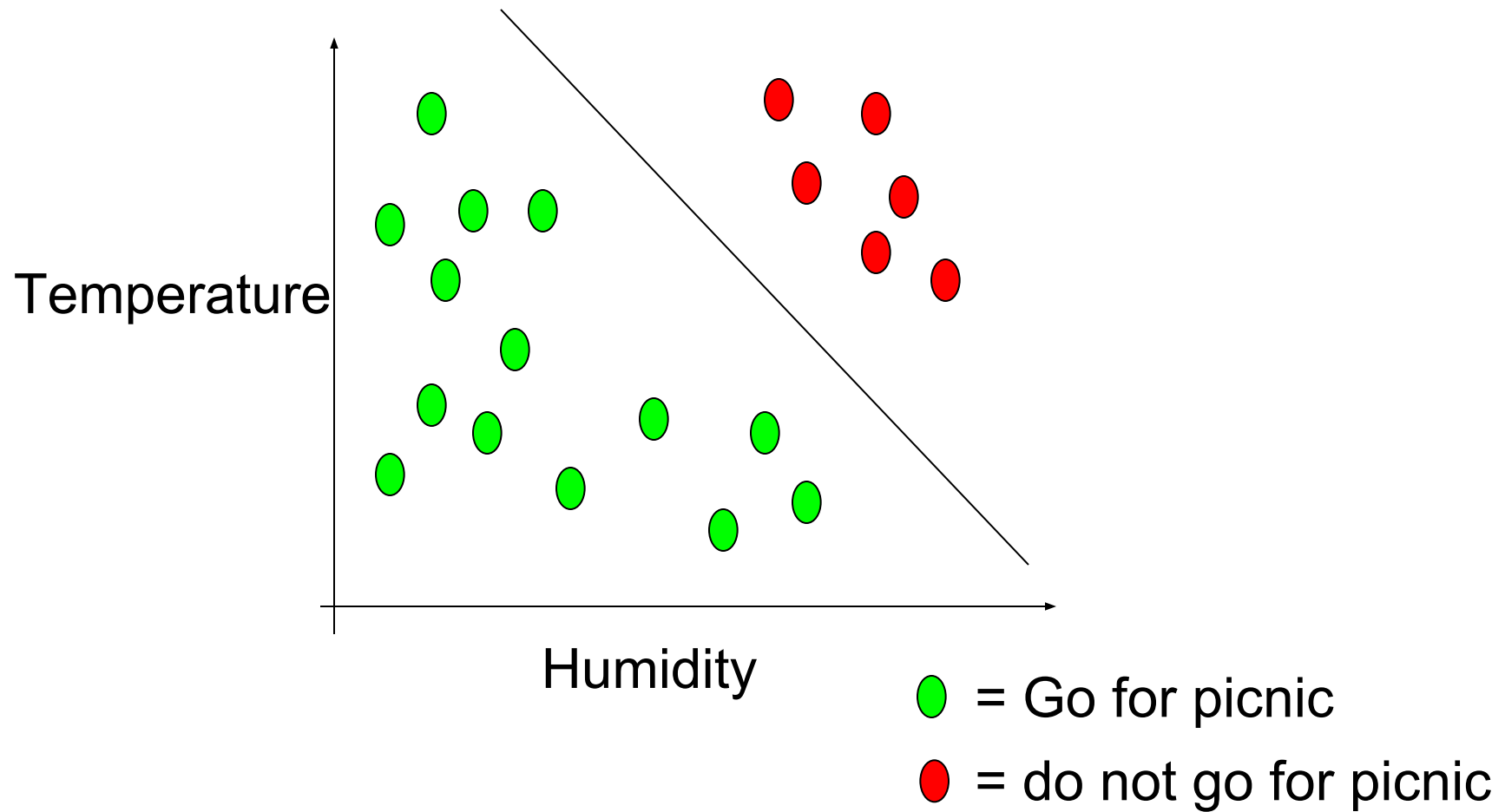# Support Vector Machines

# Main Ideas

- Supervised, binary, linear, discriminative
  - Used for classification or regression
- Optimal hyperplane
  - Outputs an optimal hyperplane to separate +ve and –ve training examples

- Maximize Margin of separation Classifier
  - Formalizes notion of the best linear separator

- Kernels
  - Projecting non-linearly separable data into higher-dimensional space makes it linearly separable
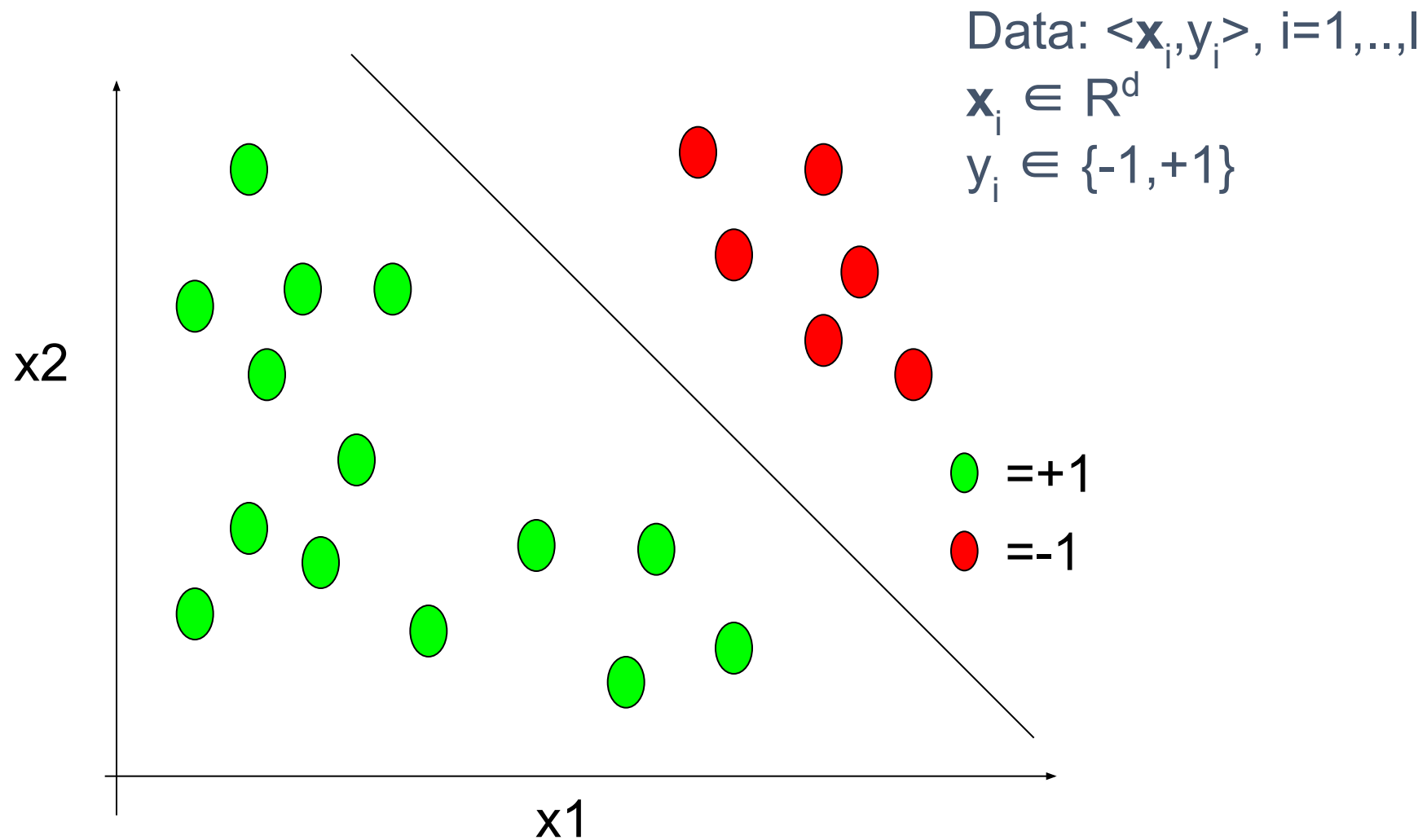
# Main Ideas

- Lagrangian Multipliers
  - Way to convert a constrained optimization problem to one that is easier to solve

- Complexity
  - Depends only on the number of training examples, not on dimensionality of the kernel space!

