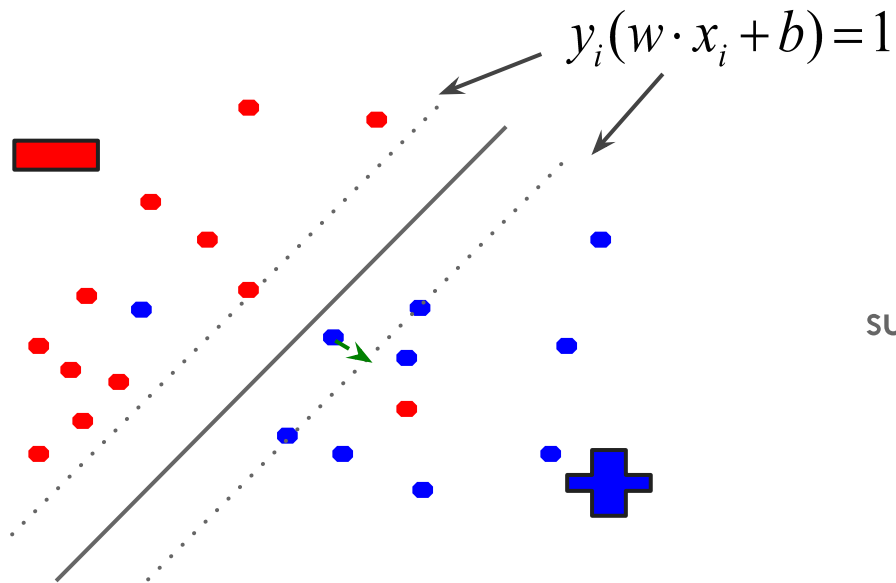


$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

slack values for points  $\geq$  margin and correctly classified is equal to:

0! The slack variables have to be greater than or equal to zero and if they are on or beyond the margin then  $y_i(w x_i + b) \geq 1$  already

# UNDERSTANDING THE SOFT MARGIN SVM

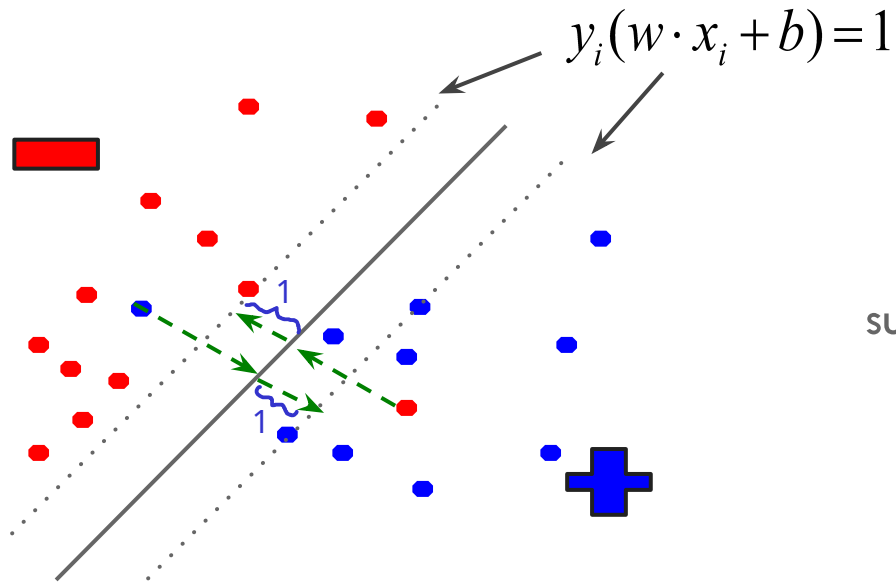


$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i \\ & \text{subject to:} \\ & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

Difference from the point to the margin, i.e.

$$\zeta_i = 1 - y_i(w \cdot x_i + b)$$

# UNDERSTANDING THE SOFT MARGIN SVM



$$\min_{w,b} \|w\|^2 + C \sum_i \zeta_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$

$$\zeta_i \geq 0$$

slack values for points incorrectly classified

“distance” to the hyperplane *plus* the “distance” to the margin

$$\zeta_i = 1 - y_i(w \cdot x_i + b)$$



# UNDERSTANDING THE SOFT MARGIN SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

---

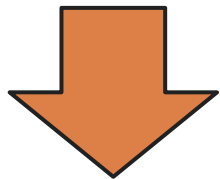
$$\varsigma_i = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & \text{otherwise} \end{cases}$$

tutti i punti classificati correttamente

tutti i punti classificati erroneamente

# UNDERSTANDING THE SOFT MARGIN SVM

$$\zeta_i = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & \text{otherwise} \end{cases}$$



