

Week 10

Course. Introduction to Machine Learning **Theory 10. Support Vector Machine**

Dr. Maria Salamó Llorente
maria.salamo@ub.edu

Dept. Mathematics and Informatics,
Faculty of Mathematics and Informatics,
University of Barcelona (UB)

Introduction to Machine Learning

Unsupervised Learning

Cluster Analysis

Factor Analysis

Visualization

K-Means,
Fuzzy C-means,
EM

PCA, ICA

Self Organized Maps (SOM) ,
Multi-Dimensional Scaling

Lazy Learning
(K-NN, IBL, CBR)

Overfitting,
model selection and feature selection

Kernel Learning

Ensemble Learning
(Trees, Adaboost)

Perceptron,
SVM

Supervised Learning

Non Linear Decision

Linear Decision

Decision Learning Theory

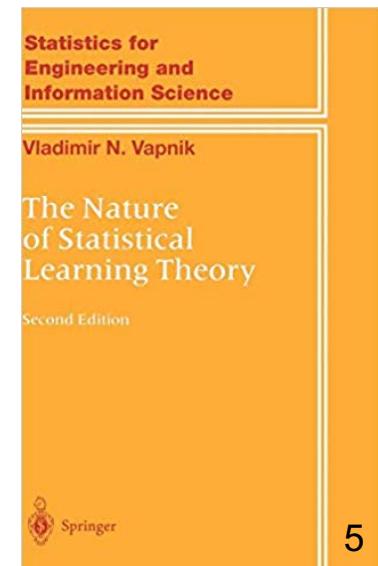
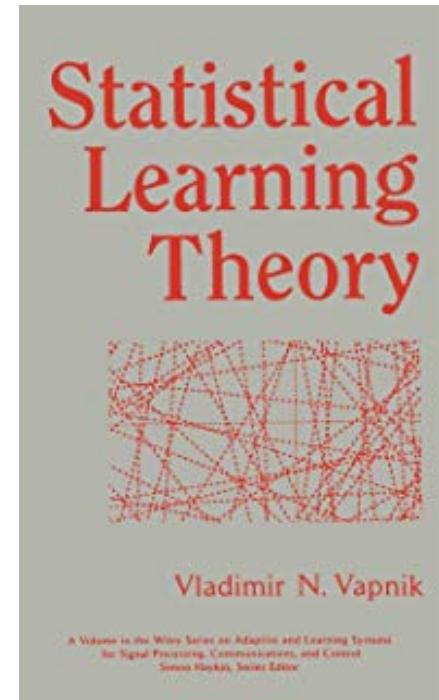
Basic concepts of Decision Learning Theory

Contents

1. Introduction to Support Vector Machine (SVM)
2. Defining SVM conceptually
3. Linear SVM mathematically
4. Non-Linear SVM mathematically

Introduction to Support Vector Machines

- Theoretically well motivated algorithm: developed from **Statistical Learning Theory** (Vapnik & Chervonenkis) since the 60s.
- SVMs were introduced in **COLT-92** by Boser, Guyon & Vapnik. Become rather popular since
- **Empirically good performance:** successful applications in many fields (bioinformatics, text, image recognition, pattern analysis, ...)



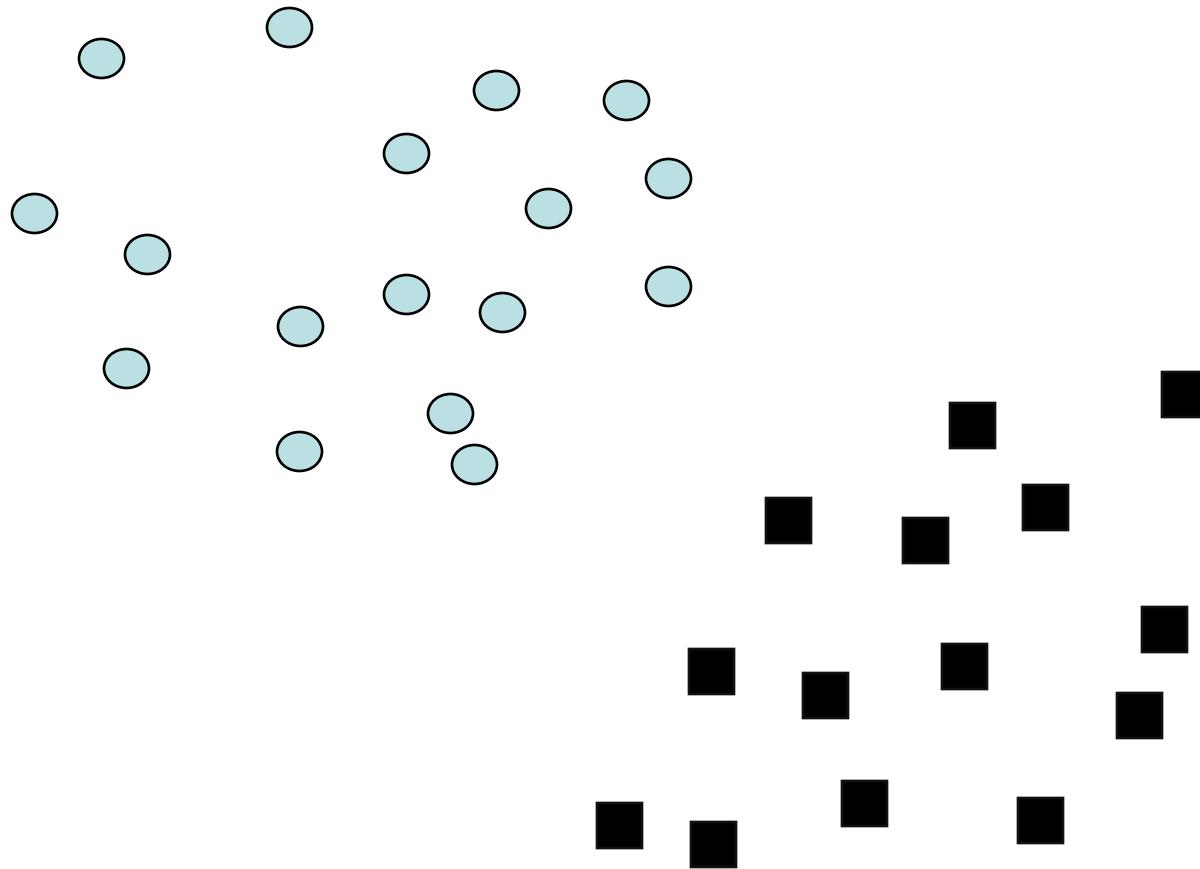
Introduction to SVM

- SVM is a **supervised machine learning algorithm** which can be used for **classification or regression**
- The **objective** of the Support Vector Machine is to find **the best splitting boundary between data**
- They belong to the category of **linear classifiers**
 - It plots each data item as a point in n-dimensional space (where n is number of features you have) with the value of each feature being the value of a particular coordinate
 - Classification is done by **finding the separating hyperplane** that differentiate the two classes very well

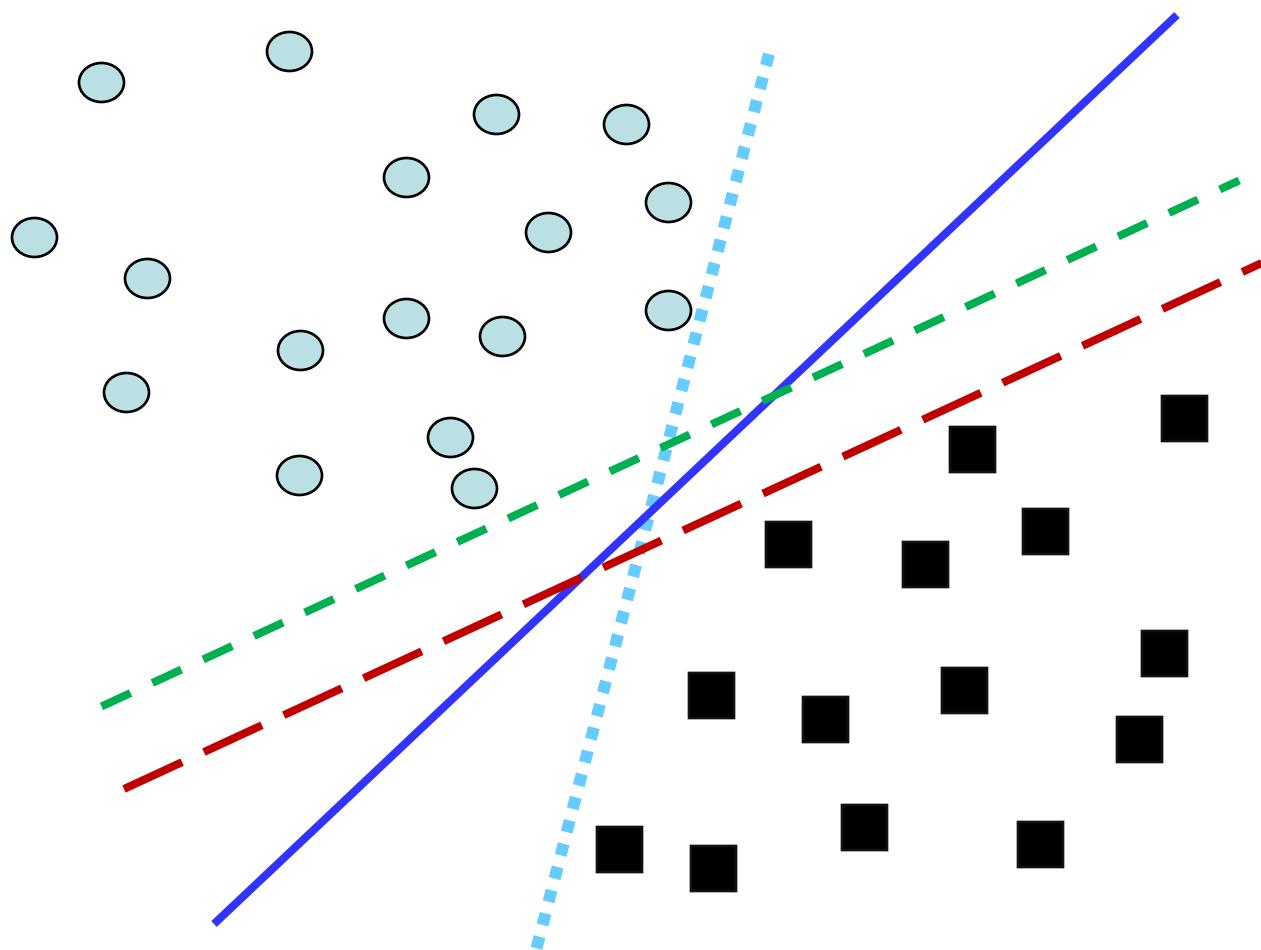
Hyperplane

- A **hyperplane** is a generalization of a plane
 - In 2 dimensions: can separate by a line
 - In 3 dimensions: can separate by a plane
 - In higher dimensions: need *hyperplanes*
- A **hyperplane** is an affine subspace of dimension $n-1$, which divides the space into two half spaces which correspond to the inputs of 2 distinct classes