- Regularization
  - Has a regularization parameter to relax optimization constraints for non separable data
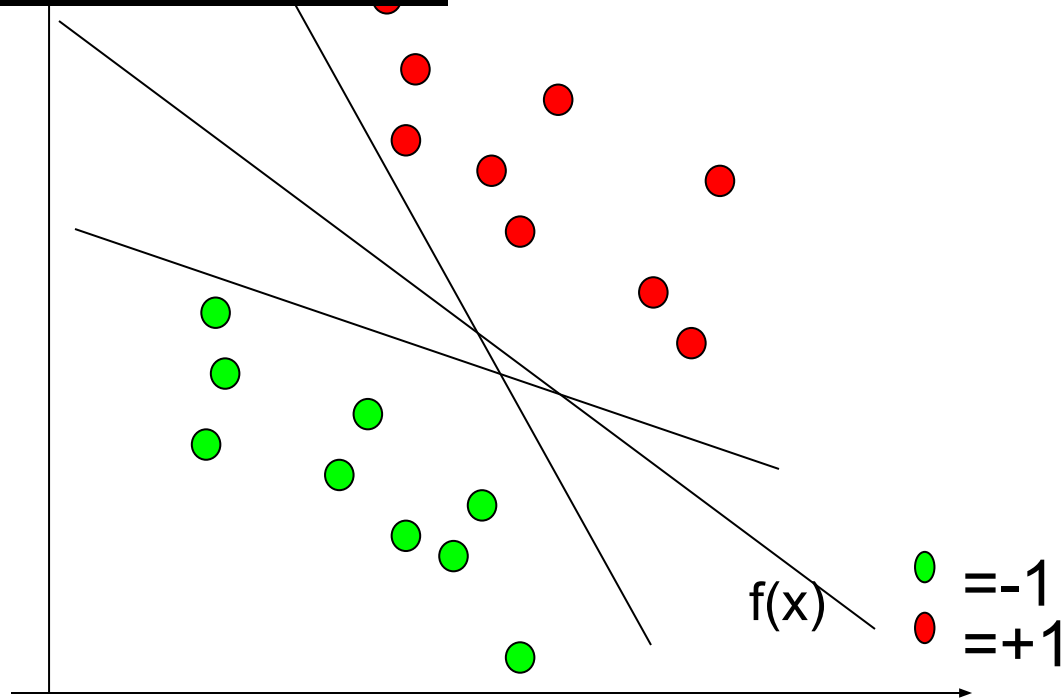
# Picnic example



Temperature

Humidity

🟢 = Go for picnic

🔴 = do not go for picnic

# Linear Support Vector Machines



Data: $<\mathbf{x}_i, y_i>$, i=1,..,l
$\mathbf{x}_i \in R^d$
$y_i \in \{-1,+1\}$

● =+1
● =-1

x2

x1

# Linear SVM 2



Data: $\langle \mathbf{x}_i, y_i \rangle$, i=1,..,l

$\mathbf{x}_i \in R^d$

$y_i \in \{-1,+1\}$

$\circ$ =-1
$\bullet$ =+1

f(x)

All hyperplanes in $R^d$ are parameterize by a vector (**w**) and a constant b.
Can be expressed as $\mathbf{w}^T \bullet \mathbf{x} + b = 0$

$\mathbf{W} = \begin{Bmatrix} \underline{w1} \\ \underline{w2} \\ \underline{\;} \end{Bmatrix}$ $\quad \boldsymbol{x} = \begin{Bmatrix} \underline{x1} \\ \underline{x2} \\ \underline{\;} \end{Bmatrix}$

Our aim is to find such a hyperplane  f(x)=sign(**w**•**x**+b),
that correctly classify our data.

# Definitions

Define the hyperplane H such that:

$x_i \cdot w + b \geq +1$ when $y_i = +1$

$x_i \cdot w + b \leq -1$ when $y_i = -1$

H1 and H2 are the planes:
H1: $x_i \cdot w + b = +1$
H2: $x_i \cdot w + b = -1$
The points on the planes
H1 and H2 are the
Support Vectors



H1

H2

$d^+$

$d^-$

H

$w \cdot x + b = +1$

$w \cdot x + b = 0$

$w \cdot x + b = -1$

d+ = the shortest distance to the closest positive point

d- = the shortest distance to the closest negative point

The <u>margin</u> of a separating hyperplane is $d^+ + d^-$.
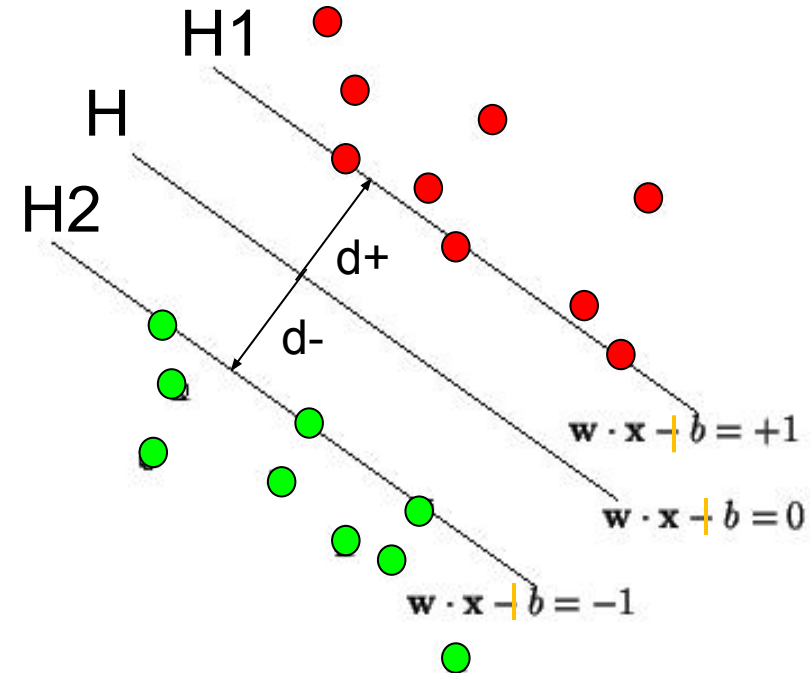
# Maximizing the margin

We want a classifier with as big margin as possible.

Recall the distance from a point$(x_0, y_0)$ to a line:
$Ax+By+c = 0$ is $|A x_0 + B y_0 + c|/sqrt(A^2+B^2)$

The distance between H and H1 is:
$|\mathbf{w} \bullet \mathbf{x} + b|/||w|| = 1/||w||$

The distance between H1 and H2 is: $2/||w||$



H1

H

H2

d+

d-

$\mathbf{w} \cdot \mathbf{x} + b = +1$

$\mathbf{w} \cdot \mathbf{x} + b = 0$

$\mathbf{w} \cdot \mathbf{x} + b = -1$

**In order to maximize the margin, we need to minimize ||w||. With the condition that there are no datapoints between H1 and H2:**