$$\begin{aligned} \zeta_i &= \max(0, 1 - y_i(w \cdot x_i + b)) \\ &= \max(0, 1 - yy') \end{aligned}$$

hinge loss

riformulazione

# HINGE LOSS

**Hinge:**

$$l(y, y') = \max(0, 1 - yy')$$

**Squared loss:**

$$l(y, y') = (y - y')^2$$

**0/1 loss:**

$$l(y, y') = 1[yy' \leq 0]$$

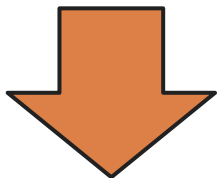
# UNDERSTANDING THE SOFT MARGIN SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

$$\zeta_i = \max(0, 1 - y_i(w \cdot x_i + b))$$



abbiamo riscritto  
SVM soft margin problem  
in una unconstrained form

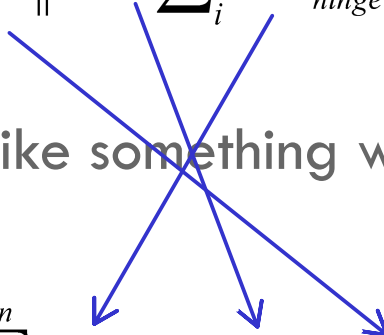
$$\min_{w,b} \quad \|w\|^2 + C \sum_i \max(0, 1 - y_i(w \cdot x_i + b))$$

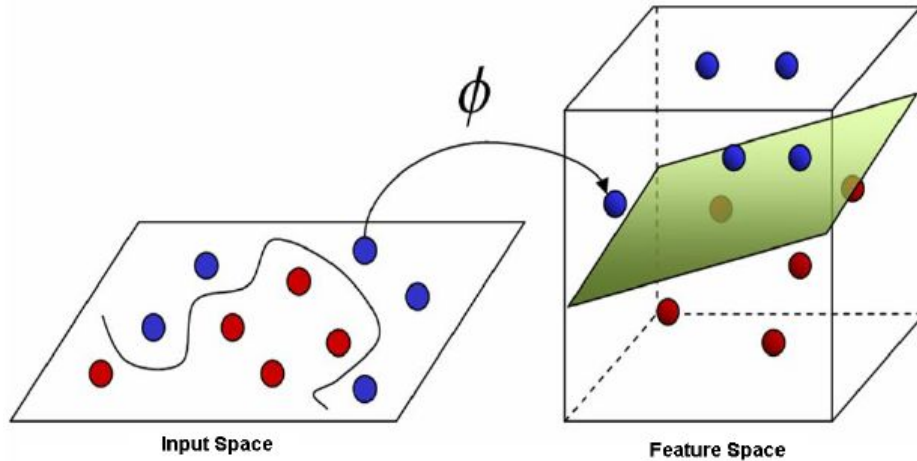
**Unconstrained problem!**

# UNDERSTANDING THE SOFT MARGIN SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \text{loss}_{\text{hinge}}(y_i, y_i')$$

Does this look like something we have seen before?

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \text{loss}(yy') + \lambda \text{ regularizer}(w, b)$$




NON LINEARLY  
SEPARABLE DATA

# SUPPORT VECTOR MACHINE PROBLEM

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

This is a version of a **quadratic optimization problem**

Maximize/minimize a quadratic function subject to a set of linear constraints

This is typically referred as **primal problem**

# RECAP: CLASSES OF OPTIMIZATION PROBLEMS

**Linear programming (LP):** linear problem, linear constraints

$$\begin{array}{ll} \min_{\mathbf{x}} & \mathbf{c}^T \mathbf{x} \text{ objective function} \\ \text{s.t.} & \mathbf{Ax} = \mathbf{b}, \quad \mathbf{x} \geq 0 \text{ constrain} \end{array}$$

**Quadratic programming (QP):** quadratic objective and linear constraints, it is convex if  $Q$  is positive semidefinite

$$\begin{array}{ll} \min_{\mathbf{x}} & \mathbf{c}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} \\ \text{s.t.} & \mathbf{Ax} = \mathbf{b}, \quad \mathbf{Cx} \geq \mathbf{d} \end{array} \quad \text{SVM}$$

**Nonlinear programming problem (NLP):** in general non-convex

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g(\mathbf{x}) = 0, \quad h(\mathbf{x}) \geq 0 \end{array}$$



# DUAL PROBLEM

- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.
- One possible solution involves constructing a dual problem where a Lagrange multiplier  $\alpha_i$  is associated with every inequality constraint in the primal (original) problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0, \forall i \end{aligned}$$

la matematica ci aiuta e ci permette di risolvere un problema equivalente

# THE SOLUTION

Given a solution  $\alpha_1 \dots \alpha_n$  to the dual problem, the solution to the primal is:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$b = y_k - \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_k$$

Each non-zero  $\alpha_i$  indicates that corresponding  $\mathbf{x}_i$  is a support vector.  
Then the classifying function is (note that we don't need  $\mathbf{w}$  explicitly):

$$f(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

# THE SOLUTION

$$f(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- Two important observations
  - The solution relies on an inner product between the test point  $\mathbf{x}$  and the support vectors  $\mathbf{x}_i$ .
  - Solving the optimization problem involves computing the inner products between all training points.