Identifying the right hyperplane

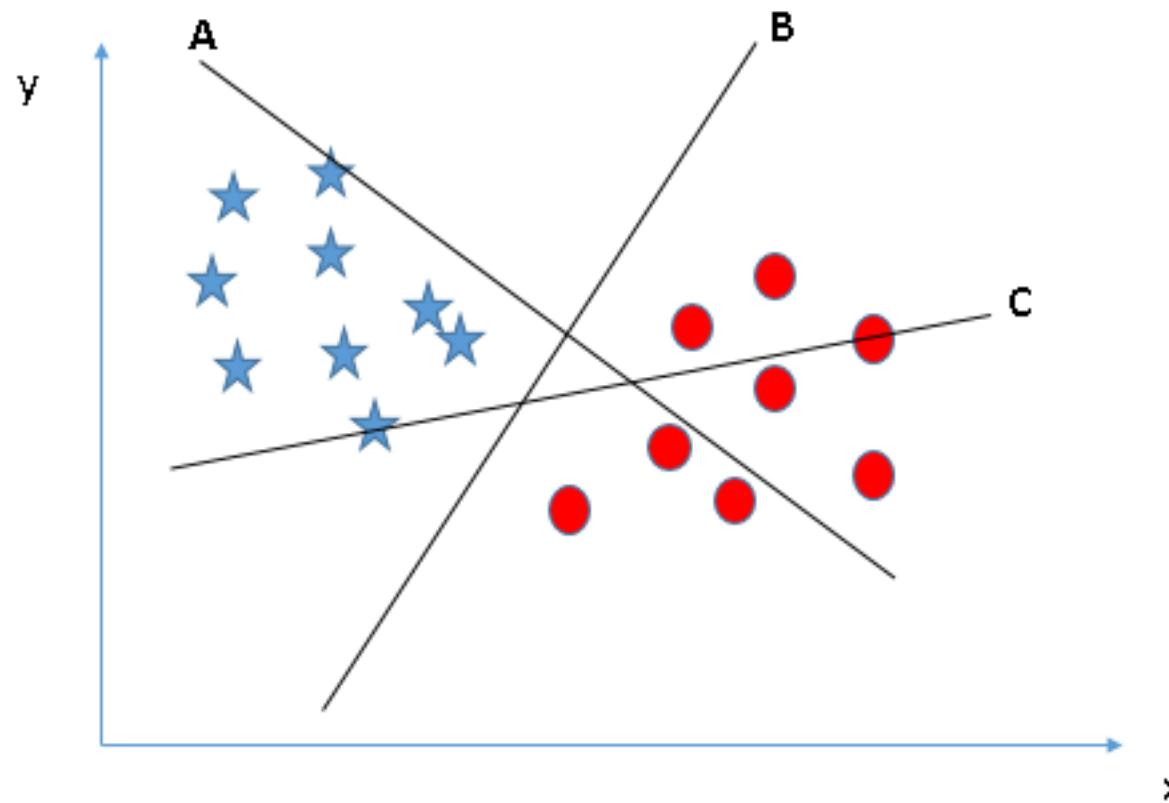


Identifying the right hyperplane



Identifying the right hyperplane

Scenario 1

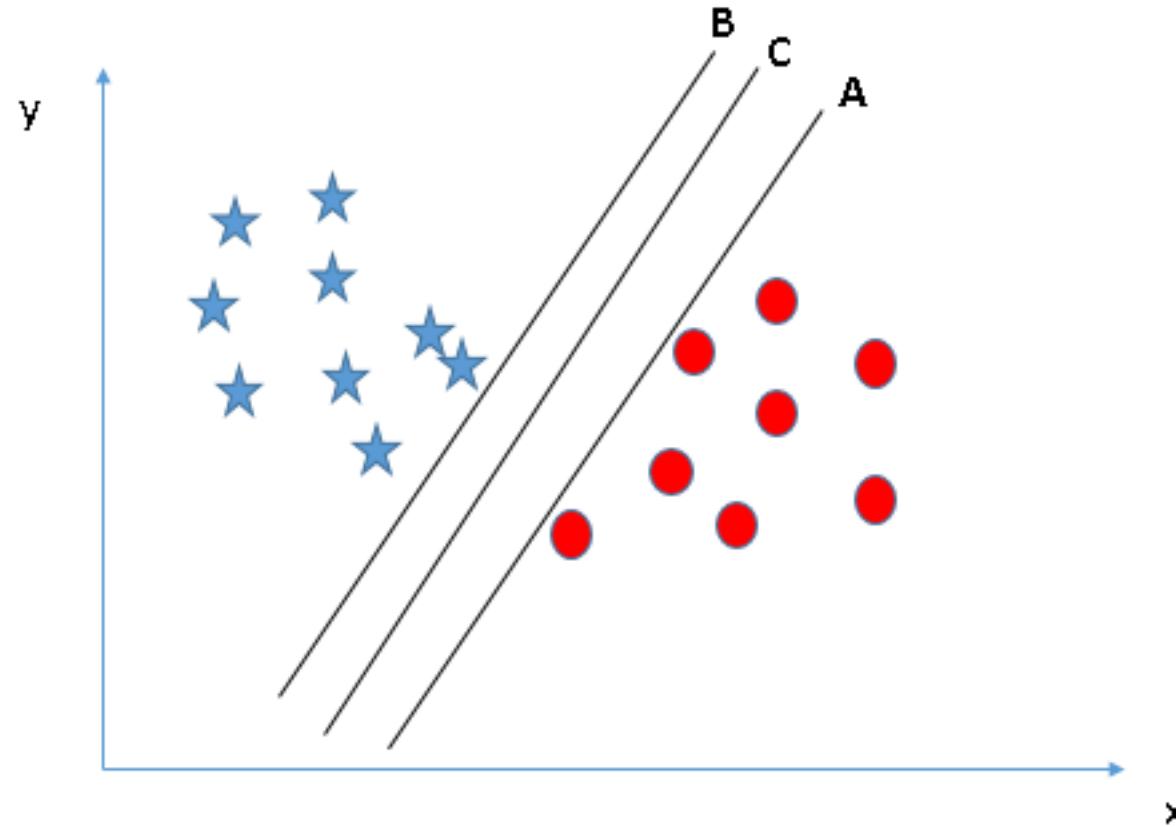


Select the hyperplane that separates/segregates the two classes better?



Identifying the right hyperplane

Scenario 2

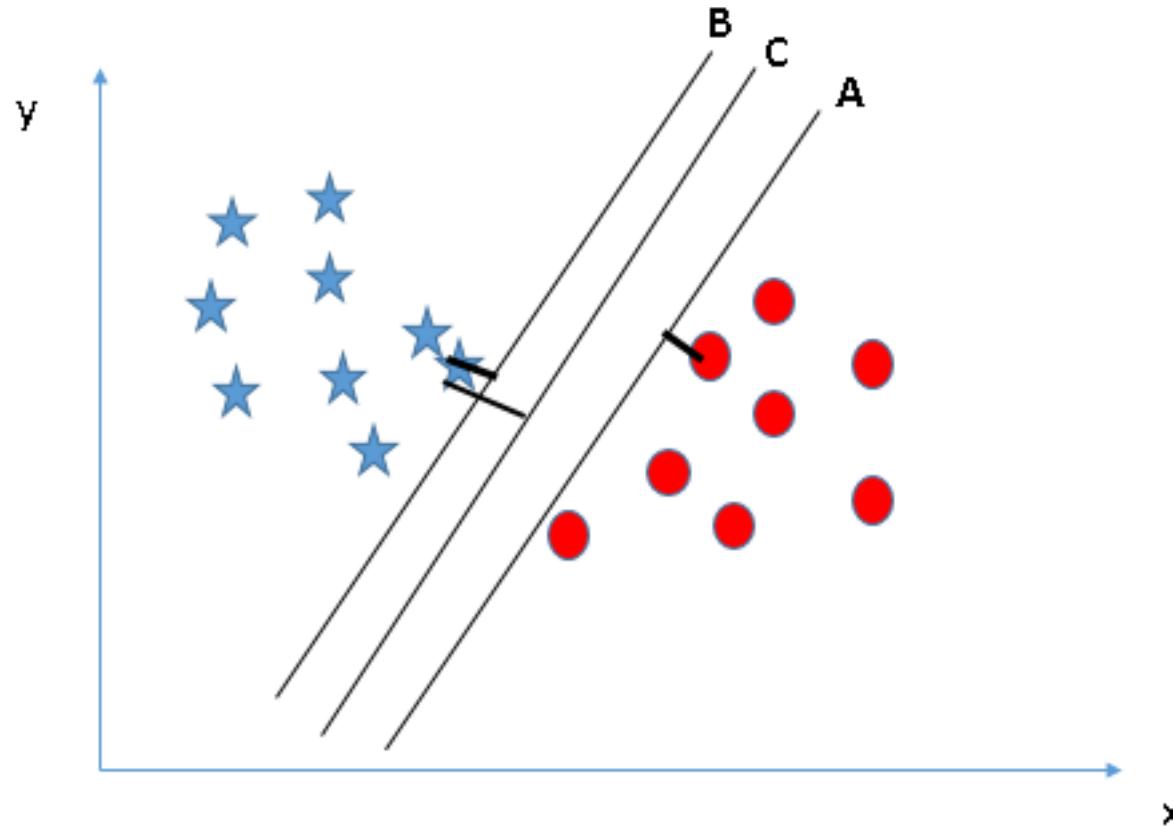


How can we identify the right hyperplane?

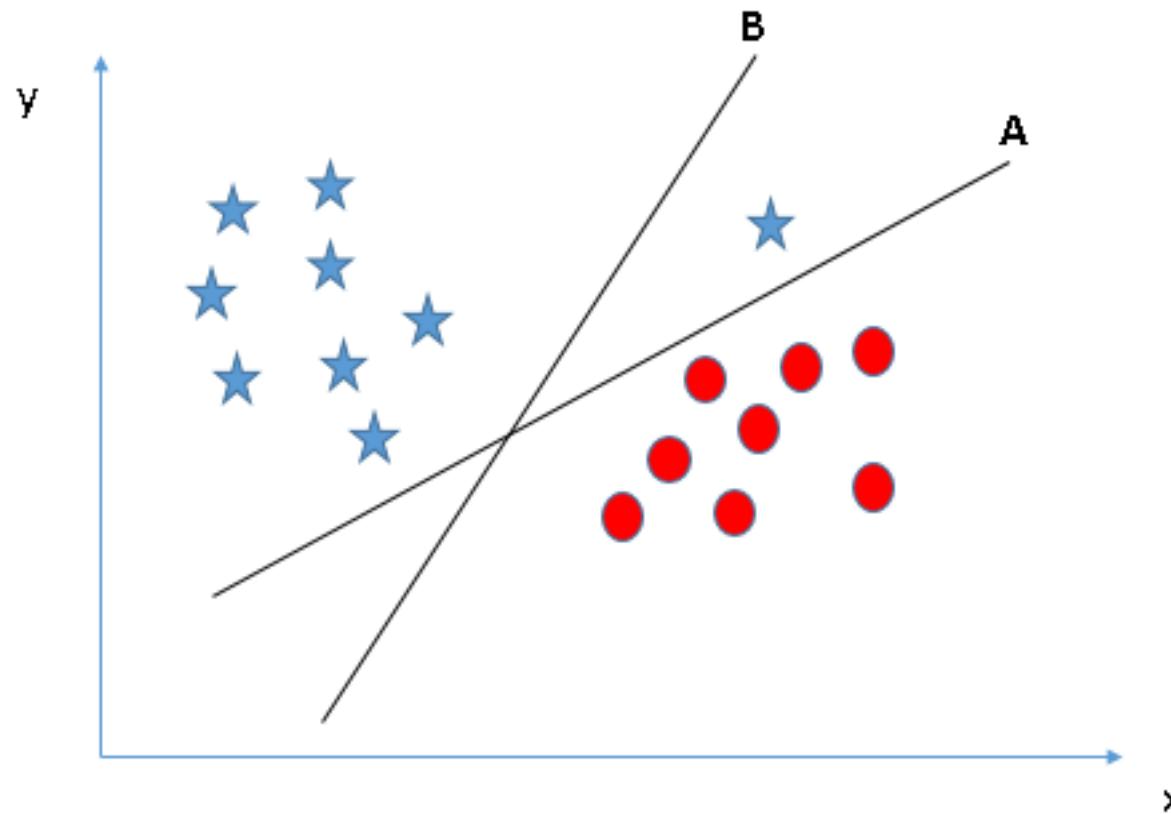


Identifying the right hyperplane

Scenario 2



Scenario 3

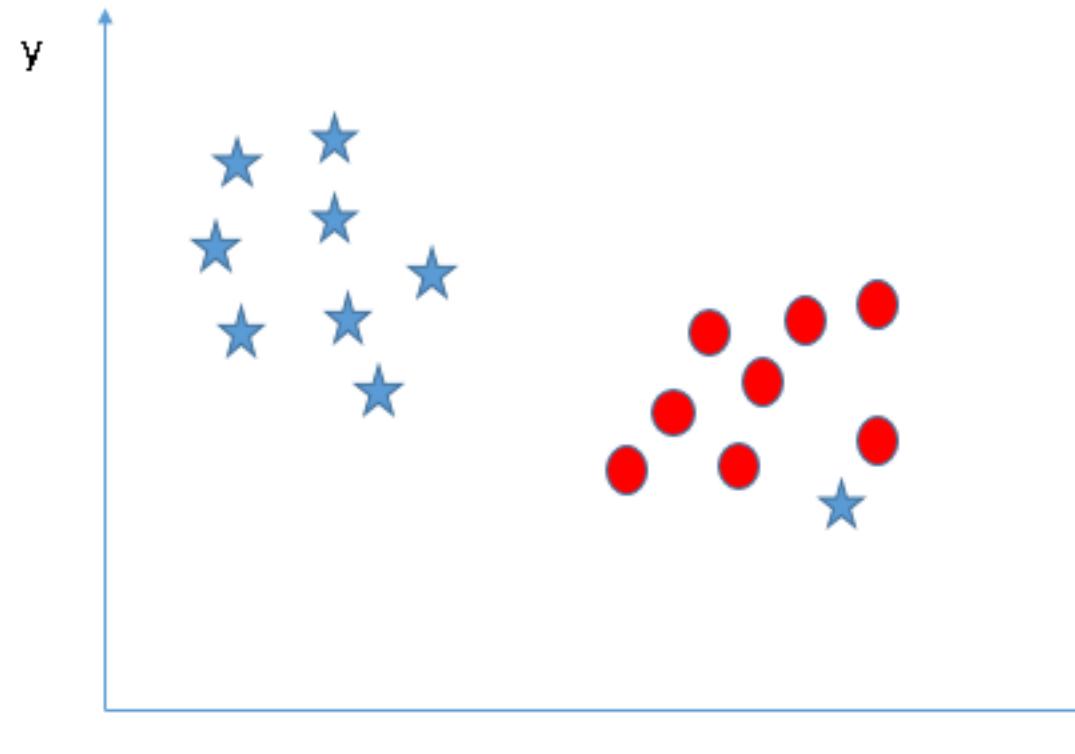


Use the rules discussed to select the right hyperplane



Identifying the right hyperplane

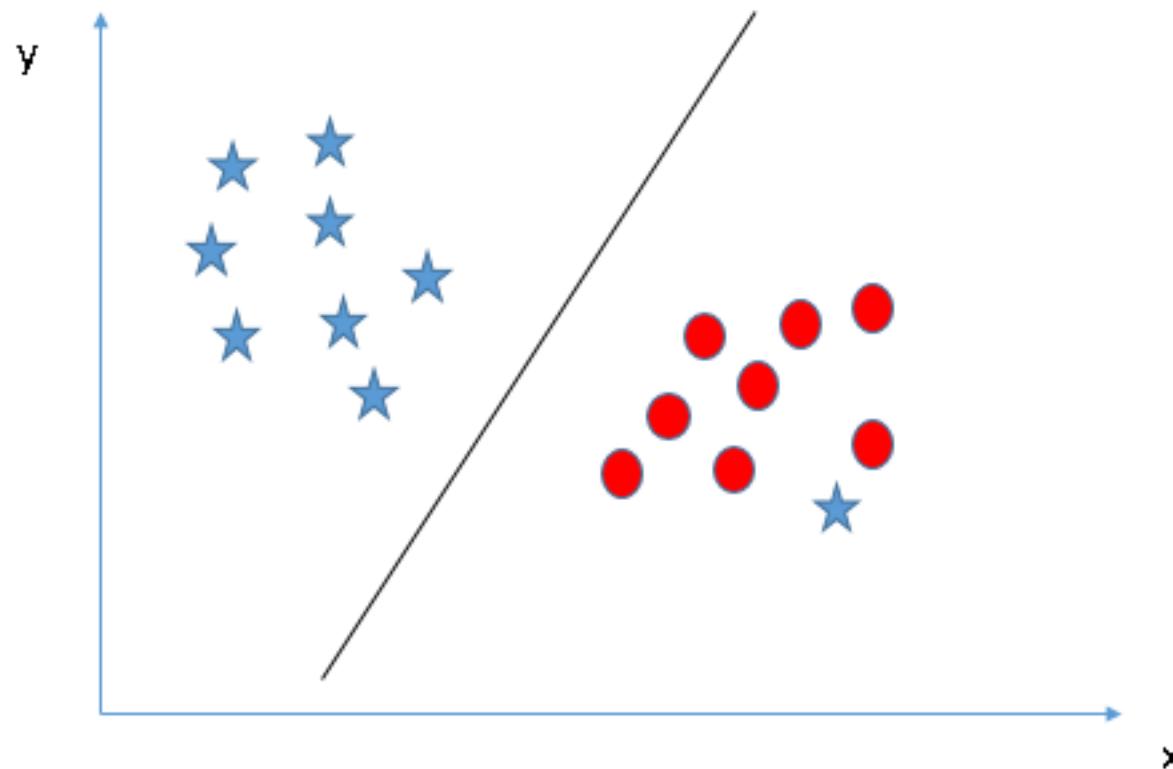
Scenario 4



Select the hyperplane that separates/segregates the two classes better?