$\mathbf{x}_i \bullet \mathbf{w} + b \geq +1$ when $y_i = +1$
$\mathbf{x}_i \bullet \mathbf{w} + b \leq -1$ when $y_i = -1$ **Can be combined into $y_i(\mathbf{x}_i \bullet \mathbf{w}) \geq 1$**
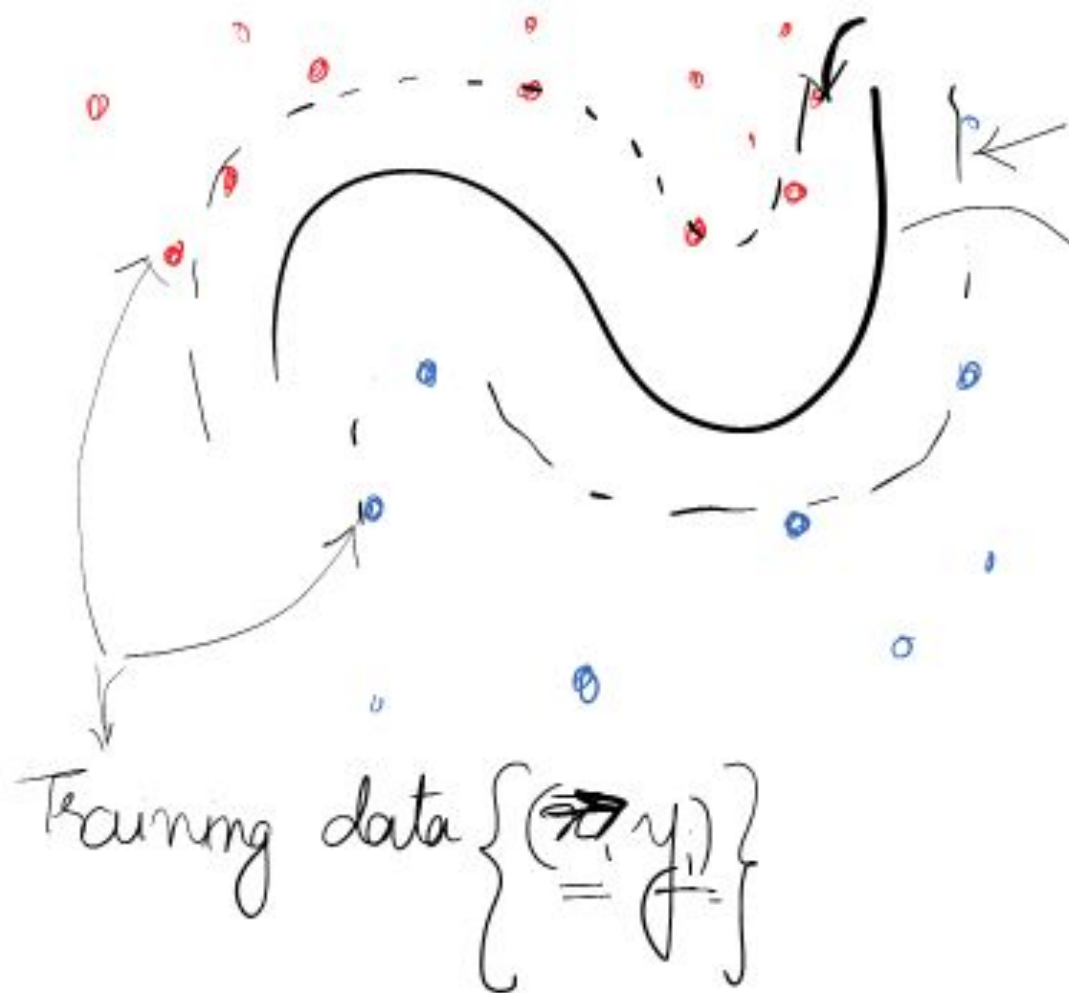
SVM – Maximize margin of separation $\rho$



Euclidean norm

$$\rho = \frac{2}{\|\omega_0\|}$$

$$\vec{\omega_0}^T \vec{x} + b = 0$$

$\vec{\omega_0}^T \vec{x}_i + b \geq 1 \quad \text{if } y_i = +1$

$\vec{\omega_0}^T \vec{x}_i + b_0 \leq -1 \quad \text{if } y_i = -1$

Training data $\{(\vec{x}_i, y_i)\}$

Optimization problem

(A) OBJECTIVE FUNCTION

$$\phi(\vec{\omega}) = \frac{1}{2}\omega^T\omega$$

$$P = \frac{2}{||\omega_{.}||}$$

(B) CONSTRAINTS FULFILL

$$y_i(\vec{\omega}^T\vec{x}_i + b) \geq 1$$
$$for\ i = 1 \dots N$$

COMBINE (A) & (B) using

Lamgrange Multipliers $\alpha$

(M) $$J(\vec{\omega}, b, \alpha) = \left(\frac{1}{2}\omega_i^T\omega_i\right) - \left(\sum_{i=1}^{N}\alpha_i\left[y_i(\omega_i^T\cdot\vec{x}_i + b) - 1\right]\right)$$

## An Optimization problem and its dual

$$\partial J/\partial w = 0 \qquad \partial J/\partial b = 0$$

yields

$$\vec{w} = \sum_{i=1}^{N} \alpha_i y_i \vec{x_i} \qquad\text{——}① $$
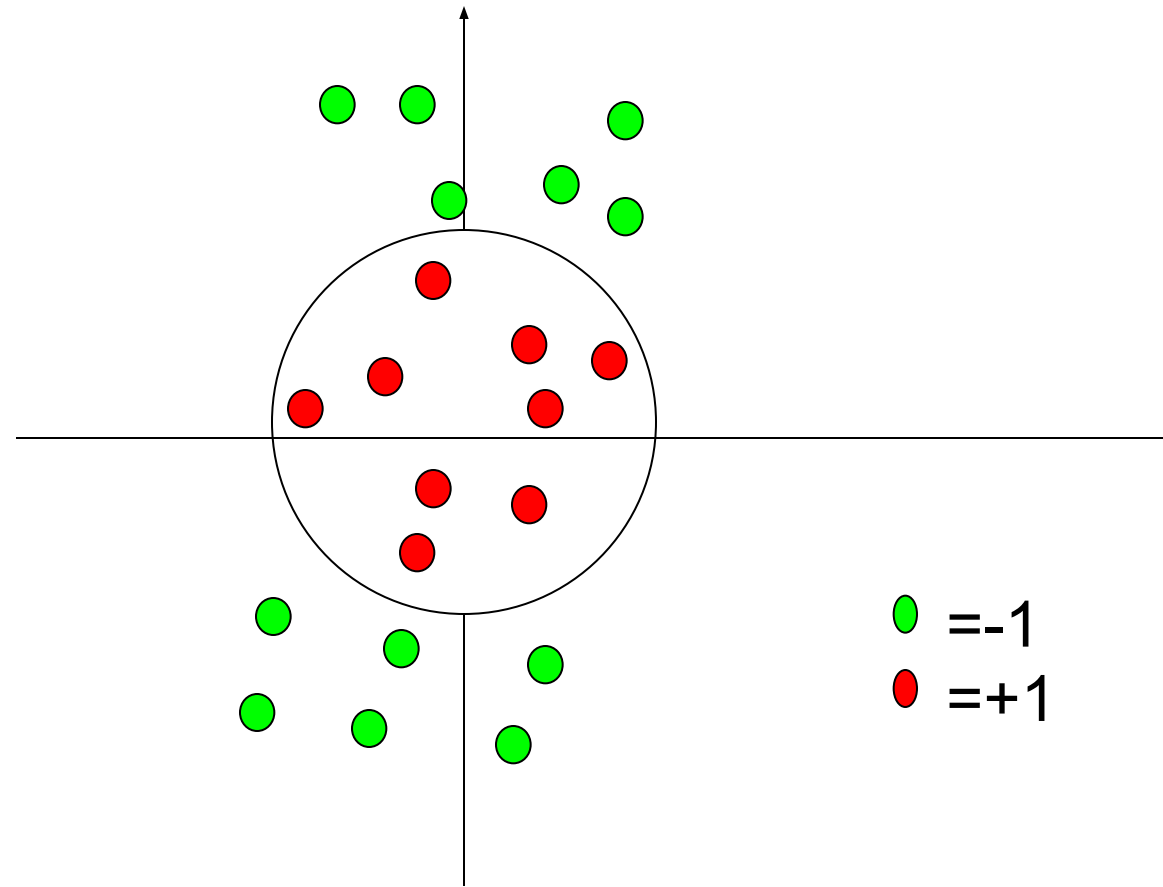
$$\sum_{i=1}^{N} \alpha_i y_i = 0 \qquad\text{——}②$$

