

Elective Course

Course Code: CS4103

Autumn 2025-26



Lecture #44

Artificial Intelligence for Data Science

Week-12:

MACHINE LEARNING (Part XII)

Support Vector Machine (SVM)

Course Instructor:

Dr. Monidipa Das

Assistant Professor

Department of Computational and Data Sciences

Indian Institute of Science Education and Research Kolkata, India 741246

Linear SVM Mathematically

The linearly separable case [Revisited]



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

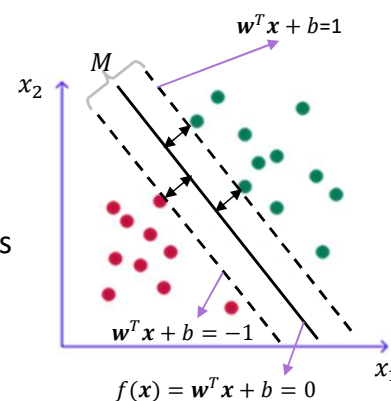
$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \quad \text{if } y_i = 1$$

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

- For support vectors, the inequality becomes an equality
- Then, since each example's distance from the hyperplane is

$$r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$

- The margin is: $M = \frac{2}{\|\mathbf{w}\|}$



Linear SVMs Mathematically (cont.) [Revisited]



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

Find \mathbf{w} and b such that

$$M = \frac{2}{\|\mathbf{w}\|} \text{ is maximized;}$$

and for all $\{(\mathbf{x}_i, y_i)\}$: $\mathbf{w}^T \mathbf{x}_i + b \geq 1$ if $y_i=1$; $\mathbf{w}^T \mathbf{x}_i + b \leq -1$ if $y_i=-1$

- A better formulation ($\min \|\mathbf{w}\| = \max \frac{1}{\|\mathbf{w}\|}$):

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized;}$$

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Solving the Optimization Problem [Revisited]



Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized;}$$

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$



Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \text{ is minimized;}$$

and for all $\{(\mathbf{x}_i, y_i)\}$: $-y_i (\mathbf{w}^T \mathbf{x}_i + b) + 1 \leq 0$

- This is now optimizing a *quadratic* function subject to *linear* constraints
- The solution involves **constructing a dual problem** where a **Lagrange multiplier** α_i is associated with every constraint in the primary problem.

– **Construct the Lagrangian:**

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1]$$

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = \mathbf{0} \Rightarrow \mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \nabla_b \mathcal{L}(\mathbf{w}, b, \alpha) = - \sum_i \alpha_i y_i = \mathbf{0} \Rightarrow \sum_i \alpha_i y_i = 0$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Solving the Optimization Problem [Revisited]



- Plugging back w and b values obtained and simplifying:

$$\mathcal{L}(w, b, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j - b \sum_i \alpha_i y_i$$

$$\Rightarrow \mathcal{L}(w, b, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

- **Dual optimization problem:**

Why do we
use the dual
formulation?

Find $\alpha_1 \dots \alpha_N$ such that
 $Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j x_i^T x_j$ is maximized and
s.t. $\alpha_i \geq 0$ for all α_i
 $\sum_i \alpha_i y_i = 0$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

The Optimization Problem Solution [Revisited]



- Solving the optimization problem involved computing the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points.
- The solution has the form:

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

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

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

- It relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i

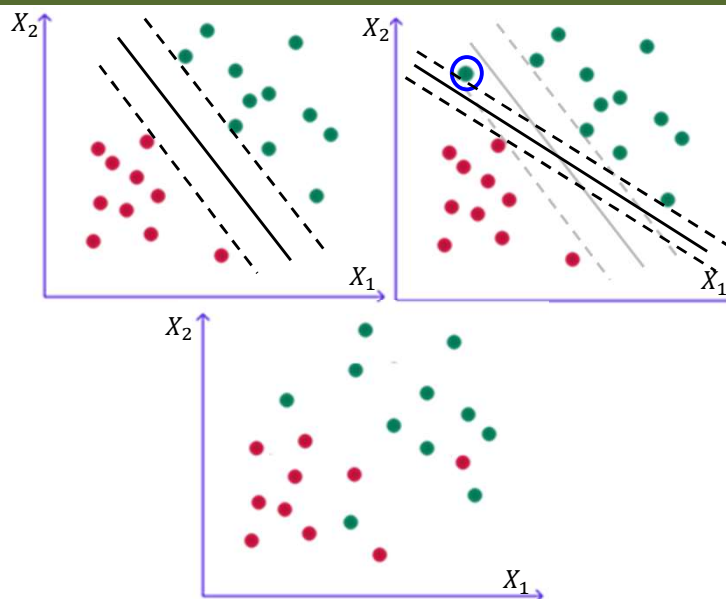
Dr. Monidipa Das, Department of CDS, IISER Kolkata