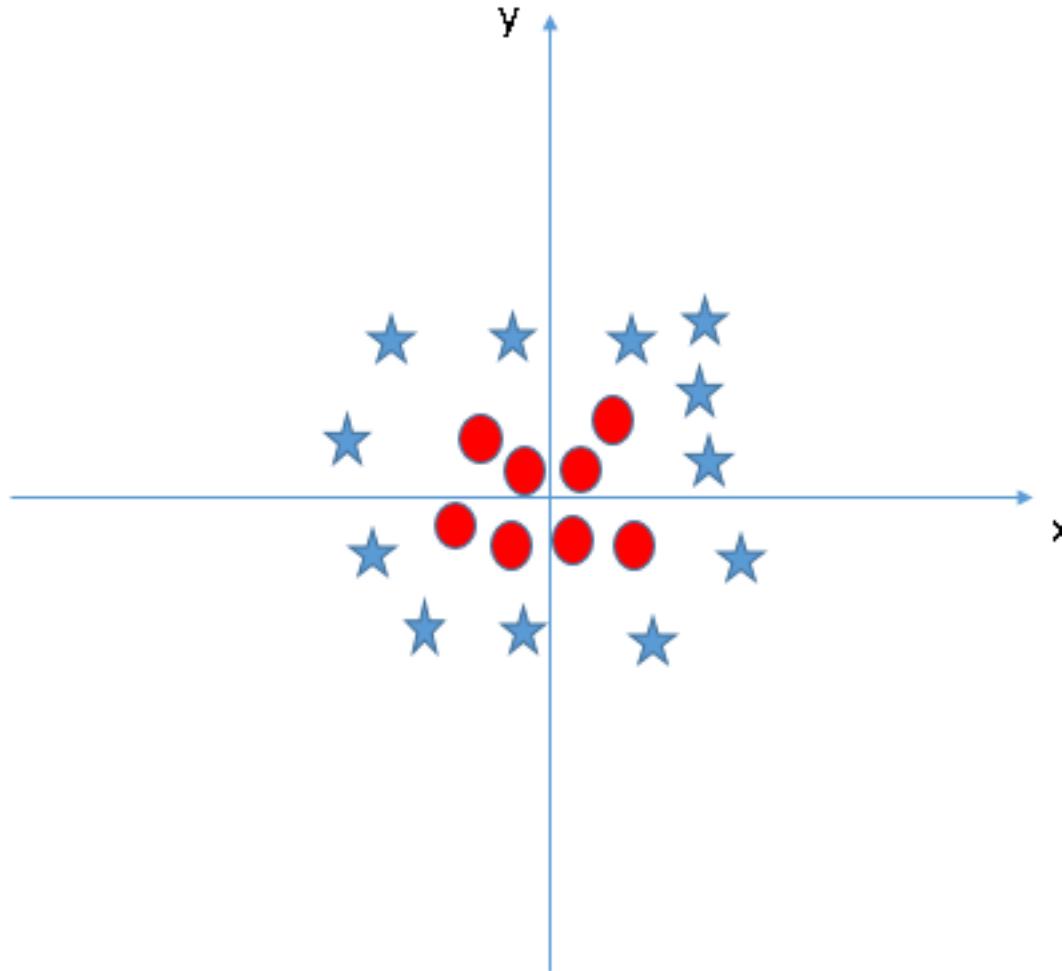


Scenario 4



Identifying the right hyperplane

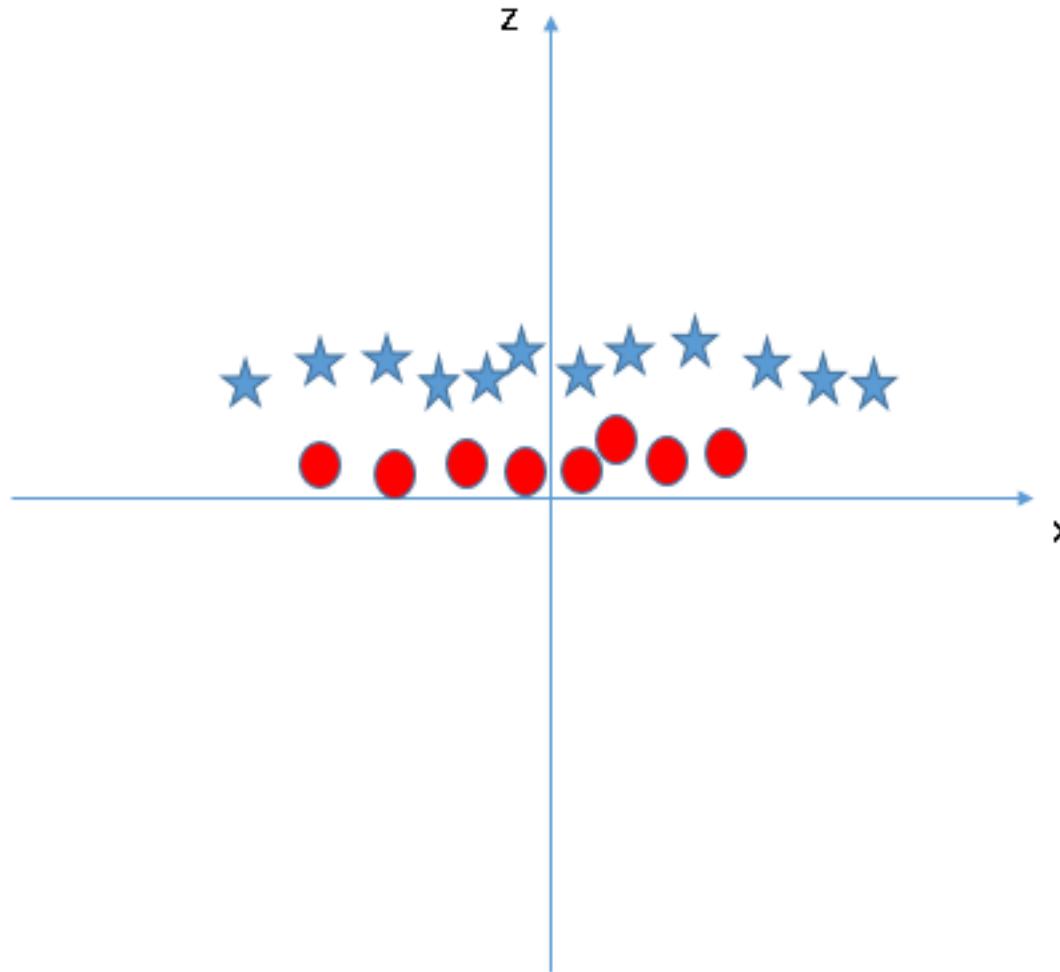
Scenario 5



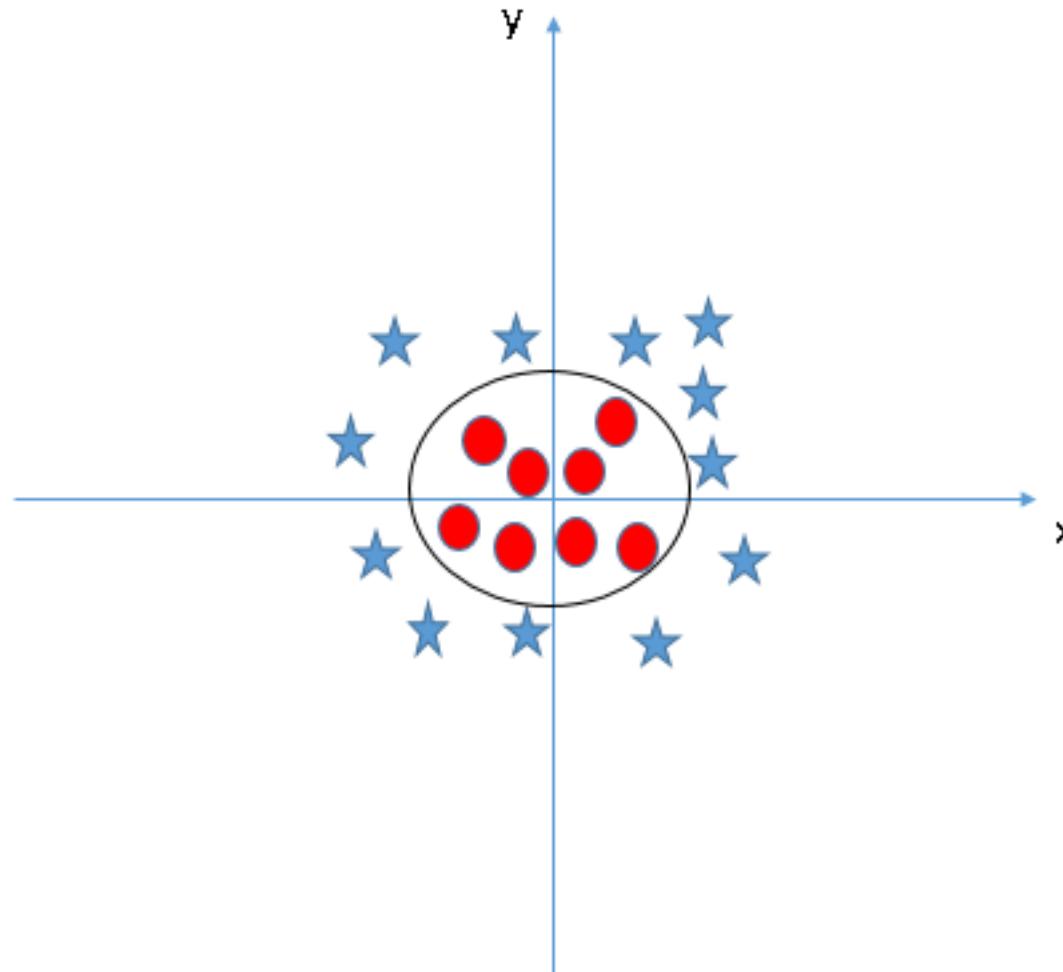
Select the hyperplane that separates/segregates the two classes better?



Scenario 5



Scenario 5

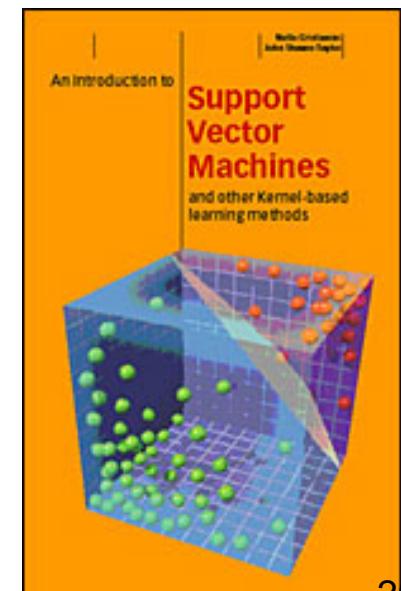


Defining SVM conceptually

Definition

Formal Definition: ‘Support Vector Machine is a system for efficiently training linear learning machines in kernel-induced feature spaces, while respecting the insights of generalization theory and exploiting optimization theory.’

- **AN INTRODUCTION TO SUPPORT VECTOR MACHINES (and other kernel-based learning methods)**
N. Cristianini and J. Shawe-Taylor
Cambridge University Press
2000 ISBN: 0 521 78019 5
- **Kernel Methods for Pattern Analysis**
John Shawe-Taylor & Nello Cristianini
Cambridge University Press, 2004



Concepts: 2D Case

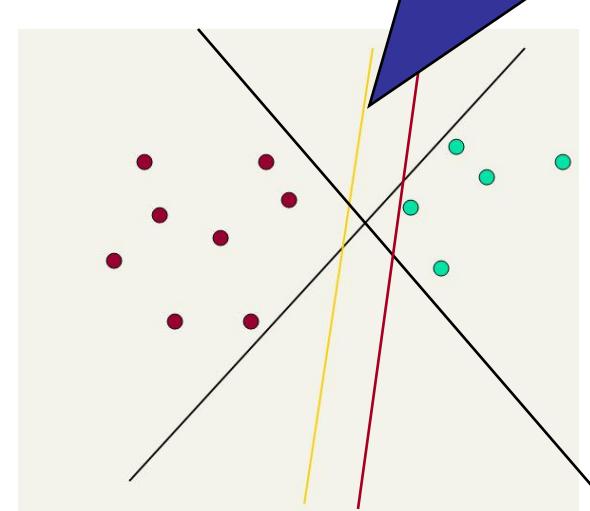
The aim of the SVM classification is to devise a computationally efficient way of learning “good” separating hyperplanes in a high dimensional feature space

Find a, b, c such that

$$ax + by \geq c \quad \text{for red points}$$

$$ax + by \leq c \quad \text{(or } < c\text{) for green points}$$

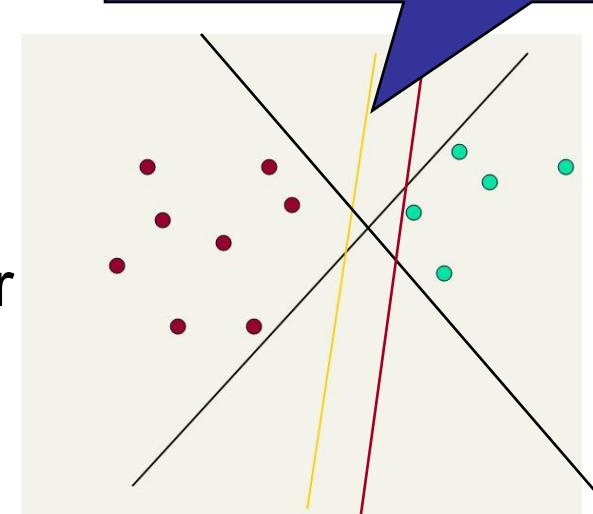
*This line represents the decision boundary:
 $ax + by - c = 0$*



Concepts of SVM

- Lots of possible solutions for a, b, c .
- Some methods find a separating hyperplane, but not the optimal one [according to some criterion of expected goodness]
 - E.g., perceptron
- Support Vector Machine (SVM) finds an optimal solution.
 - Maximizes the distance between the hyperplane and the “difficult points” close to decision boundary
 - **One intuition:** if there are no points near the decision surface, then there are no very uncertain classification decisions