# Quadratic Programming

- Why is this reformulation a good thing?
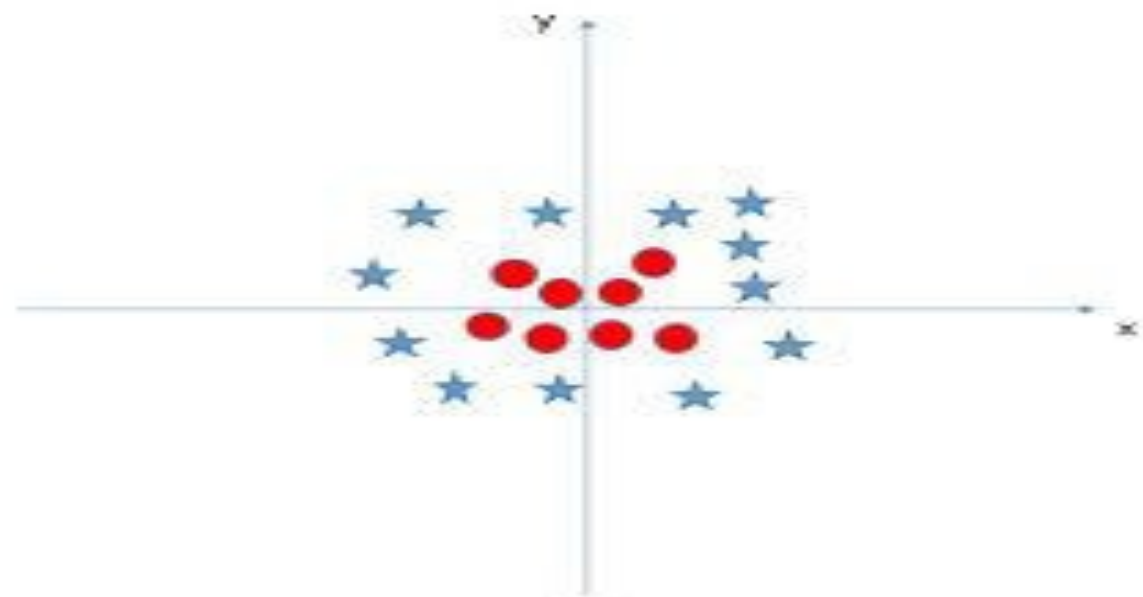
- The problem

$$\text{Maximize} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i \cdot \mathbf{x}_j \rangle$$

$$\text{subject to} \sum_i y_i \alpha_i = 0 \ \text{ and } \ \alpha_i \geq 0$$
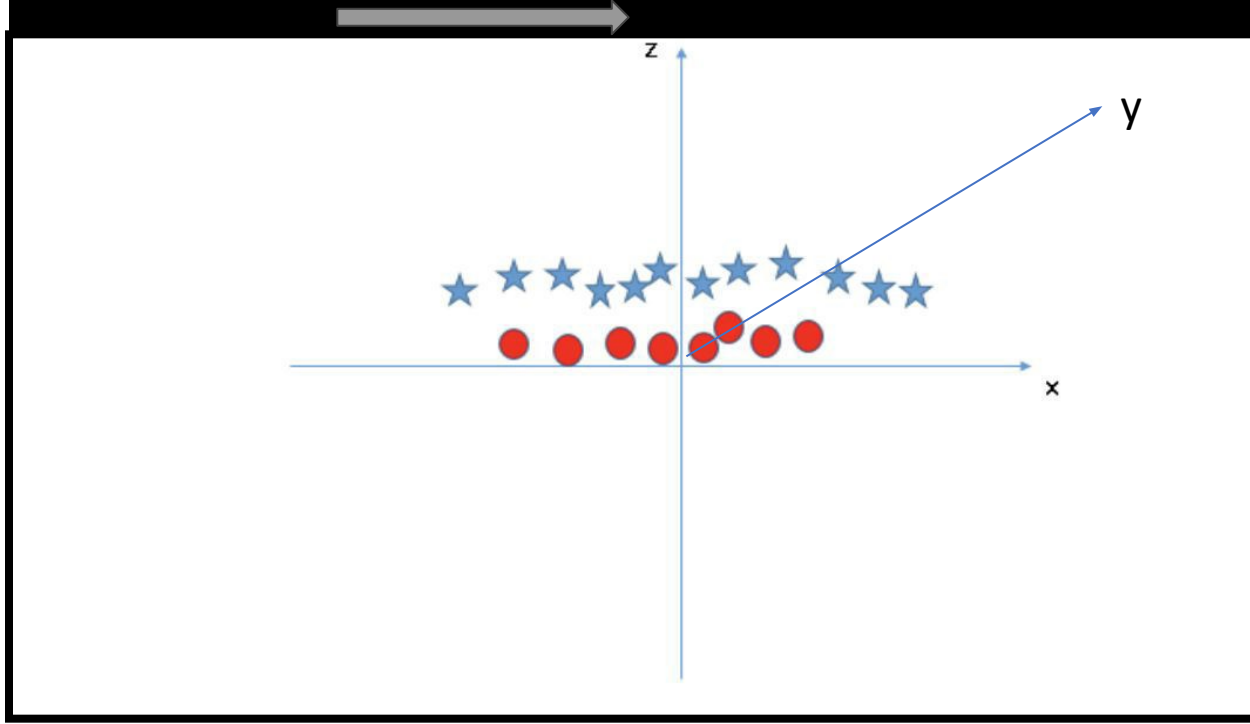
- is an instance of what is called a positive, semi-definite programming problem

- For a fixed real-number accuracy, can be solved in time order = $O(|D|^2 \log |D|^2)$

# Problems with linear SVM



=-1
=+1

What if the decision function is not a linear?