Maximal Margin Classifier: Limitations

[Revisited]

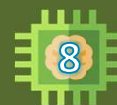


- Although the maximal margin classifier is often successful, it can also lead to overfitting when n is large.
- In many cases ***no separating hyperplane exists***, hence, no maximal margin classifier
- **Remedy:** Idea of soft margin

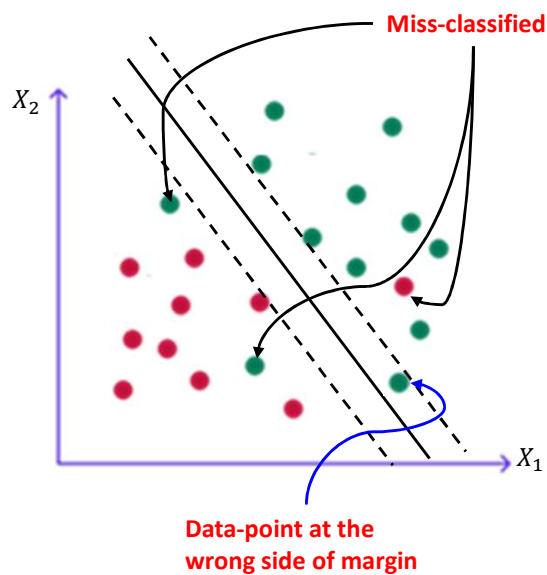


Dr. Monidipa Das, Department of CDS, IISER Kolkata

Soft Margin

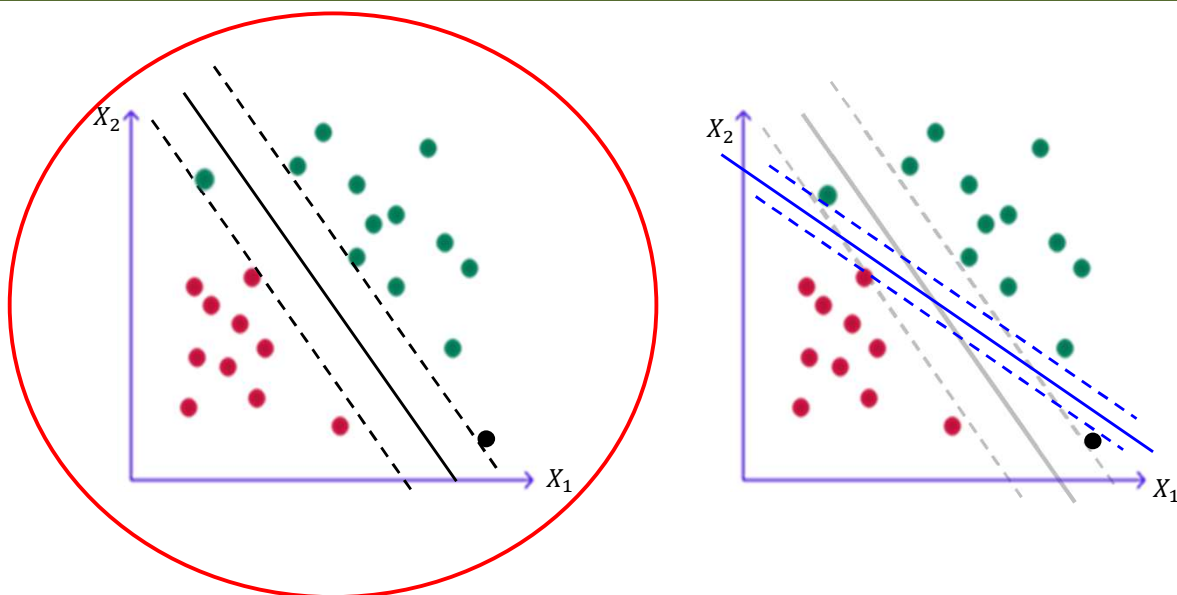


- Allowing some data points to be misclassified or to violate the margin
 - **Greater robustness** to individual observations
 - Better classification of **most** of the training observations



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Which one is more robust?



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Support Vector Classifier



- The generalization of the maximal margin classifier to the non-separable case
- Soft margin classifier
 - The hyperplane is chosen to correctly separate most of the training observations into the two classes, but may misclassify a few observations

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Soft Margin Classification Mathematically



- The old formulation:

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\} \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