*This line represents the decision boundary:
 $ax + by - c = 0$*



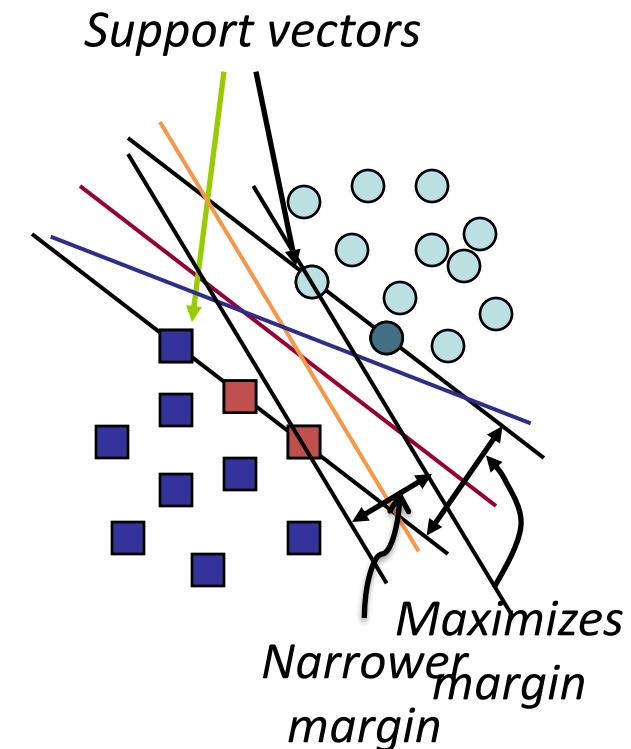
- By “*good hyperplanes*” we will understand the ones optimizing the generalization bounds
 - Different generalization bounds exist
 - Optimize the maximal margin
 - Margin distribution
 - Number of support vectors, etc...
 - The most common and well established is to minimizing the norm of the weight vector
- The simplest model of SVM is a **maximal margin classifier**

Maximal Margin Classifier

- It works only for data that is linearly separable
 - Hence can not be used in many real-world problems
 - However, it is the most simple to understand the building blocks of SVM
- The maximal margin classifier optimizes the generalization error by separating the data with the maximal margin hyperplane,
 - Given that the generalization does not depend on the dimensionality of the space, this separation can be sought in any kernel-induced feature space

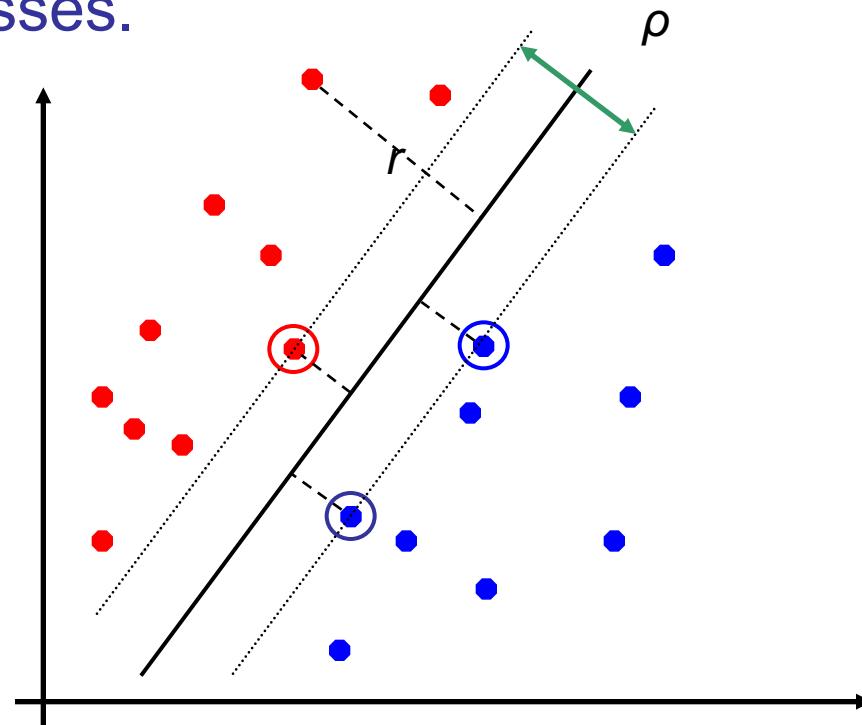
Margin and Support Vectors

- SVMs maximize the **margin** around the separating hyperplane.
 - A.k.a. large margin classifiers
- The decision function is fully specified by a subset of training samples, *the support vectors*.
- Solving SVMs is a **quadratic programming** problem
- Seen by many as the most successful current classification method



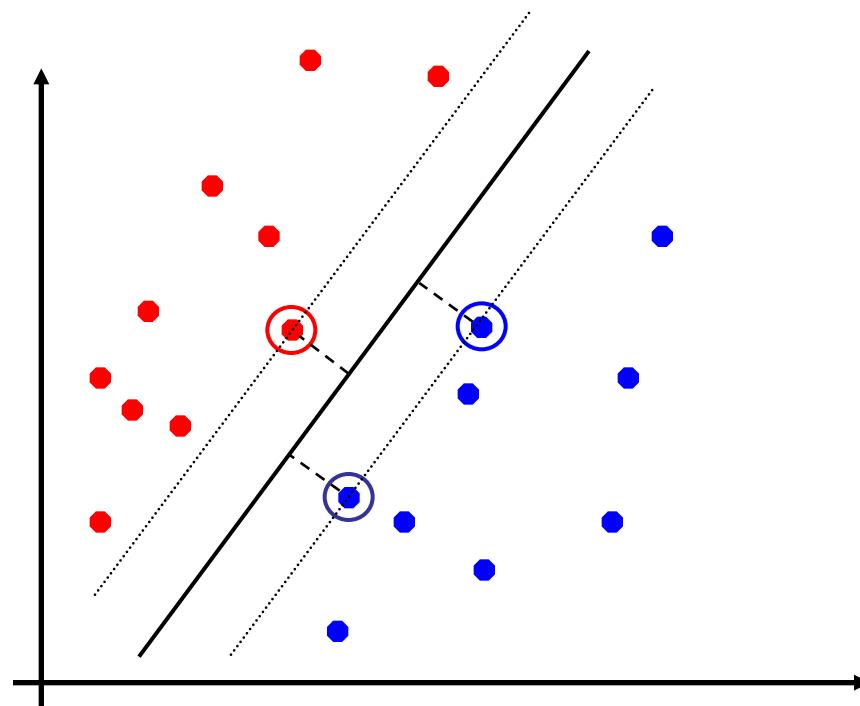
Classification Margin

- Distance from example data, x , to the separator is $r = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- Data closest to the hyperplane are ***support vectors***.
- ***Margin*** ρ of the separator is the width of separation between classes.



Maximum Margin Classification

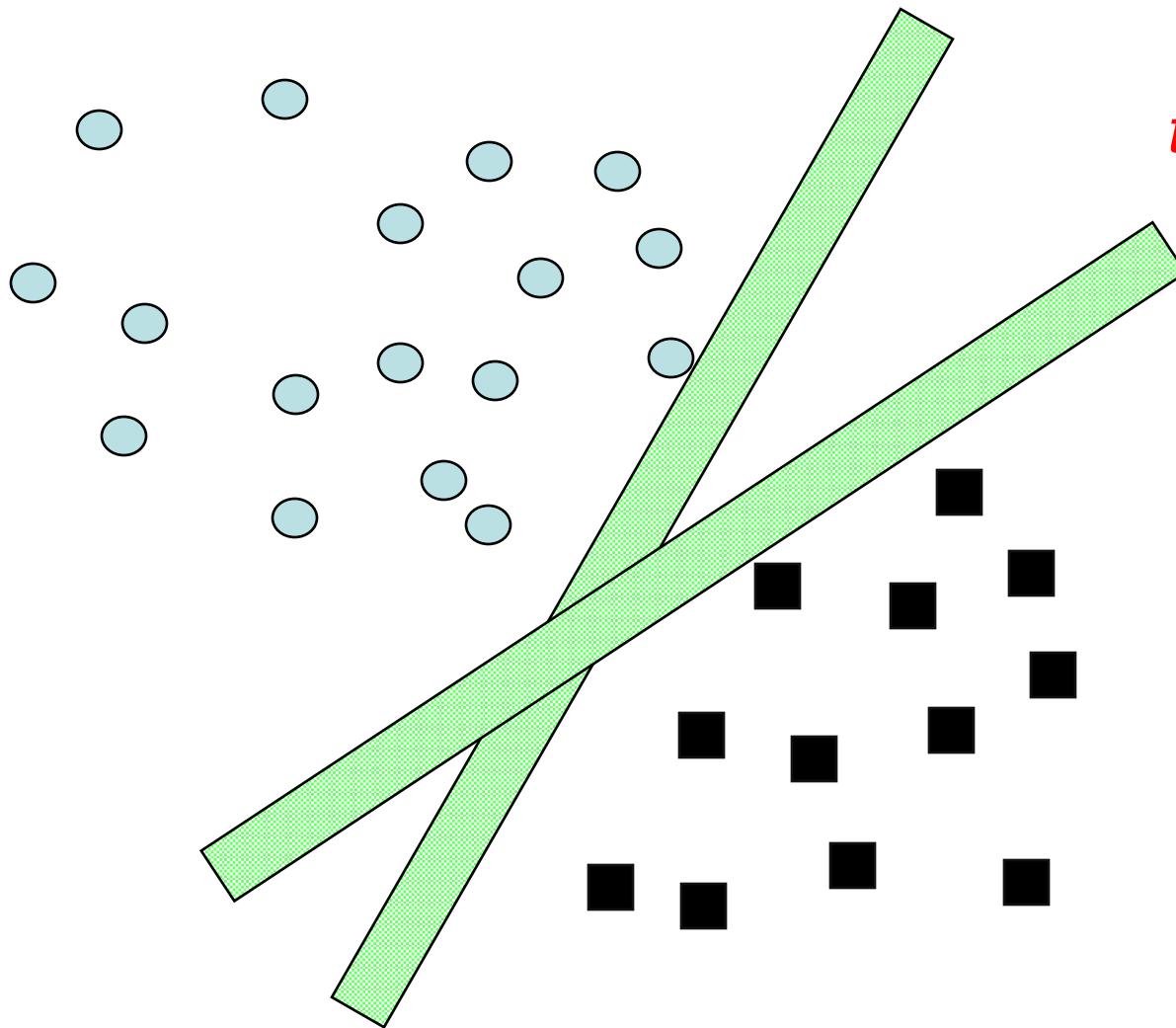
- Maximizing the margin is good according to intuition and theory.
- Implies that only support vectors are important; other training examples are ignorable.



Linear SVM mathematically

- **Generalization theory** gives guidance to prevent overfitting
 - Misclassification error and the function complexity bound **generalization error**.
 - Maximizing margins minimizes complexity.
 - “Eliminates” **overfitting**.
- **Optimization theory** provides mathematical techniques to find hyperplanes
 - Solution depends only on *Support Vectors* not on the number of attributes.

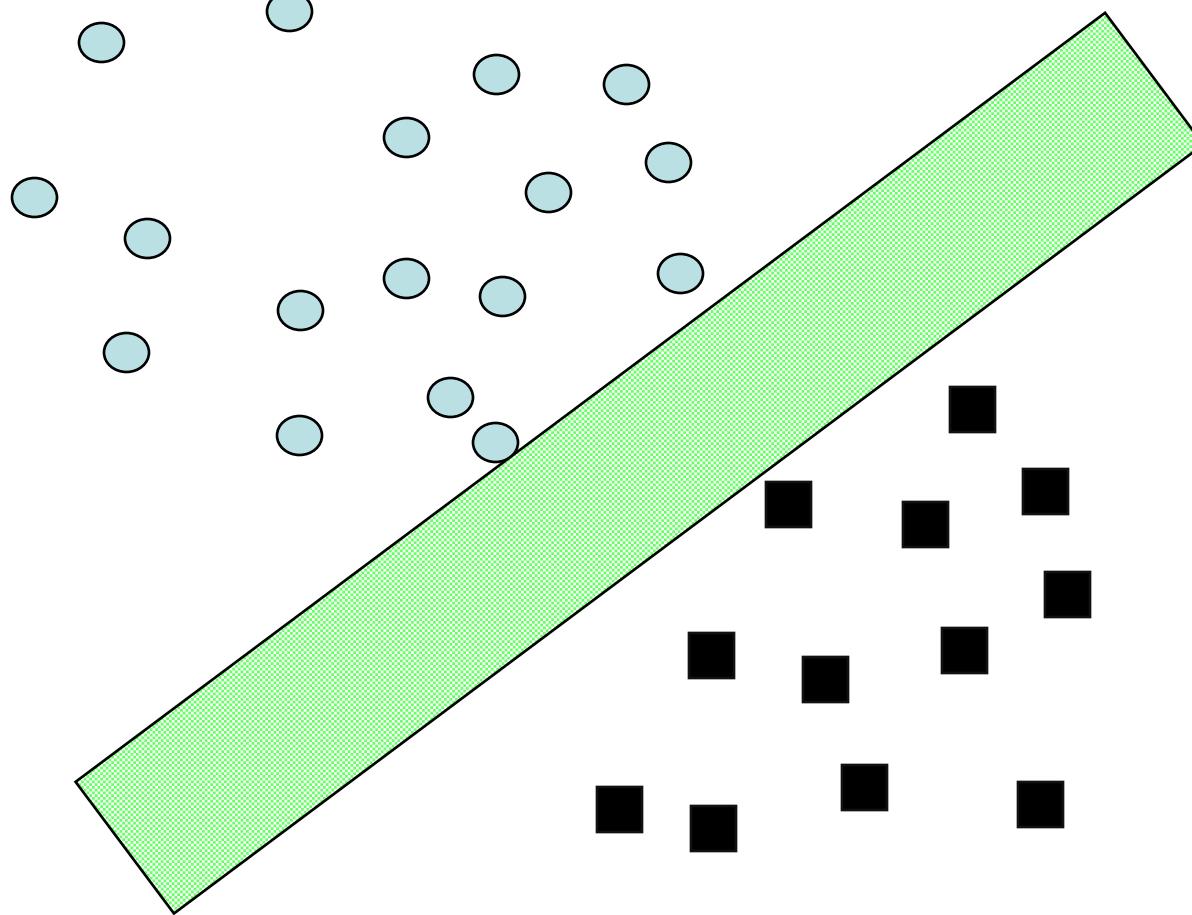
Margins and Complexity



*Skinny margin
is more flexible
thus more complex*

Margins and Complexity

*Fat margin
is less complex*



- The Vapnik-Chervonenkis dimension
 - Model complexity determines the PERFORMANCE/COST on both the TRAINING and TEST sets