# The kernel trick

- For many mappings from a low-D space to a high-D space, there is a simple operation on two vectors in the low-D space that can be used to compute the scalar product of their two images in the high-D space.
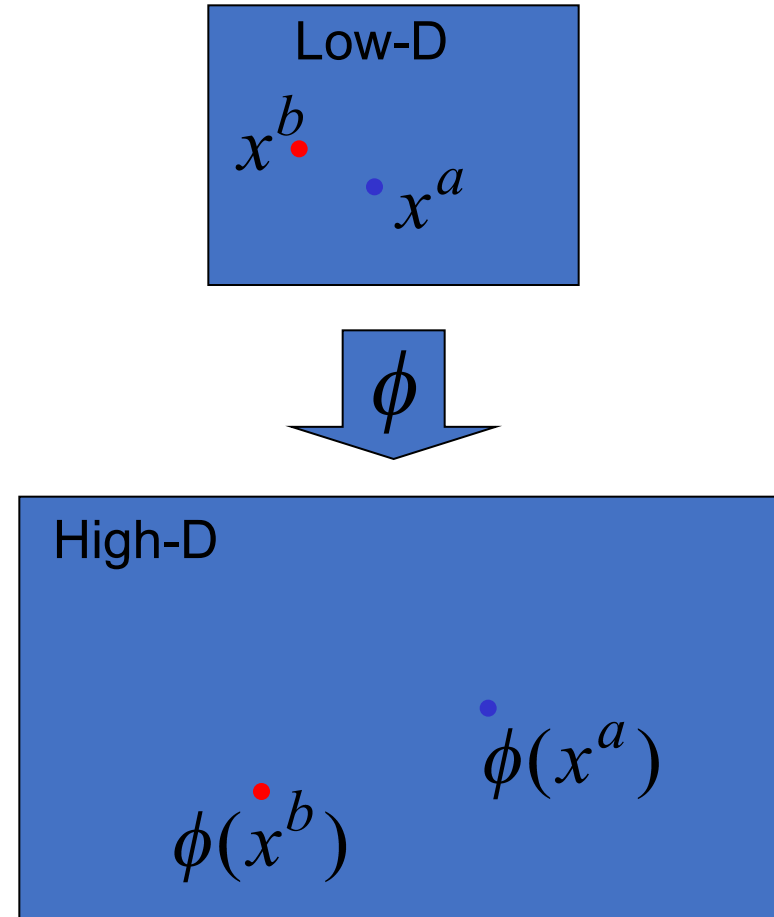
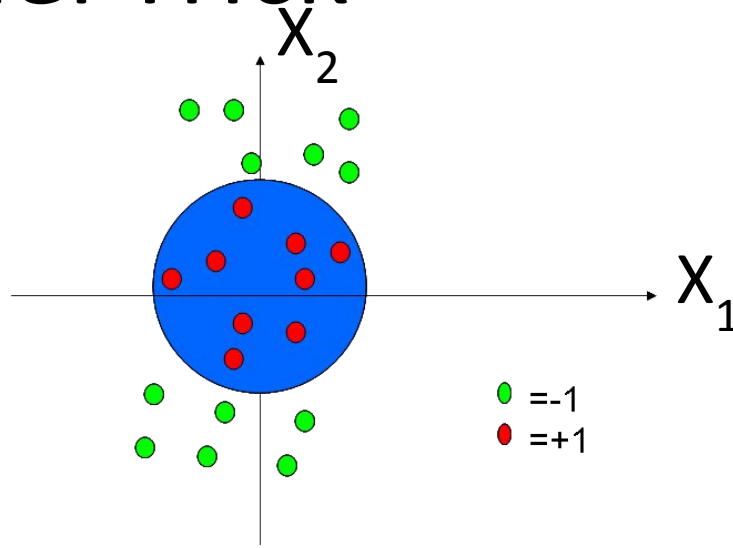$$K(x^a, x^b) = \phi(x^a) . \phi(x^b)$$

<span style="color:blue">↑</span>    <span style="color:blue">↑</span>
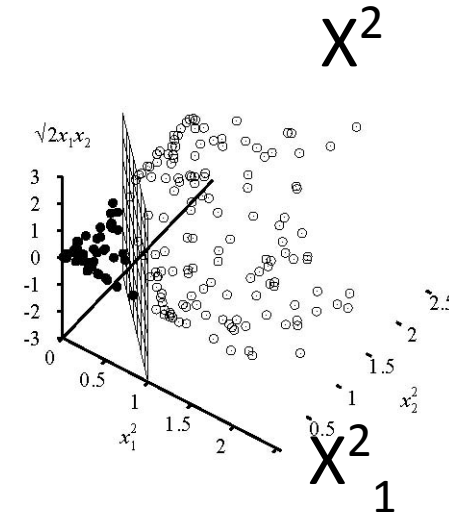
<span style="color:blue">Letting the kernel do the work</span>    <span style="color:blue">doing the scalar product in the obvious way</span>

# Kernel Trick



$X_2$

$X_1$

● =-1
● =+1

$X^2$

$\sqrt{2}x_1x_2$

$X^2_1$

Data points are linearly separable

in the space $(x_1^2, x_2^2, \sqrt{2}x_1x_2)$

We want to maximize $\sum_i \alpha_i - \frac{1}{2}\sum_{i,j} y_i y_j \alpha_i \alpha_j \langle F(\mathbf{x}_i) \cdot F(\mathbf{x}_j)\rangle$

Define $K(\mathbf{x}_i, \mathbf{x}_j) = \langle F(\mathbf{x}_i) \cdot F(\mathbf{x}_j)\rangle$

$K$ is often easy to compute directly!

Here,

$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i \cdot \mathbf{x}_j\rangle^2$
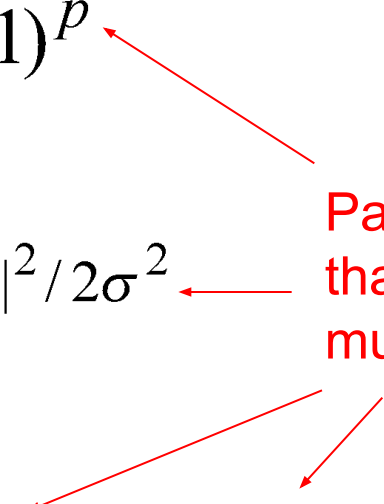
# Some commonly used kernels

Polynomial: $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}.\mathbf{y} + 1)^p$

Gaussian radial basis function $K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2 / 2\sigma^2}$

Neural net: $K(\mathbf{x}, \mathbf{y}) = \tanh\left(k\,\mathbf{x}.\mathbf{y} - \delta\right)$

Parameters that the user must choose

Kernel M should be   1) real   2) symmetric

3) +ve, semi-definite   $z^T M z$   is strictly +ve

E.g.   Identity Matrix

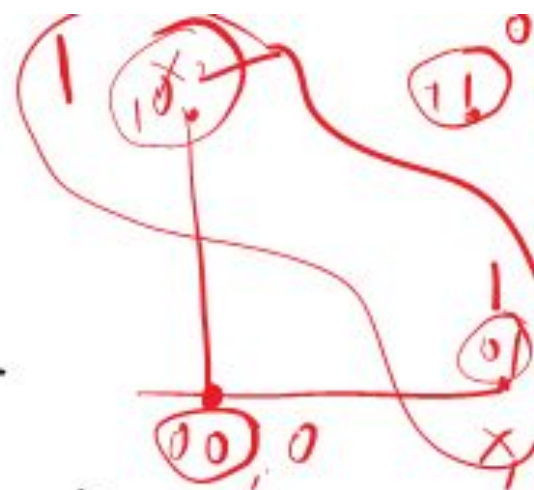$$[a \quad b] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a^2 + b^2$$

$z^T$   $z$   $\rightarrow$ +ve

XOR example

$$Input\ vector\ \vec{x}$$

| | $x_1$ | $x_2$ | label $y$ |
|---|---|---|---|
| $x^1$ | $-1$ | $-1$ | $-1$ |
| $x^2$ | $-1$ ✓ | $+1$ ✓ | $+1$ |
| $x^3$ | $+1$ | $-1$ | $+1$ |
| $x^4$ | $+1$ | $+1$ | $-1$ |

Let

$$K(\vec{x}, \vec{x_i})$$

$$= \left(1 + \vec{x}^T \vec{x_i}\right)^2$$

$$\Rightarrow \left(1 + \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}\right)^2 \Rightarrow \left(1 + \overline{x_1 x_{i1} + x_2 x_{i2}}\right)^2 \Rightarrow$$

$$\overline{1 + \left(x_1 x_{i1} + x_2 x_{i2}\right)^2 + 2\left(x_1 x_{i1} + x_2 x_{i2}\right)}$$

$$\Rightarrow 1 + \boxed{x_1^2}\boxed{x_{i1}^2} + 2 x_1 x_2 x_{i1} x_{i2} + x_2^2 x_{i2}^2 + 2 x_1 x_{i1} + 2 x_2 x_{i2}$$

Image of input vector x in high dimensional feature space

$$\varphi(\vec{x}) = \left[1, \quad x_1^2, \sqrt{2}\, x_1 x_2, \quad x_2^2, \quad \sqrt{2}\, x_1, \quad \sqrt{2}\, x_2\right]^T$$