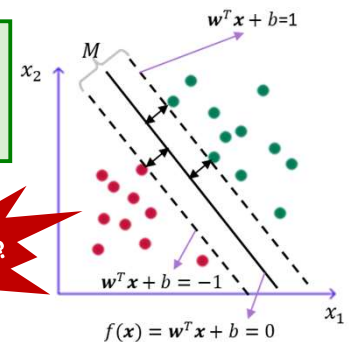
- The new formulation incorporating **slack variables** ξ_i :

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i \text{ is minimized and}$$

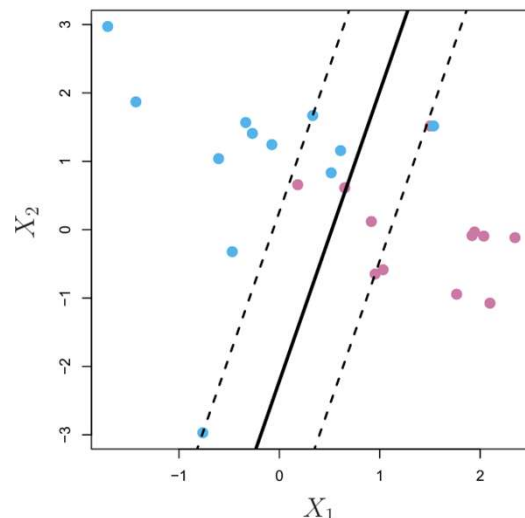
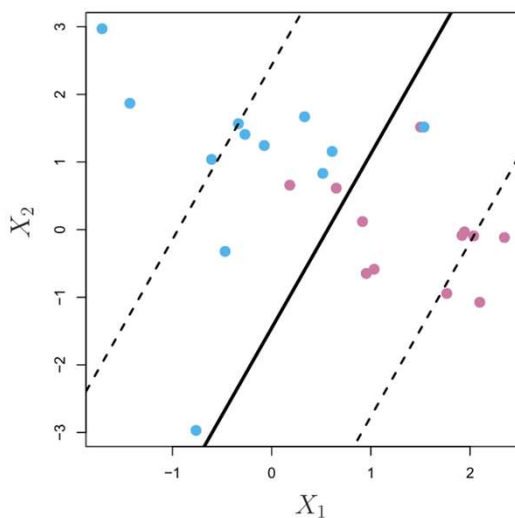
$$\text{for all } \{(\mathbf{x}_i, y_i)\} \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0 \text{ for all } i$$

- Parameter C can be viewed as
 - a way to control weighting between the twin goals
 - a way to control overfitting
 - a regularization term



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Soft Margin Classification: Impact of C



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Soft Margin Classification – Solution



- The dual problem for soft margin classification:



Find $\alpha_1 \dots \alpha_N$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

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

(2) $0 \leq \alpha_i \leq C$ for all α_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:

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

$$b = y_k (1 - \xi_k) - \mathbf{w}^T \mathbf{x}_k \text{ where } k = \underset{k'}{\operatorname{argmax}} \alpha_k$$

\mathbf{w} is not needed explicitly for classification!