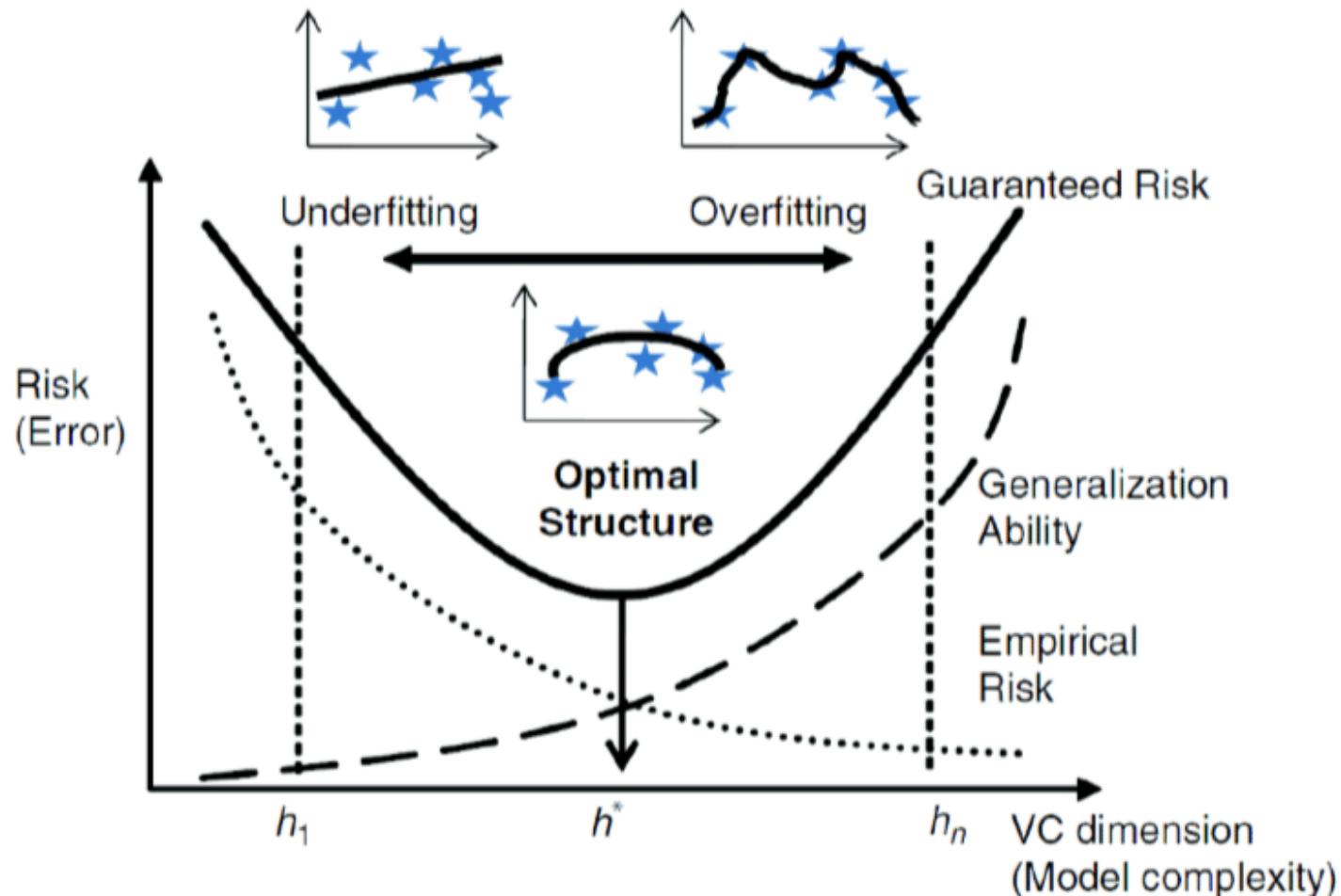
$$P\left(\text{test}_{\text{error}} \leq \text{training}_{\text{error}} + \sqrt{\frac{h \left(\log\left(\frac{2N}{h}\right) + 1 \right) - \log(\eta/4)}{N}}\right) = 1 - \eta$$

true risk → empirical risk → failure probability

- h quantifies the complexity of the model (VC dimension)
- N size on training data
- The expression shows that the error calculated on the test set has an upper bound

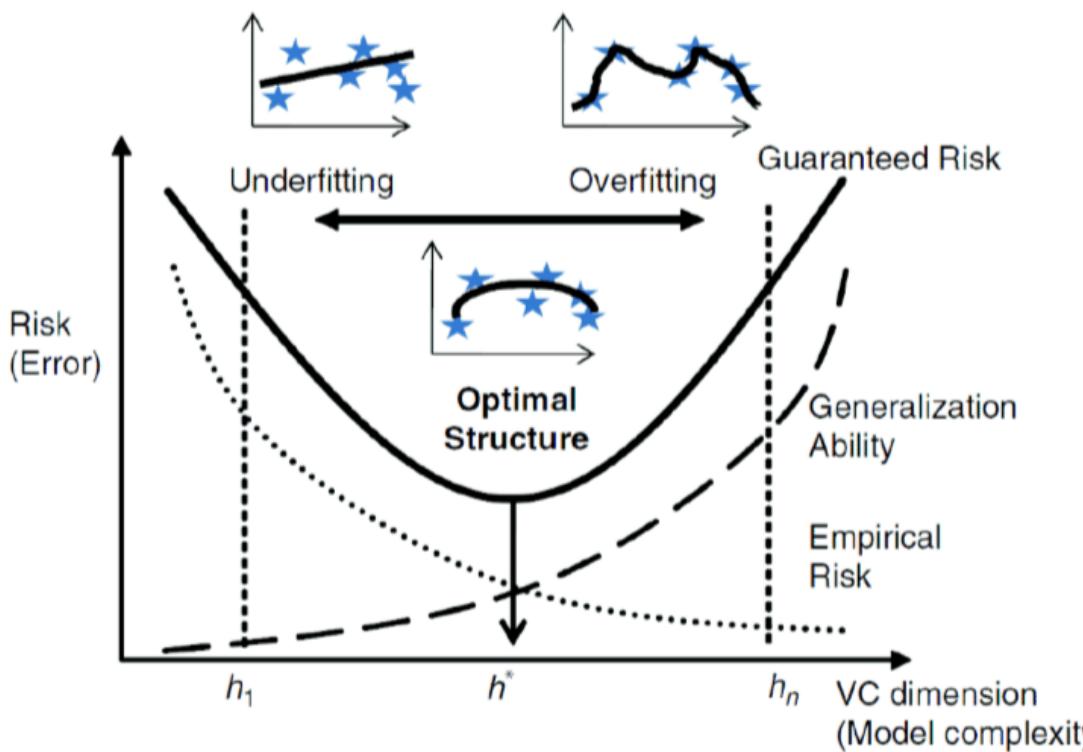
VC dimension

$$P \left(\text{test}_{\text{error}} \leq \text{training}_{\text{error}} + \sqrt{\frac{h \left(\log \left(\frac{2N}{h} \right) + 1 \right) - \log(\eta / 4)}{N}} \right) = 1 - \eta$$



VC dimension

$$P\left(\text{test}_{error} \leq \text{training}_{error} + \sqrt{\frac{h \left(\log\left(\frac{2N}{h}\right) + 1 \right) - \log(\eta/4)}{N}}\right) = 1 - \eta$$



For **LINEAR CLASSIFIERS**

m: dimensions/features

$$h = m + 1$$

Given a **linearly separable** training sample

$$S = ((x_1, y_1), \dots, (x_n, y_n))$$

where $x_i \in \mathbb{R}^n$ and $y_i \in \{+1, -1\}$

the hyperplane (w, b) that solves the optimization problem

$$\begin{aligned} & \text{minimize}_{w,b} && < w \cdot w >, \\ & \text{subject to} && y_i(< w \cdot x_i > + b) \geq 1, \\ & && i = 1, \dots, n \end{aligned}$$

realizes the maximal margin hyperplane with geometric margin $\rho = \frac{1}{\|w\|}$

Hyperplane restrictions

- Assume that all data is at least distance 1 from the hyperplane, then the following two constraints follow for a training set $\{(\mathbf{x}_i, y_i)\}$

$$\begin{aligned} w^T x_i + b &\geq 1 \quad \text{if } y_i = 1 \\ w^T x_i + b &\leq -1 \quad \text{if } y_i = -1 \end{aligned}$$

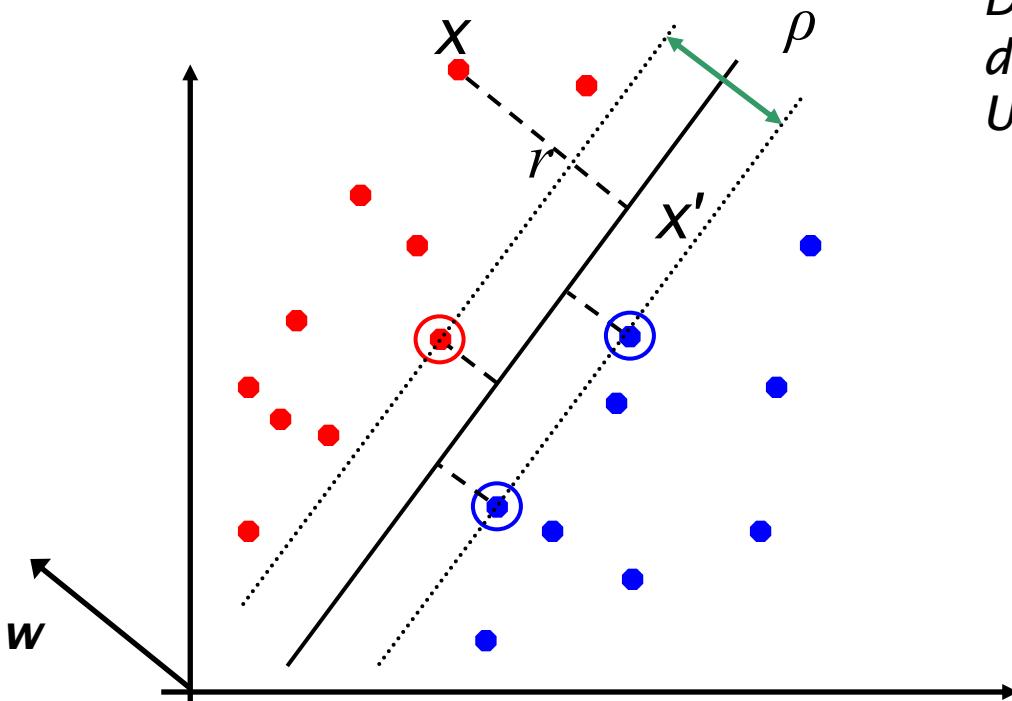
- For **support vectors**, the inequality becomes an equality
- Then, since each example's distance from the hyperplane is $r = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$ and the margin is: $\rho = \frac{1}{\|\mathbf{w}\|}$

Maximum Margin: Formalization

- w : decision hyperplane normal vector
- x_i : data point i
- y_i : class of data point i (+1 or -1) NB: Not 1/0
- Classifier is: $f(x_i) = \text{sign}(w^T x_i + b)$
- Functional margin of x_i is: $y_i (w^T x_i + b)$
 - But note that we can increase this margin simply by scaling w, b
- Functional margin of dataset is twice the minimum functional margin for any point

Geometric Margin

- Distance from example to the separator is $r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- **Margin ρ** of the separator is the width of separation between support vectors of classes.



Derivation of finding r :

Dotted line $x' - x$ is perpendicular to decision boundary so parallel to w . Unit vector is $w/\|w\|$, so line is $rw/\|w\|$.

$$x' = x - yrw/\|w\|.$$

x' satisfies $\mathbf{w}^T x' + b = 0$.

$$\text{So } \mathbf{w}^T(x - yrw/\|w\|) + b = 0$$

Recall that $|w| = \sqrt{w^T w}$.

$$\text{So } \mathbf{w}^T x - yr/\|w\| + b = 0$$

So, solving for r gives:

$$r = y(\mathbf{w}^T x + b)/\|w\|$$

Linear Support Vector Machine

- **Hyperplane**

$$\mathbf{w}^T \mathbf{x} + b = 0$$

- **Extra scale constraint:**

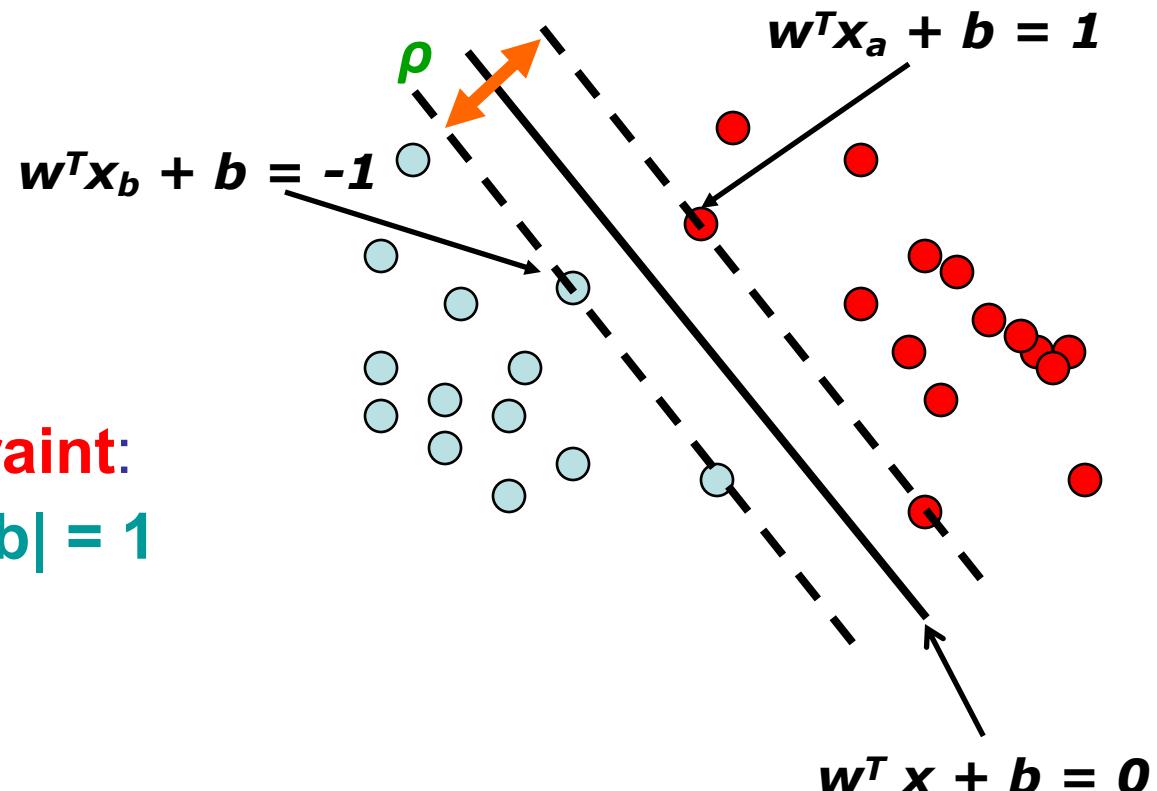
$$\min_{i=1,\dots,n} |\mathbf{w}^T \mathbf{x}_i + b| = 1$$