$\varphi(x_i) \cdot \varphi(x_j)$

$$\varphi(\vec{x_i}) = \left[1 \quad x_{i1}^2, \quad \sqrt{2} x_{i1} x_{i2}, \quad x_{i2}^2, \quad \sqrt{2} x_{i1}, \quad \sqrt{2} x_{i2}\right]$$

$$K[1,1] = 1 + (-1)^2(-1)^2 + 2(-1)(-1)(-1)(-1) + (-1)^2(-1)^2 + 2(-1)(-1) + 2$$

$x_1 \; x_1$
$$= 9 \qquad\qquad\qquad (-1)(-1)$$

Likewise  calculate   $K[1,2], \quad K[1,3] \quad K[1,4] \quad K[2,2], K[2,3]$

$x_1 \;(-1 \quad -1) \quad -1 \qquad\qquad K[2,4], \quad K[3,3], \; K[3,4], \; K[4,4]$

$K[1,2] = 1 + (-1)^2(-1)^2 + 2\ (-1)(-1)(-1)(+1) + (-1)^2(1)^2 + 2(-1)(-1) + 2(-1)1$

$$K = \begin{bmatrix} 9 & 1 & 1 & 1 \\ 1 & 9 & 1 & 1 \\ 1 & 1 & 9 & 1 \\ 1 & 1 & 1 & 9 \end{bmatrix} \longrightarrow \text{KERNEL}$$

$$Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i \, y_j \, K(x_i, x_j)$$

Now $Q(\alpha) = (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) - \frac{1}{2} \Big( 9\alpha_1^2 - 2\alpha_1\alpha_2 - 2\alpha_1\alpha_3$

$$\alpha_1 \alpha_1 (-1)(-1) \quad 9$$
$$\alpha_i \alpha_i \, y_i \, y_i \, K(x_i, x_i)$$

$$+ 2\alpha_1\alpha_4 + 9\alpha_2^2 + 2\alpha_2\alpha_3 - 2\alpha_2\alpha_4 + 9\alpha_3^2 - 2\alpha_3\alpha_4 + 9\alpha_4^2$$

Note $Q$ can be found only in terms of $\alpha_s$

optimize $Q(\alpha)$ — w.r.t. $\alpha_1 \, \alpha_2 \, \alpha_3 \, \alpha_4$ :

$$9\alpha_1 - \alpha_2 - \alpha_3 + \alpha_4 = 1 \quad —①$$

$$-\alpha_1 + 9\alpha_2 + \alpha_3 - \alpha_4 = 1 \quad —②$$

$$-\alpha_1 + \alpha_2 + 9\alpha_3 - \alpha_4 = 1 \quad —③$$

$$\alpha_1 - \alpha_2 - \alpha_3 + 9\alpha_4 = 1 \quad —④$$

Solving $\quad \alpha_{0,1} = \alpha_{0,2} = \alpha_{0,3} = \alpha_{0,4} = 1/8$

$$Q_0(\alpha) = \frac{1}{4}$$

We need to calculate equation of plane $\rightarrow$ $\omega$ $b$ $Eq$.

$$Q(\alpha) = \frac{1}{2}\|w_0\|^2 = \frac{1}{4} \text{ :because constraints vanish.}$$

$$\Rightarrow \quad \|w_0\| = \frac{1}{\sqrt{2}} \quad \text{———————} \quad \textcircled{1}$$

$N_\omega$

$$\vec{w_0} = \sum \alpha_{0i} y_i \varphi(x_i) \qquad \overset{w_0}{\longrightarrow}$$

$$\frac{1}{8}\left[-1\begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \\ -\sqrt{2} \\ -\sqrt{2} \end{bmatrix} +1\begin{bmatrix} 1 \\ -\sqrt{2} \\ 1 \\ -\sqrt{2} \\ 2 \end{bmatrix} +1\begin{bmatrix} 1 \\ -\sqrt{2} \\ 1 \\ \sqrt{2} \\ -\sqrt{2} \end{bmatrix} -1\begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \\ \sqrt{2} \end{bmatrix}\right] = \begin{bmatrix} 0 \\ 0 \\ -1/\sqrt{2} \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

The Optimal hyperplane is

$$w_o^T$$

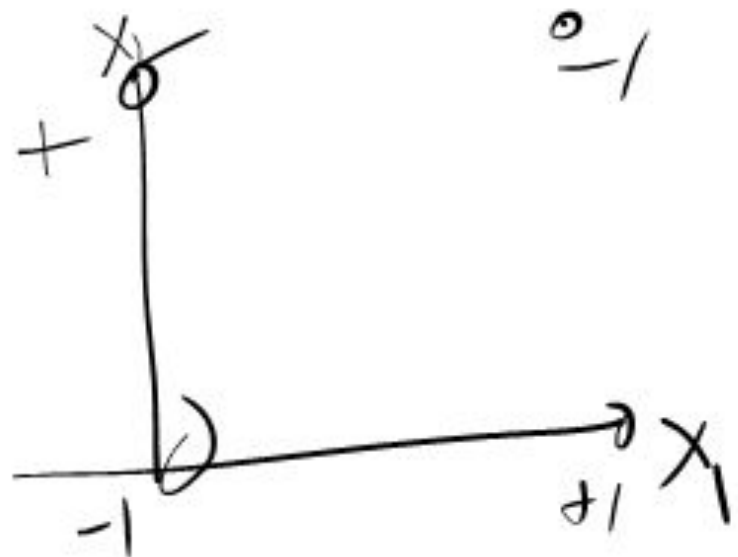$$w_o^T \cdot \varphi(x) = 0$$

$$\varphi(x) = 0 \qquad y = \boxed{mx + b}$$

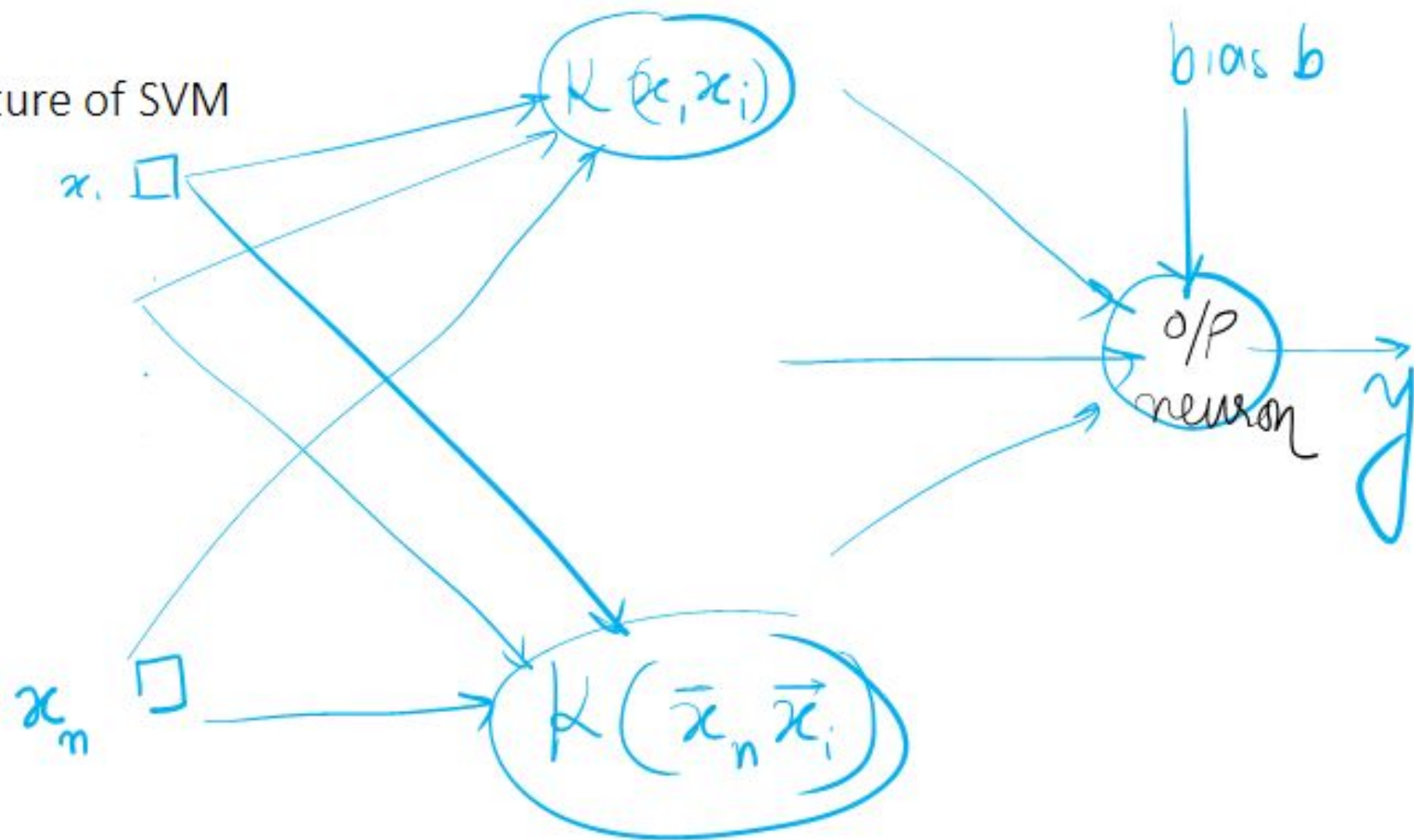$$\begin{bmatrix} 0 & 0 & -\dfrac{1}{\sqrt{2}} & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ x_1^2 \\ \sqrt{2}\,x_1 x_2 \\ x_2^2 \\ \sqrt{2}\,x_1 \\ \sqrt{2}\,x_2 \end{bmatrix} = 0 \quad \Leftrightarrow \quad \underline{-x_1 x_2 = 0}$$