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

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Classification with SVMs

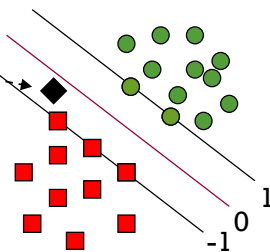


- Given a new point \mathbf{x} , we can score its projection onto the hyperplane normal:
 - I.e., compute score: $\mathbf{w}^T \mathbf{x} + b = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$
 - Decide class based on whether $<$ or > 0
 - Can set confidence threshold t .

Score $> t$: yes

Score $< -t$: no

Else: don't know



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Linear SVMs: Summary



- The classifier is a *separating hyperplane*.
- The most “important” training points are the support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points \mathbf{x}_i are support vectors with non-zero Lagrange 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
 $Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and
 (1) $\sum \alpha_i y_i = 0$
 (2) $0 \leq \alpha_i \leq C$ for all α_i

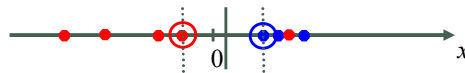
$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

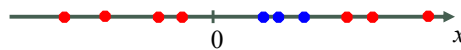
Non-linear SVMs



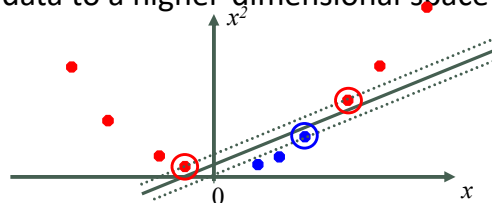
- Datasets that are linearly separable (with some noise) work out great:



- But if the dataset is just too hard?



- What if we map data to a higher-dimensional space?:

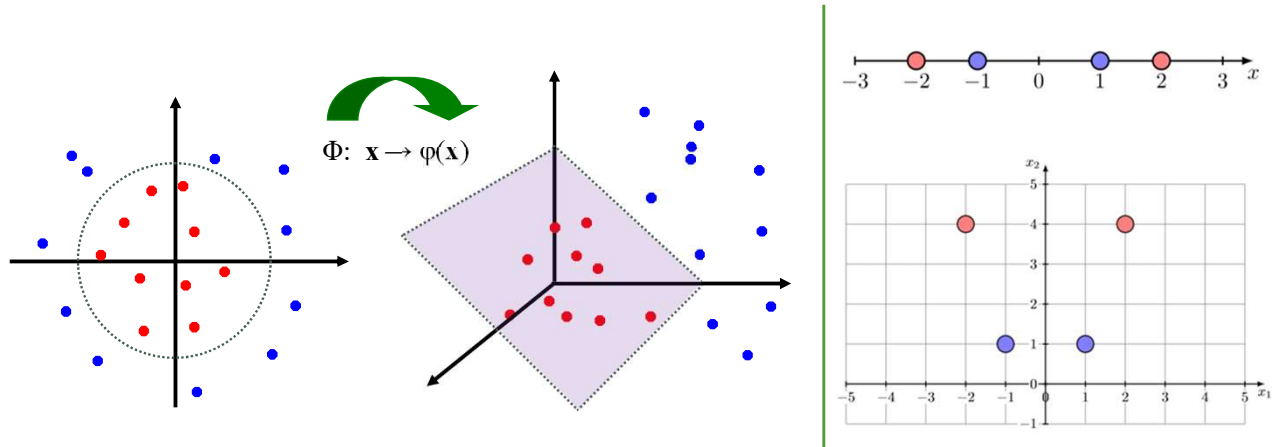


Dr. Monidipa Das, Department of CDS, IISER Kolkata

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:



Dr. Monidipa Das, Department of CDS, IISER Kolkata

The “Kernel Trick”



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

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$$
- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.

- **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 $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$:

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

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Kernels



- Why to use kernels?
 - Map data into better representational space
 - Make non-separable problem separable.
- Common kernels
 - Linear Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
 - Polynomial Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^d$
 - Radial basis function (RBF) Kernel / Gaussian Kernel (infinite dimensional space)

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2} \quad K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2}$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata



Questions?

Dr. Monidipa Das, Department of CDS, IISER Kolkata