- This implies:

$$\mathbf{w}^T (\mathbf{x}_a - \mathbf{x}_b) = 2$$

$$\rho = \|\mathbf{x}_a - \mathbf{x}_b\|_2 = 2/\|\mathbf{w}\|_2$$

The margin is $1/\|\mathbf{w}\|_2$



The optimization problem

- Formulation of the problem

Find w and b such that $\rho = \frac{1}{\|w\|}$

is maximized and for all $\{x_i, y_i\}$ satisfying the following constraints

$$w^T x + b \geq +1 \text{ for } y_i = 1 ; w^T x + b \leq -1 \text{ for } y_i = -1$$

- The **optimization theory** establishes that an optimization problem, called **primal**, has a **DUAL** form if the function to be optimized and the restrictions are strictly **convex**
 - Under these circumstances, solving the dual problem let obtaining the solution to the primal problem

The optimization problem

- The next function

$$\min f(w) = \frac{1}{2} \|w\|^2 = \frac{1}{2} \langle w, w \rangle$$

Satisfies the criterion of convexity and thus, it has a DUAL

- **STEPS** to transform our optimization problem:
 1. We build an optimization problem without restrictions using **Lagrange function**
 2. We apply **Karush-Kuhn-Tucker (KKT) conditions**

The optimization problem

- Then we can formulate the ***quadratic optimization problem:***

Find w and b such that $\rho = \frac{1}{\|w\|}$

is maximized and for all $\{x_i, y_i\}$ satisfying the following constraints
 $w^T x + b \geq +1$ **for** $y_i = 1$; $w^T x + b \leq -1$ **for** $y_i = -1$

A better formulation:

Find w and b such that $\Phi(w) = \frac{1}{2} w^T w$ is minimized and

for all $\{x_i, y_i\}$ satisfying the following constraints

$$y_j (w^T x_i + b) \geq 1$$

Solving the Optimization Problem

Find w and b such that $\Phi(w) = \frac{1}{2}w^T w$ is minimized and

for all $\{x_i, y_i\}$ satisfying the following constraints

$$y_j (w^T x_i + b) \geq 1$$

- Need to optimize a **quadratic** function subject to *linear* constraints.
- Quadratic optimization problems are a well-known class of mathematical programming problems, and many (intricate) algorithms exist for solving them (with many special ones built for SVMs).
- The solution involves constructing a **dual problem** where a *Lagrange multiplier* α_i is associated with every constraint in the **primary problem**:

Find $\alpha_1 \dots \alpha_N$ such that

$$Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y_i y_j x_i \cdot x_j \text{ is maximized s.t.}$$

$$(1) \sum_i \alpha_i y_i = 0$$

$$(2) \alpha_i \geq 0 \text{ for all } \alpha_i$$

The Optimization Problem Solution

- The solution has the form:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \mathbf{b} = y_k - \mathbf{w}_k^T \mathbf{x}_k \text{ for any } x_k \text{ such that } \alpha_k \neq 0$$

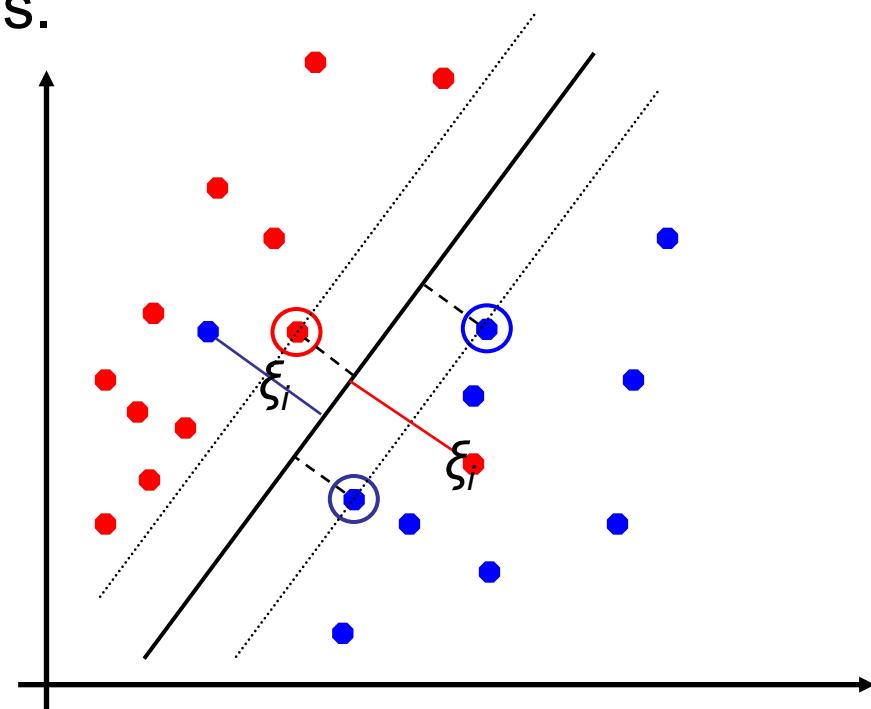
- Each non-zero α_i indicates that corresponding \mathbf{x}_i is a support vector. Then the classifying function will have the form:

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

- Notice that it relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points!

Soft Margin Classification

- What to do if the training set is not linearly separable?
 - *Slack variables* ξ_i can be added to allow misclassification of difficult or noisy examples.
 - Allow some errors
 - Let some points be moved to where they belong to, at a cost
 - Still, try to minimize training set errors, and, to place hyperplane “far” from each class (large margin)



Soft Margin Mathematical definition

Given a linearly separable training sample

$$S = ((x_1, y_1), \dots, (x_n, y_n))$$

where $x_i \in \mathbb{R}^n$ and $y_i \in \{+1, -1\}$

the hyperplane (w, b) that solves the optimization problem

$$\text{minimize}_{w,b} \quad \langle w \cdot w \rangle,$$

$$\begin{aligned} \text{subject to} \quad & y_i(\langle w \cdot x_i \rangle + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, i = 1, \dots, n \end{aligned}$$

realizes the maximal margin hyperplane with geometric margin $\rho = \frac{1}{\|w\|}$

- The old formulation:

Find w and b such that $\Phi(w) = \frac{1}{2}w^T w$ is minimized and

for all $\{x_i, y_i\}$ satisfying the following constraints

$$y_j (w^T x_i + b) \geq 1$$

- The new formulation incorporating **slack variables**:

Find w and b such that $\Phi(w) = \frac{1}{2}w^T w + C \sum \xi_i$ is minimized

and for all $\{x_i, y_i\}$ satisfying the following constraints

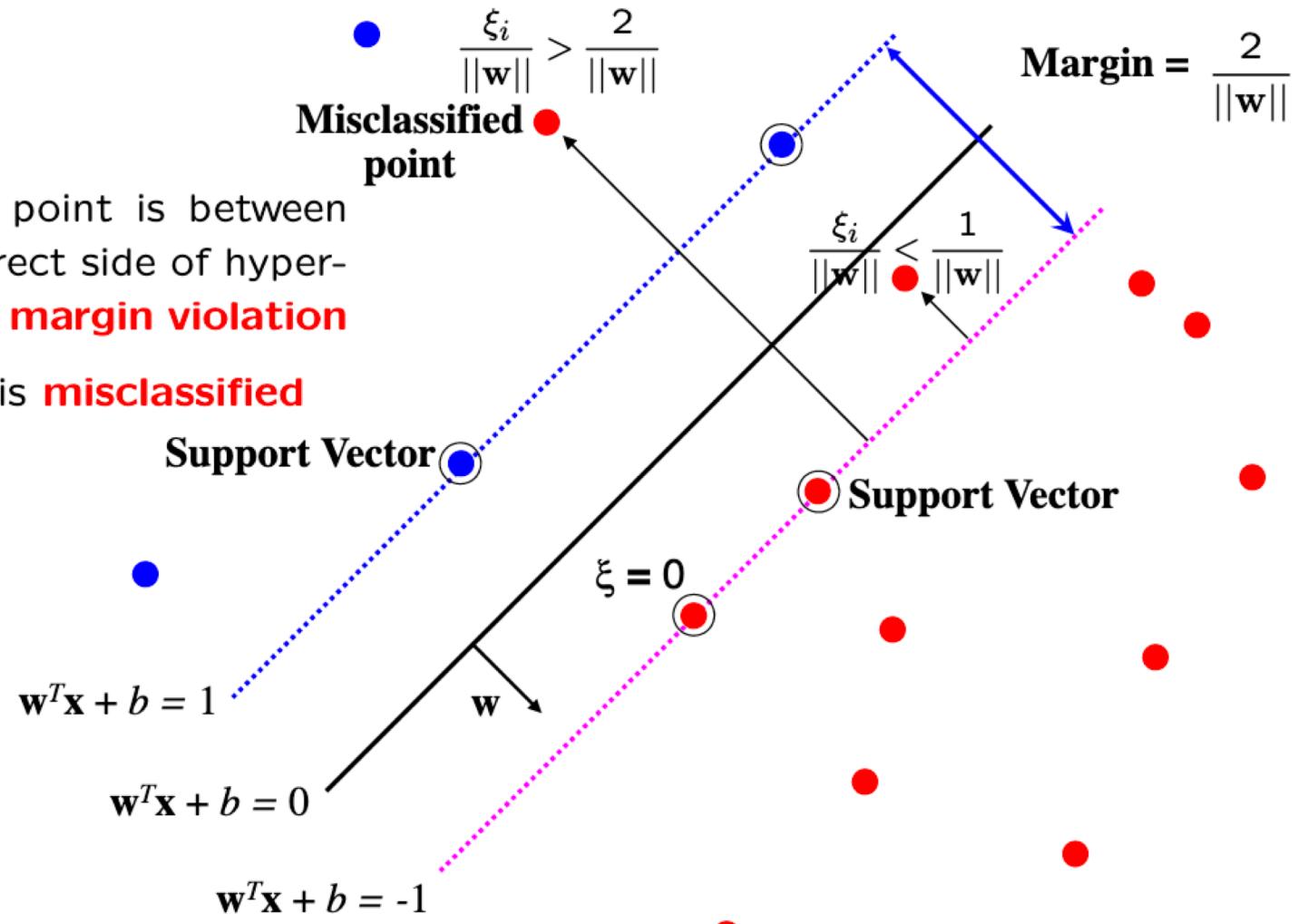
$$y_j (w^T x_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

- Parameter C can be viewed as a way to control overfitting (a regularization term)

"Slack" variables

$$\xi_i \geq 0$$

- for $0 < \xi \leq 1$ point is between margin and correct side of hyperplane. This is a **margin violation**
- for $\xi > 1$ point is **misclassified**



- The dual problem for soft margin classification:

Find $\alpha_1 \dots \alpha_N$ such that

$$Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y_i y_j x_i^T x_j \text{ is maximized s.t.}$$

$$(1) \sum_i \alpha_i y_i = 0$$

$$(2) 0 \leq \alpha_i \leq C \text{ for all } \alpha_i$$

- Neither slack variables ξ_i nor their Lagrange multipliers appear in the dual problem!
 - Again, \mathbf{x}_i with non-zero α_i will be **support vectors**.
 - Solution to the dual problem is:
- w and b are not needed explicitly for classification!*

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

$$\mathbf{b} = y_k(1 - \xi_k) - \mathbf{w}_k^T \mathbf{x}_k \text{ where } k = \arg \max_k \alpha_k$$

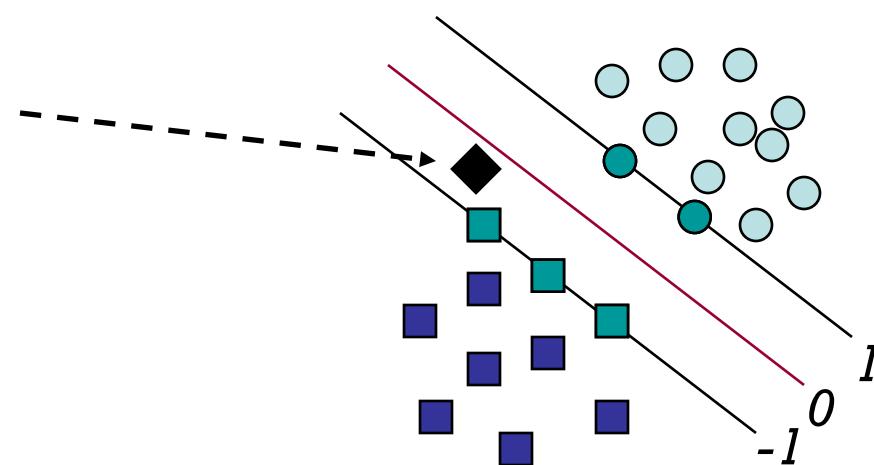
$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

“Soft” margin solution

- C is a **regularization** parameter:
 - Small C allows constraints to be easily ignored → large margin
 - Large C makes constraints hard to ignore → narrow margin
 - $C = \infty$ enforces all constraints: hard margin
- This is still a quadratic optimization problem and there is a unique minimum. Note there is only one parameter, C

- Given a new point \mathbf{x} , we can score its projection onto the hyperplane normal:
 - I.e., compute score: $w^T \mathbf{x} + b = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$
 - Decide class based on whether < 0 or > 0
 - Can set confidence threshold t .

Score $> t$: yes
Score $< -t$: no
Else: don't know



Linear SVMs: 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 \mathbf{x}_i are support vectors with **non-zero Lagrangian multipliers α_i** .
- Both in the dual formulation of the problem and in the solution training points appear only inside inner products:

Find $\alpha_1 \dots \alpha_N$ such that

$$Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y_i y_j x_i^T x_j \text{ is maximized s.t.}$$

$$(1) \sum_i \alpha_i y_i = 0$$

$$(2) 0 \leq \alpha_i \leq C \text{ for all } \alpha_i$$

$$f(x) = \sum \alpha_i y_i x_i^T x + b$$

What else?



Look at the video

<https://youtu.be/efR1C6CvhmE>



TO READ:



A Tutorial on Support Vector Machines for Pattern Recognition
Christopher J.C. Burges

Sections 2 and 3.