$$\overset{0}{-1}$$

$$\dfrac{1 \cdot 0 \quad \left( \begin{matrix} -1, & +1 \\ +1, & -1 \end{matrix} \right)}{\wedge}$$

$$-x_1 x_2 = 0$$

$$\underline{\vee} \quad +0 \quad \left( \begin{matrix} +1 & +1 \\ -1 & -1 \end{matrix} \right)$$

Architecture of SVM



$K(\vec{x}, x_i)$

$x_i \ \square$
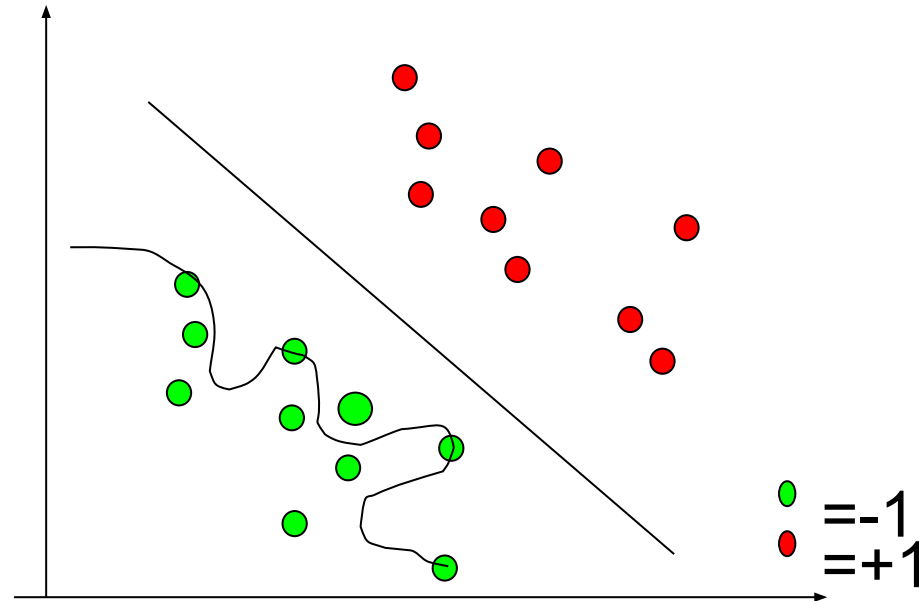
bias $b$

o/p
neuron

$y$

$x_n \ \square$

$K(\vec{x}_n, \vec{x}_i)$

# Overtraining/overfitting

A well known problem with machine learning methods is overtraining.
This means that we have learned the training data very well, but
we can not classify unseen examples correctly.

An example: A environmentalist knows rivers very well.
Everytime he sees a new river, he says it is not a river!

# Overtraining/overfitting 2

A measure of the risk of overtraining with SVM (there are also other measures).

It can be shown that: The portion, n, of unseen data that will be missclassified is bounded by:

$$n \leq \text{Number of support vectors / number of training examples}$$

Ockham´s razor principle: Simpler system are better than more complex ones. In SVM case: fewer support vectors mean a simpler representation of the hyperplane.
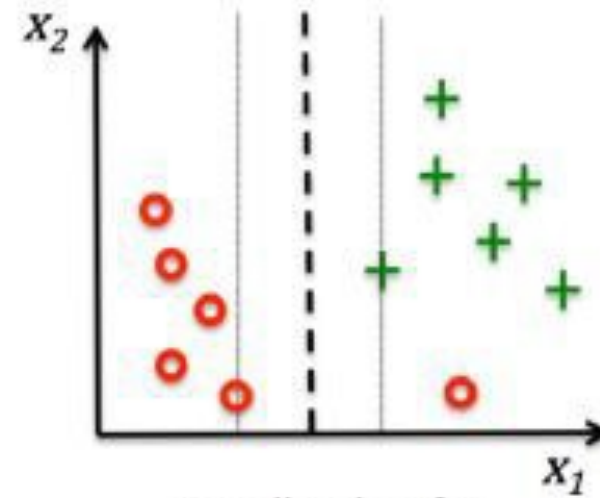
Example: Understanding a certain cancer if it can be described by one gene is easier than if we have to describe it with 5000.
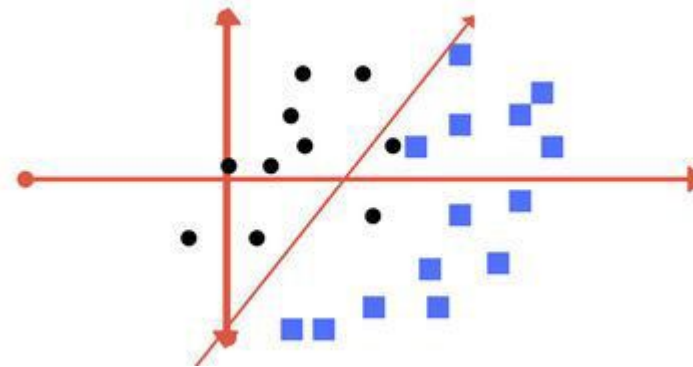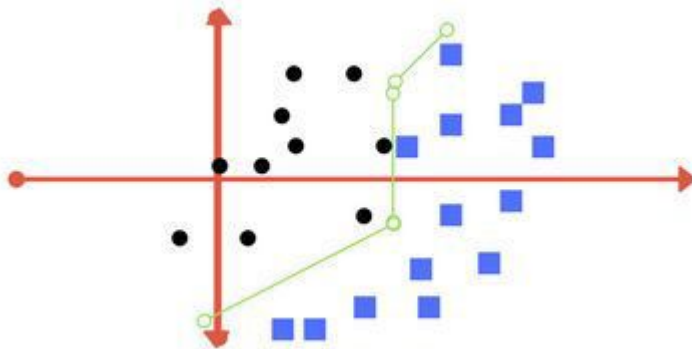
# Regularization

- Also the ' C ' parameter in Python's SkLearn Library
- Optimises SVM classifier to avoid misclassifying the data.