# DUAL PROBLEM WITH SOFT MARGIN

- Dual problem is similar in the non separable case but notice the constraints.

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

- Again,  $\mathbf{x}_i$  with non-zero  $\alpha_i$  will be support vectors.

# LINEAR SVM SUMMARY

- The classifier is a separating hyperplane.
- Most “important” training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points are support vectors with non-zero Lagrangian multipliers  $\alpha_i$ .

In inference phase soli i support vector points vengono considerati

# LINEAR SVM SUMMARY

- Both in the dual formulation of the problem and in the solution training points appear only inside inner products:

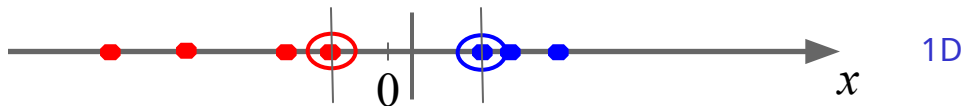
$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

$$f(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

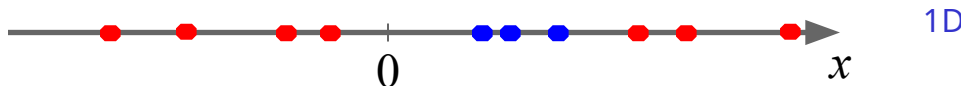
# NON LINEAR SVM

passando ad uno spazio a più dimensioni possiamo trovare un piano che separi dei punti che nello spazio attuale non sono linearmente separabili

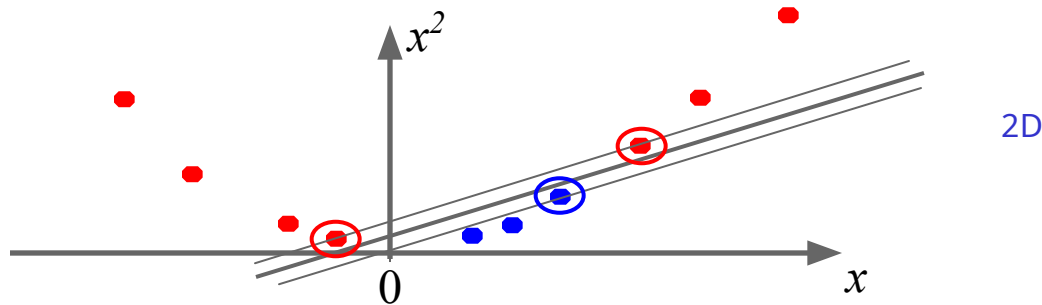
- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?