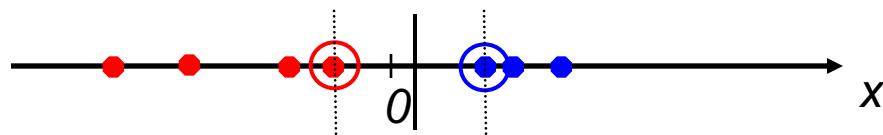
With the exception of Section 4, the remaining of the sections are optional

Non-linear SVM mathematically

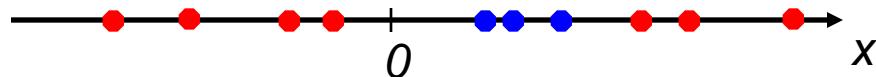
- So far, we have only considered large-margin classifier with a linear decision boundary
- How to generalize it to become nonlinear?
- **Key idea:** transform x_i to a higher dimensional space to “make life easier”
 - **Input space:** the space the point x_i are located
 - **Feature space:** the space of $\phi(x_i)$ after transformation

Non-linear SVMs

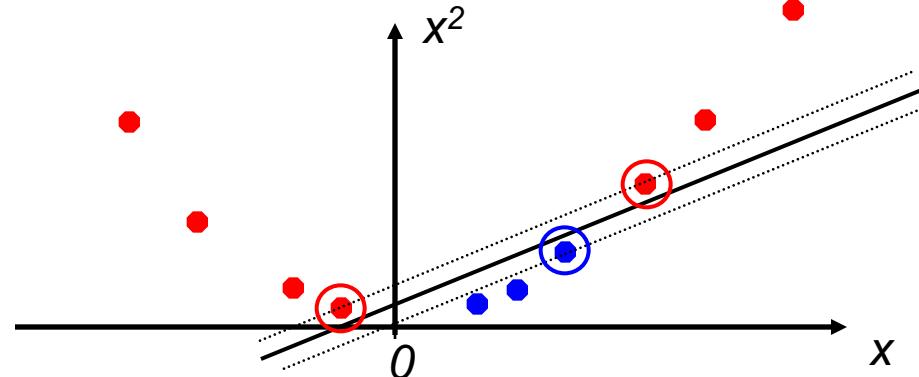
- Datasets that are linearly separable with some noise work out great:



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

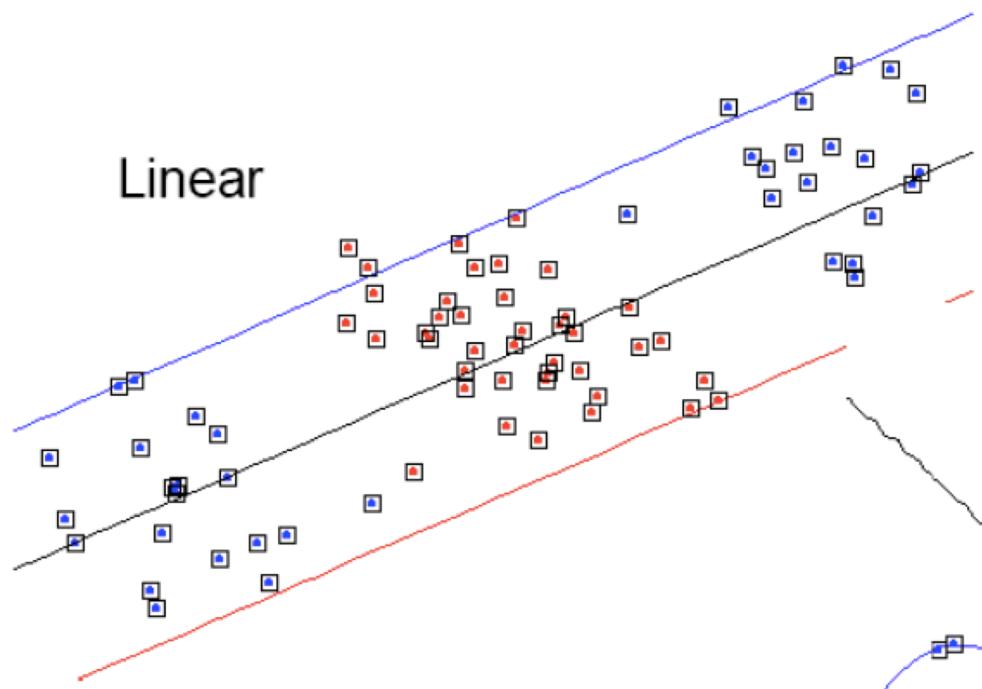


- How about... mapping data to a higher-dimensional space:

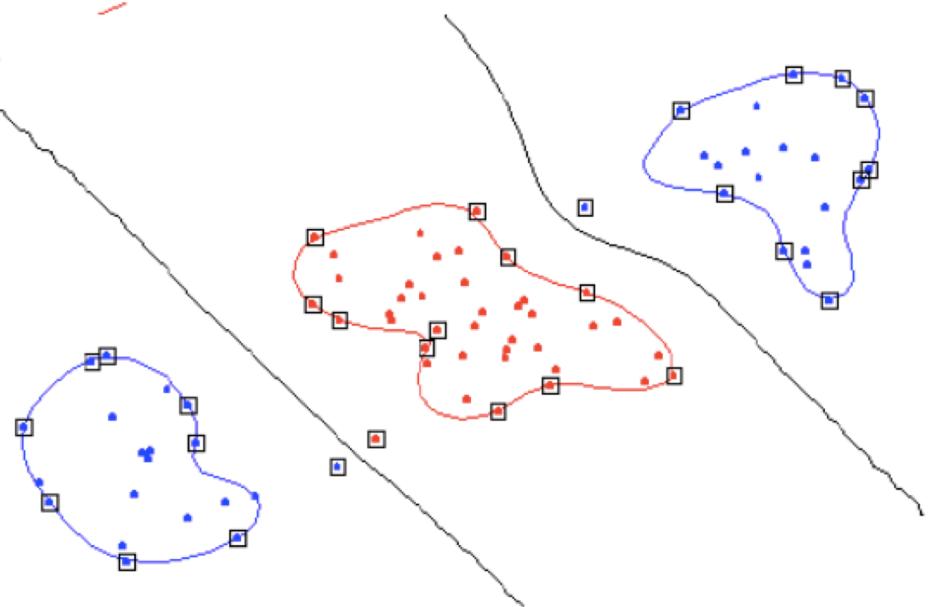


Non-linear SVMs

Linear

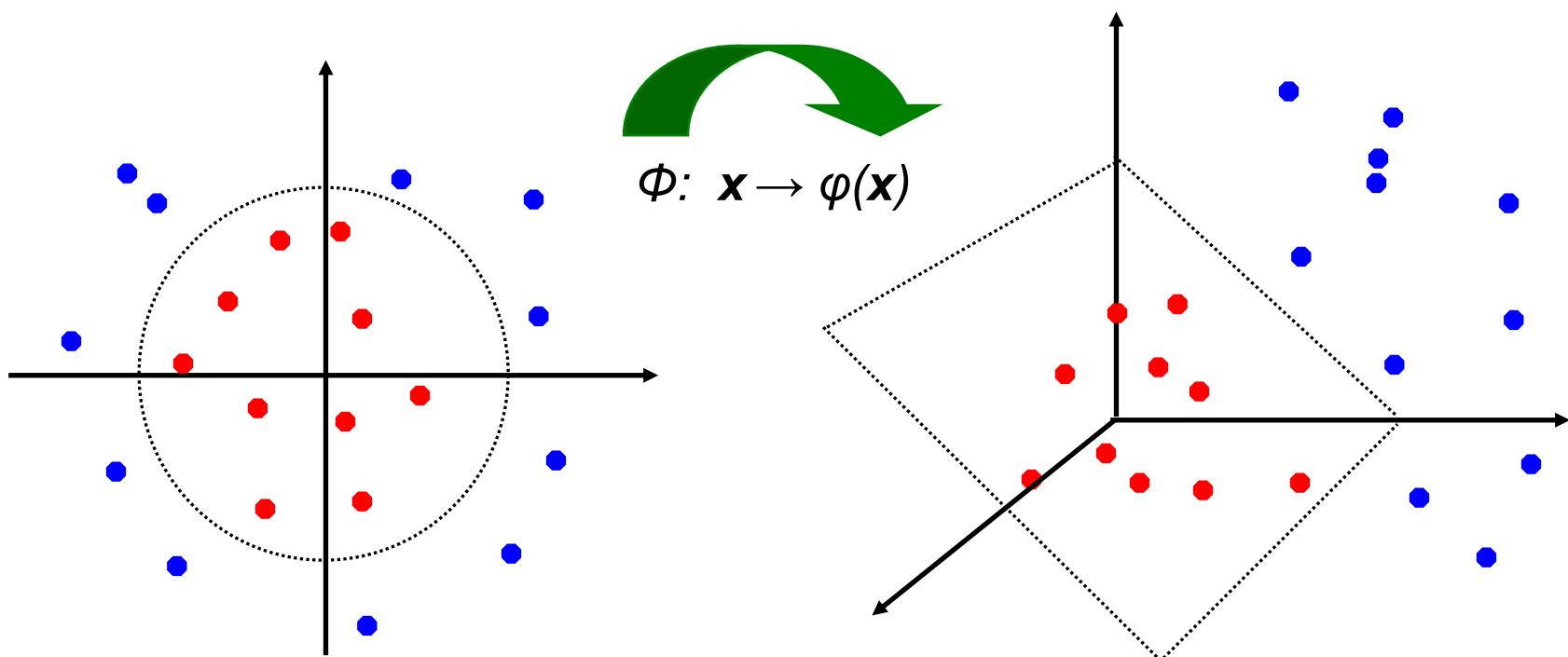


Gaussian

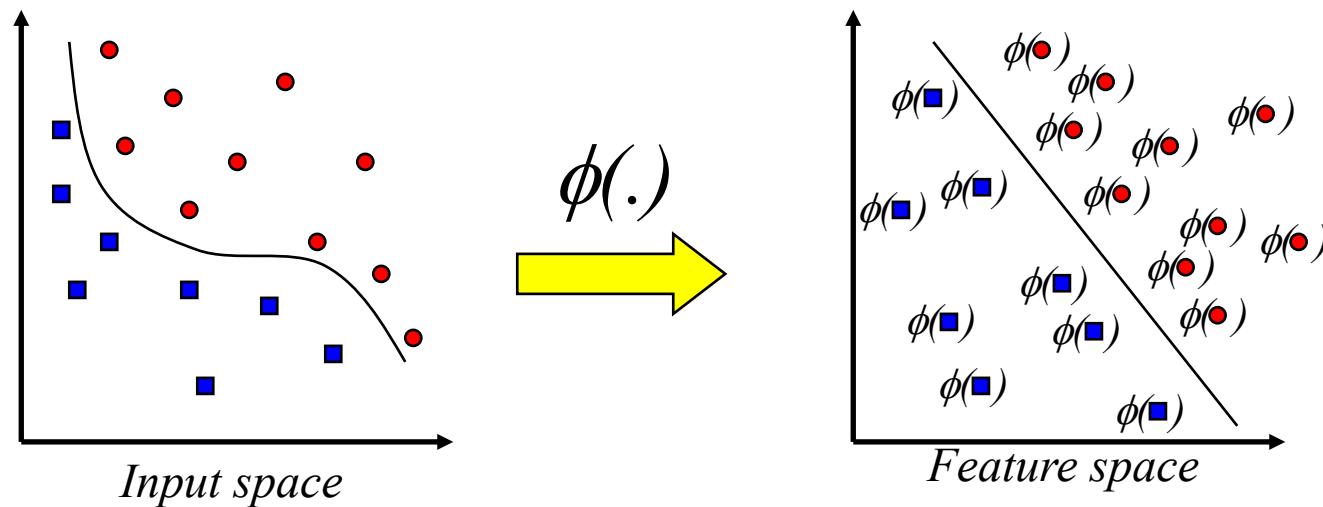


Non-linear SVMs: Feature spaces

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



Transforming the Data



Note: feature space is of higher dimension than the input space in practice

- Computation in the feature space can be costly because it is high dimensional
 - The feature space is typically infinite-dimensional!
- The kernel trick comes to rescue

The Kernel Trick

- Recall the SVM optimization problem

Find $\alpha_1 \dots \alpha_N$ such that

$$Q(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

is maximized s.t.

$$(1) \sum_i \alpha_i y_i = 0, \quad (2) 0 \leq \alpha_i \leq C \text{ for all } \alpha_i$$

- The data points only appear as **inner product**
- As long as we can calculate the inner product in the feature space, we do not need the mapping explicitly
- Many common geometric operations (angles, distances) can be expressed by inner products
- Define the kernel function K by

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

An Example for $\phi(\cdot)$ and $K(\cdot, \cdot)$

- Suppose $\phi(\cdot)$ is given as follows

$$\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- An inner product in the feature space is

$$\langle \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right), \phi\left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) \rangle = (1 + x_1y_1 + x_2y_2)^2$$

- So, if we define the kernel function as follows, there is no need to carry out $\phi(\cdot)$ explicitly

$$K(\mathbf{x}, \mathbf{y}) = (1 + x_1y_1 + x_2y_2)^2$$

- This use of kernel function to avoid carrying out $\phi(\cdot)$ explicitly is known as the **kernel trick**

More on Kernel Functions

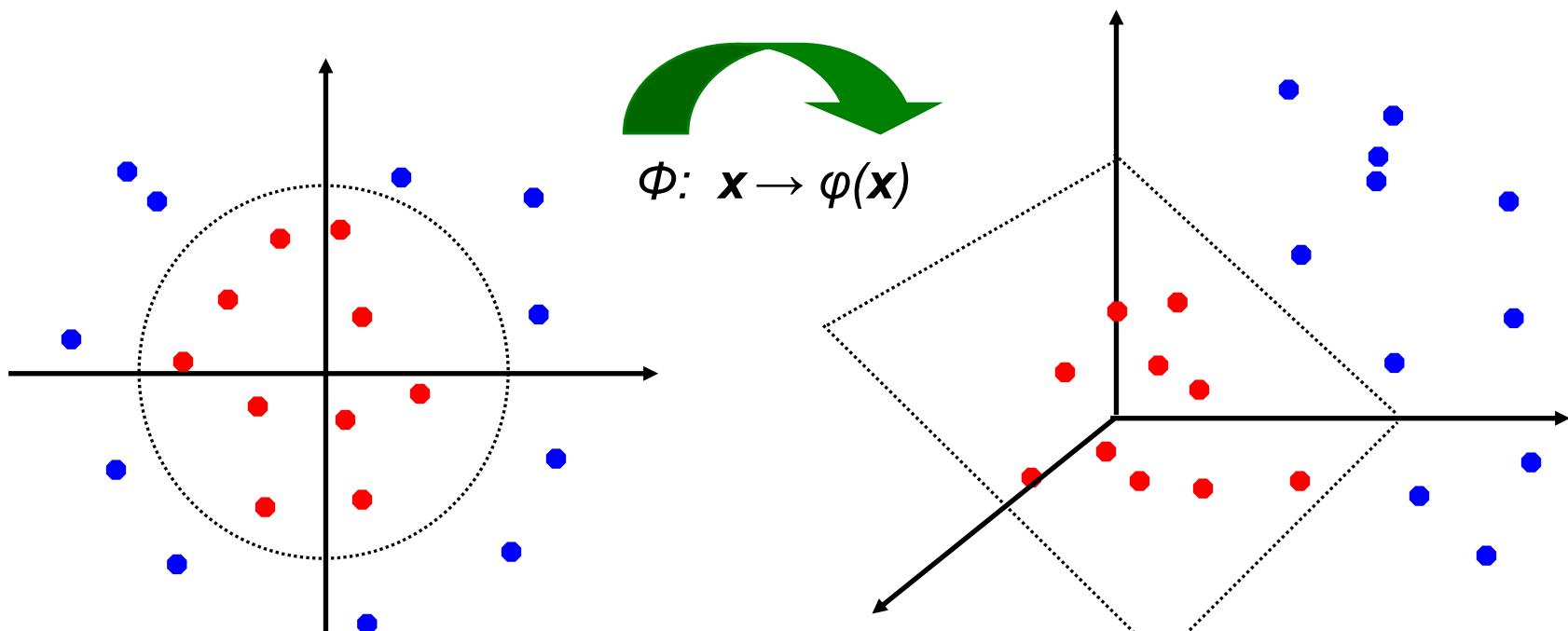
- Not all similarity measures can be used as kernel function, however
 - The kernel function needs to satisfy the **Mercer function**,
 - *Every positive definite symmetric function is a kernel*
- This implies that
 - the n by n kernel matrix,
 - in which the (i,j) -th entry is the $K(\mathbf{x}_i, \mathbf{x}_j)$, is always positive definite
- This also means that optimization problem can be solved in polynomial time!