- C → large
- C → small

Margin of hyperplane → small

Margin of hyperplane → large

- misclassification(possible)

1. C ---> large , chance of overfit
2. C ---> small , chance of underfitting
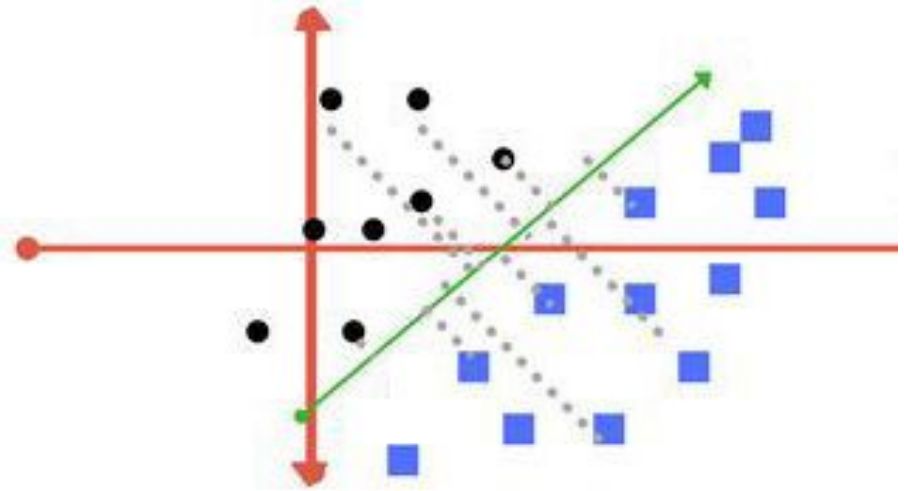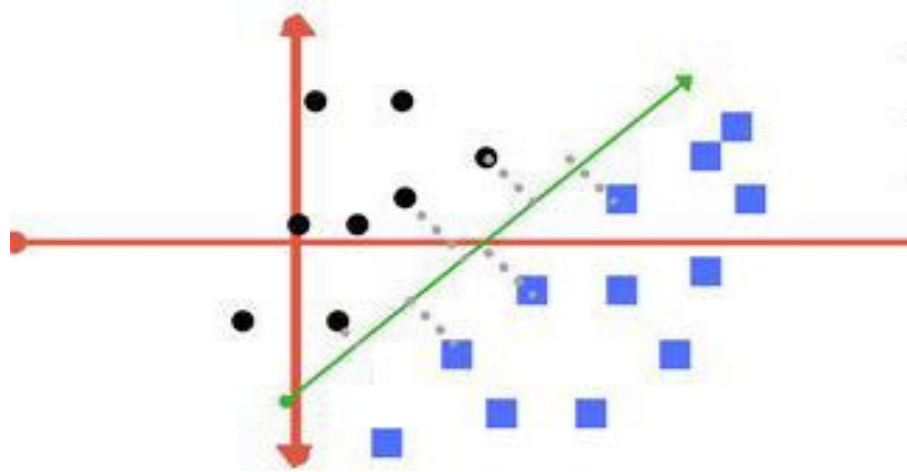
Large value for parameter C
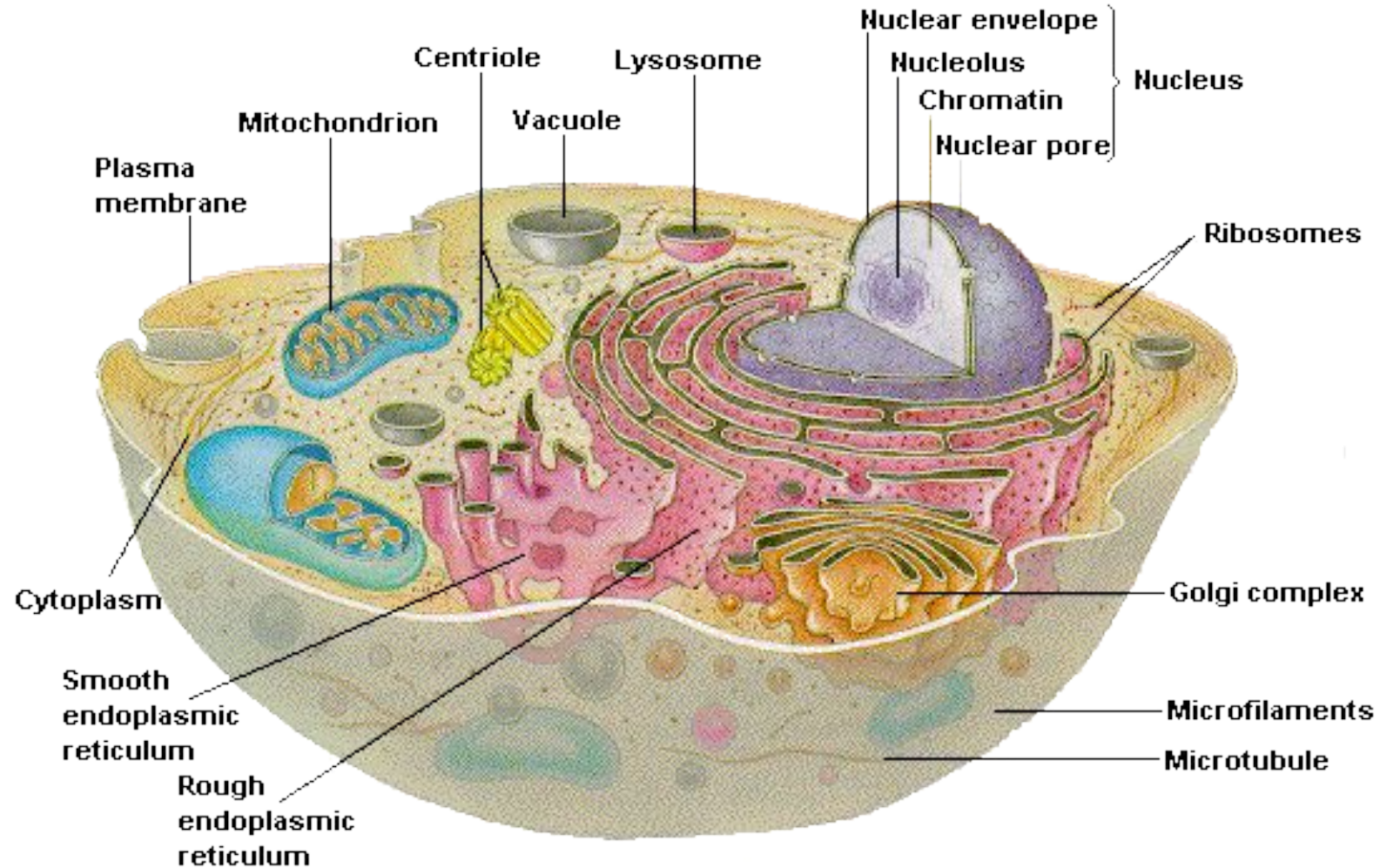
Small value for parameter C

# Gamma

- Defines how far influences the calculation of plausible line of separation.

- Low gamma    -----> points far from plausible line are considered for calculation

- High gamma    -----> points close to plausible line are considered for calculation

# A practical example, protein localization

- Proteins are synthesized in the cytosol.

- Transported into different subcellular locations where they carry out their functions.

- Aim: To predict in what location a certain protein will end up!!!

# Subcellular Locations

# Method

- Hypothesis: The amino acid composition of proteins from different compartments should differ.

- Extract proteins with know subcellular location from SWISSPROT.

- Calculate the amino acid composition of the proteins.

- Try to differentiate between: cytosol, extracellular, mitochondria and nuclear by using SVM

# Input encoding

Prediction of nuclear proteins:

Label the known nuclear proteins as +1 and all others as –1.

The input vector xi represents the amino acid composition.

Eg xi =(4.2,6.7,12,….,0.5)

     A , C , D,….., Y)

Nuclear → SVM → Model

All others →

# Caution on overfitting



Image classification of tanks. Autofire when an enemy tank is spotted.

Input data: Photos of own and enemy tanks.

Worked really good with the training set used.

In reality it failed completely.

Reason: All enemy tank photos taken in the morning. All own tanks in dawn.

The classifier could recognize dusk from dawn!!!!