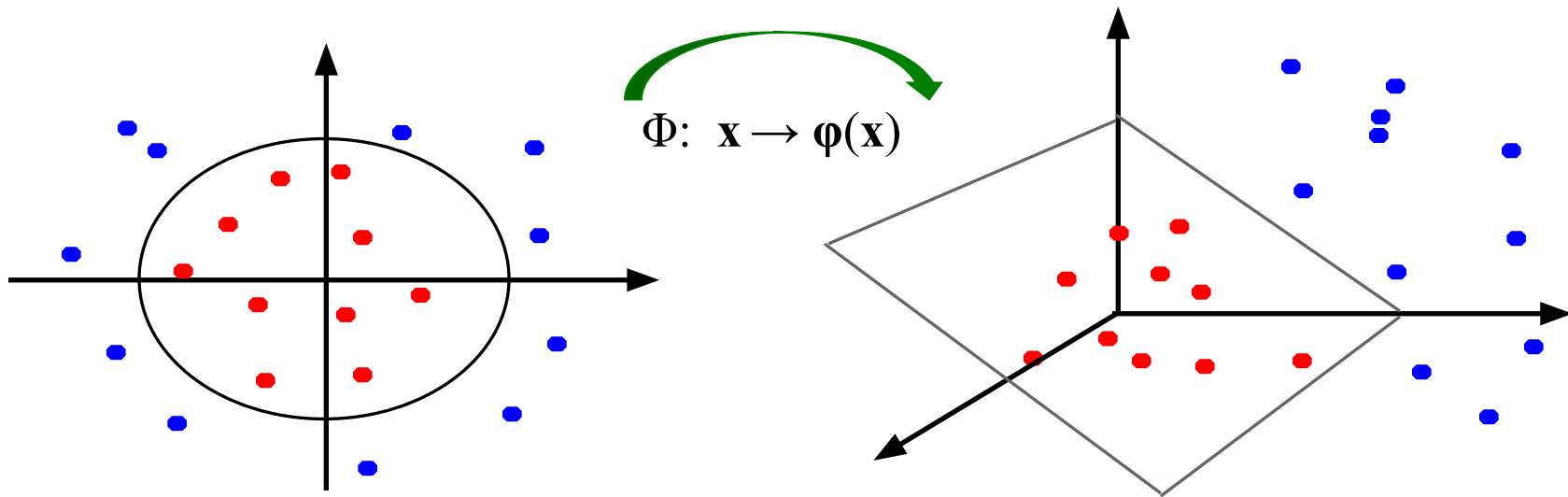


- How about... mapping data to a higher-dimensional space?



# NON LINEAR SVM: FEATURE SPACES

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:





# KERNEL TRICK

- The linear classifier relies on inner product between vectors  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- If every datapoint is mapped into high-dimensional space via some , transformation  $\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ , the inner product becomes:

$\phi$  rimane sembra implicita

noi lavoriamo con il kernel dato dalla definizione  $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$

- A kernel function is a function that is equivalent to an inner product in some feature space.

Noi scegliamo solo la kernel function

# KERNEL TRICK

Example:

2-dimensional vectors  $\mathbf{x}=[x_1 \ x_2]$ ; let  $K(\mathbf{x}_i, \mathbf{x}_j)=(1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ ,

Need to show that  $K(\mathbf{x}_i, \mathbf{x}_j)=\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ :

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 = 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} = \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] = \\ &= \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j), \quad \text{where } \varphi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2] \end{aligned}$$

- A kernel function **implicitly** maps data to a high-dimensional space (without the need to compute each  $\varphi(\mathbf{x})$  explicitly).

# KERNELS

FYI: Ricerca continua di valid kernel function

- For some functions  $K(\mathbf{x}_i, \mathbf{x}_j)$  checking that  $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$  can be cumbersome.
- Mercer's theorem:
  - Every positive semidefinite symmetric function is a kernel
  - A positive semidefinite symmetric functions correspond to a positive semidefinite symmetric Gram matrix:

K=

$K(\mathbf{x}_1, \mathbf{x}_1)$	$K(\mathbf{x}_1, \mathbf{x}_2)$	$K(\mathbf{x}_1, \mathbf{x}_3)$	...	$K(\mathbf{x}_1, \mathbf{x}_n)$
$K(\mathbf{x}_2, \mathbf{x}_1)$	$K(\mathbf{x}_2, \mathbf{x}_2)$	$K(\mathbf{x}_2, \mathbf{x}_3)$		$K(\mathbf{x}_2, \mathbf{x}_n)$
...	...	...	...	...
$K(\mathbf{x}_n, \mathbf{x}_1)$	$K(\mathbf{x}_n, \mathbf{x}_2)$	$K(\mathbf{x}_n, \mathbf{x}_3)$	...	$K(\mathbf{x}_n, \mathbf{x}_n)$

- Recap: A symmetric matrix is positive semidefinite if and only if all eigenvalues are non-negative

# KERNELS

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial of power  $p$ :  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
- Gaussian (radial-basis function):  $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$ 
  - Mapping  $\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ , where  $\varphi(\mathbf{x})$  is *infinite-dimensional*

# NON LINEAR SVM PROBLEM

- Dual problem formulation:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0, \forall i \end{aligned}$$

- The solution is:

$$f(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

- Optimization techniques for finding  $\alpha_i$ 's remain the same!

# SVM REMARKS

- Most popular optimization algorithms for SVMs use decomposition to hill-climb over a subset of  $\alpha_i$ 's at a time, e.g. SMO [Platt '99] and [Joachims '99]
- Tuning SVMs remains a black art: selecting a specific kernel and parameters is usually done in a try-and-see manner (grid search)

SVM è diventato così tanto prevalente perchè:

- ha garanzie teoriche molto potenti di generalizzazione
- è super flessibile e può essere usato per tanti scopi (anche oltre la classification)

# SVM APPLICATIONS

## Pedestrian detection in Computer Vision

Objective: detect (localize) standing humans in an image



- reduces object detection to binary classification
- does an image window contain a person or not?

# QUESTIONS?



Some slides are taken from David Kauchak