- Why use kernels?
 - Make non-separable problem separable
 - Map data into better representational space
- Common Kernels
 - **Linear**: $K(x_i, x_j) = \langle x_i \cdot x_j \rangle$
 - **Polynomial of power p**: $K(x_i, x_j) = (\tau + \gamma x_i \cdot x_j)^p$
 - **Gaussian** (radial-basis function network): $K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}, \gamma > 0$
 - **Sigmoidal**: $K(x_i, x_j) = \tanh(\gamma \langle x_i \cdot x_j \rangle + \tau)$

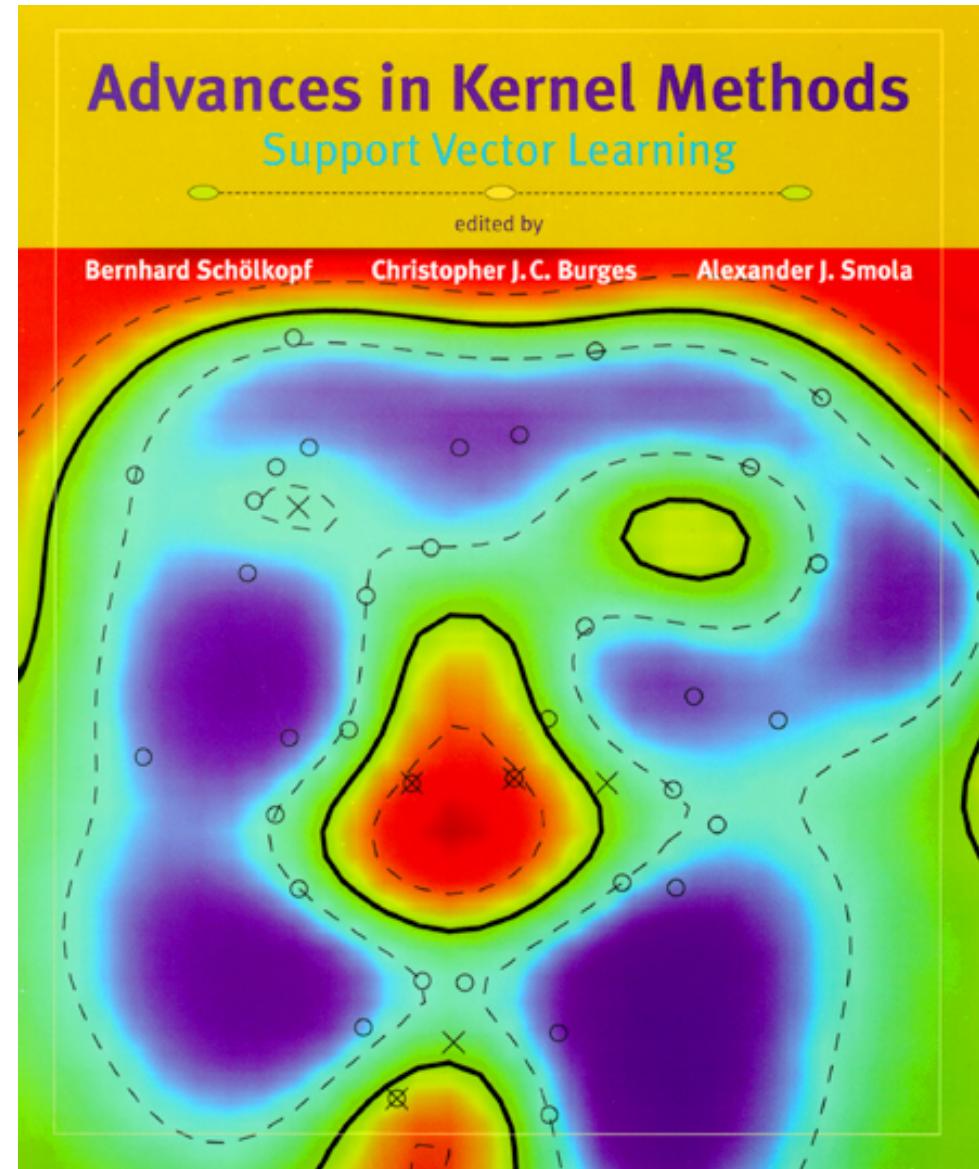
The parameters γ , τ , and p are called the parameters of the kernel.

Non-linear SVMs: Feature spaces

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



$$\begin{aligned}\Phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^3 \\ (x_1, x_2) &\mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}x_1x_2, x_2^2)\end{aligned}$$



- Dual problem formulation:

Find $\alpha_1 \dots \alpha_N$ such that

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

is maximized s.t.

$$(1) \sum_i \alpha_i y_i = 0, \quad (2) 0 \leq \alpha_i \text{ for all } \alpha_i$$

- The solution is:

$$f(x) = \sum \alpha_i y_i K(x_i^T x) + b$$

- Optimization techniques for finding α_i 's remain the same!

Example

- Suppose we have 5 one-dimensional data points
 - $x_1=1, x_2=2, x_3=4, x_4=5, x_5=6$, with 1, 2, 6 as class 1 and 4, 5 as class 2 $\Rightarrow y_1=1, y_2=1, y_3=-1, y_4=-1, y_5=1$
- We use the polynomial kernel of degree 2
 - $K(x,y) = (xy+1)^2$
 - C is set to 100
- We first find α_i ($i=1, \dots, 5$) by

Find $\alpha_1 \dots \alpha_5$ such that

$$Q(\alpha) = \sum_{i=1}^5 \alpha_i - \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \alpha_i \alpha_j y_i y_j (x_i \cdot x_j + 1)^2 \text{ is maximized}$$

s.t. (1) $\sum_{i=1}^5 \alpha_i y_i = 0$, (2) $0 \leq \alpha_i \leq 100$

Example

- By using a Quadratic Problem solver, we get

- $\alpha_1=0, \alpha_2=2.5, \alpha_3=0, \alpha_4=7.333, \alpha_5=4.833$

- Note that the constraints are indeed satisfied

- The support vectors are $\{x_2=2, x_4=5, x_5=6\}$

- The discriminant function is

$$f(z)$$

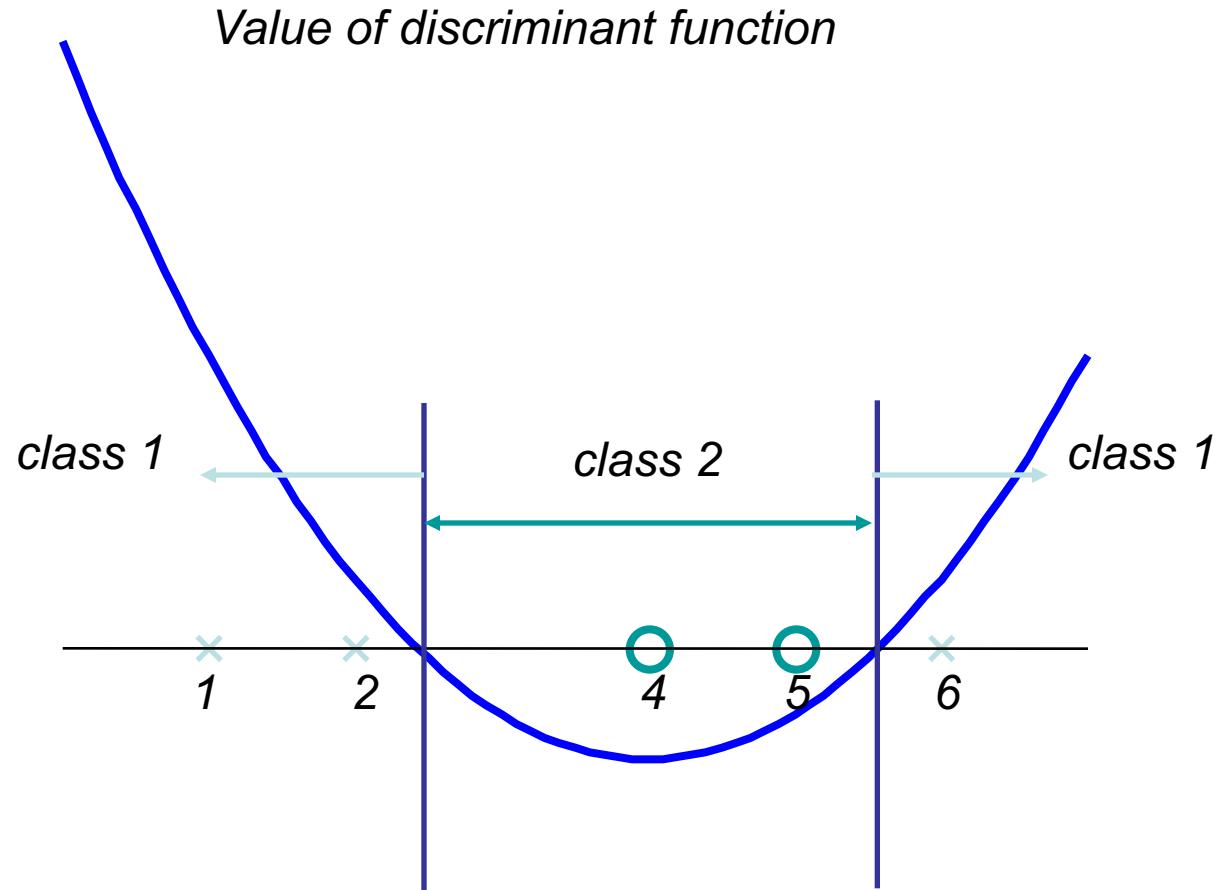
$$= 2.5(1)(2z + 1)^2 + 7.333(-1)(5z + 1)^2 + 4.833(1)(6z + 1)^2 + b$$

$$= 0.6667z^2 - 5.333z + b$$

- b is recovered by solving $f(2)=1$ or by $f(5)=-1$ or by $f(6)=1$, as x_2 and x_5 lie on the line $\phi(\mathbf{w})^T \phi(\mathbf{x}) + b = 1$ and x_4 lies on the line $\phi(\mathbf{w})^T \phi(\mathbf{x}) + b = -1$
- All three give $b=9$  $f(z) = 0.6667z^2 - 5.333z + 9$

$$\begin{array}{c} \alpha_5 \\ \downarrow \\ y_5 \\ \downarrow \\ K(z, x_5) \end{array}$$

Example





TO READ:



A Tutorial on Support Vector Machines for Pattern
Recognition
Christopher J.C. Burges

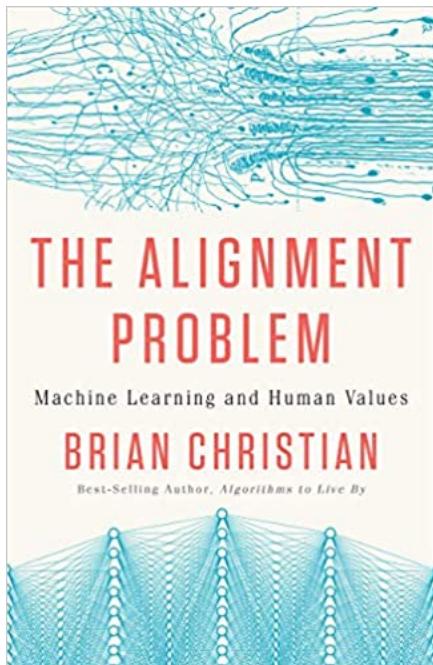
Section 4

With the exception of Sections 2 and 3, the remaining of
the sections are optional

A final reading



The Alignment Problem: Machine Learning and Human Values by Brian Christian



Today's "machine-learning" systems, trained by data, are so effective that we've invited them to see and hear for us—and to make decisions on our behalf. But alarm bells are ringing. Recent years have seen an eruption of concern as the field of machine learning advances. When the systems we attempt to teach will not, in the end, do what we want or what we expect, ethical and potentially existential risks emerge. Researchers call this the alignment problem.

Systems cull résumés until, years later, we discover that they have inherent gender biases. Algorithms decide bail and parole—and appear to assess Black and White defendants differently. We can no longer assume that our mortgage application, or even our medical tests, will be seen by human eyes. And as autonomous vehicles share our streets, we are increasingly putting our lives in their hands.

The mathematical and computational models driving these changes range in complexity from something that can fit on a spreadsheet to a complex system that might credibly be called "artificial intelligence." They are steadily replacing both human judgment and explicitly programmed software.

Week 10

Course. Introduction to Machine Learning

Theory 10. Support Vector Machine

Dr. Maria Salamó Llorente
maria.salamo@ub.edu

Dept. Mathematics and Informatics,
Faculty of Mathematics and Informatics,
University of Barcelona (UB)