

Support Vector Machine

Bùi Tiến Lên

2022



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

Contents



1. Linear Support Vector Machines
2. Kernels Support Vector Machines
3. Multi-class SVM

Notation



symbol	meaning
$a, b, c, N \dots$	scalar number
$\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$	column vector
$\mathbf{X}, \mathbf{Y} \dots$	matrix
\mathbb{R}	set of real numbers
\mathbb{Z}	set of integer numbers
\mathbb{N}	set of natural numbers
\mathbb{R}^D	set of vectors
$\mathcal{X}, \mathcal{Y}, \dots$	set
\mathcal{A}	algorithm

operator	meaning
\mathbf{w}^T	transpose
$\mathbf{X}\mathbf{Y}$	matrix multiplication
\mathbf{X}^{-1}	inverse
$\mathbf{x} \cdot \mathbf{y}$	dot



Linear Support Vector Machines

- The Separable Case
- The Non-Separable Case



Problem Statement

- Training set:

$$(\mathbf{x}_i, y_i)_{i=1 \dots N} \in \mathbb{R}^D \times \{-1, 1\}$$

- We would like to find an hyperplane

$$\mathbf{w}^T \mathbf{x} + b = 0, \quad (\mathbf{w} \in \mathbb{R}^D, b \in \mathbb{R})$$

which **separates** the two classes.

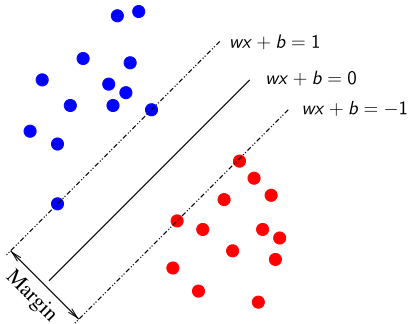


Margin

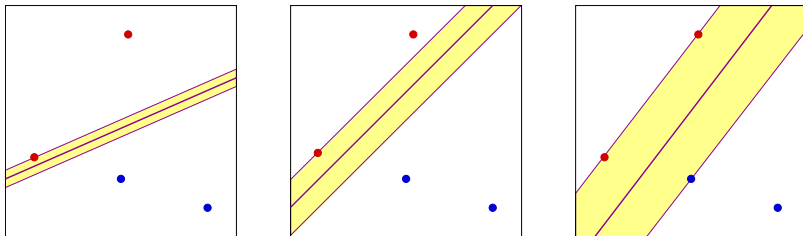
The Separable Case

The Non-Separable
CaseKernels
Support Vector
MachinesMulti-class
SVM

- Let d_+ be the shortest distance from the hyperplane to the closest positive example.
- Let d_- be the shortest distance from the hyperplane to the closest negative example.
- Define the **margin** of the hyperplane to be $d_+ + d_-$.



Better Linear Separation



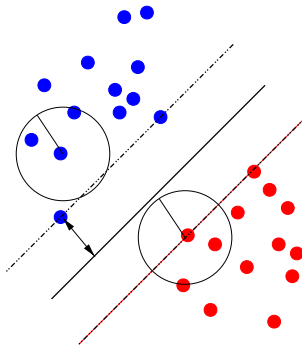
Two questions:

1. Why is bigger margin better?
2. Which w , b maximizes the margin



Why is it Good to Maximize the Margin?

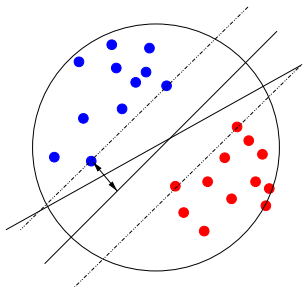
- If training and test data come from the same distribution and all test data are within some Δ distance from the training points. If all points lie at a distance of at least Δ from the separator, and all points are in a bounded sphere, then a small perturbation of the definition of the separator will not hurt.



Why is it Good to Maximize the Margin? (cont.)



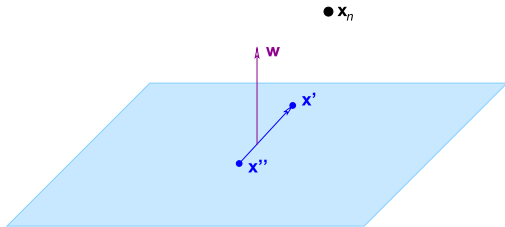
- One can use less bits to encode the separating hyperplane (**Minimum Description Length** principle)





Finding w with large margin

- Let x_n be the nearest data point to the plane $w^T x + b = 0$. How far is it?



- The distance between x_n and the plane $w^T x + b = 0$ where $|w^T x_n + b| = 1$ (normalize w and b)

$$distance = \frac{1}{|w|} \quad (1)$$



A Constrained Optimization Problem

- Representation of hypothesis set

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) \quad (2)$$

- Evaluation

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (3)$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N \quad (4)$$



The Dual Formulation

- Representation of hypothesis set

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i (\mathbf{x} \cdot \mathbf{x}_i) + b \right) \quad (5)$$

- Evaluation

$$\arg \min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^N \alpha_i \quad (6)$$

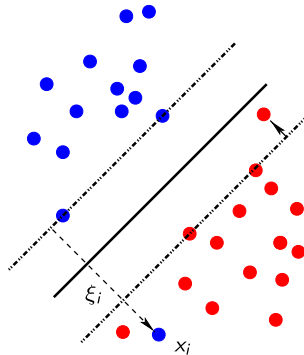
$$\text{subject to} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad (7)$$

$$\alpha_i \geq 0, \quad i = 1, 2, \dots, N \quad (8)$$

- This can be solved using classical **quadratic programming optimization**



- This minimization problem does not have any solution if the two classes are not separable.





Fixing The Bug: “Soft” Margin

- Relax the constraints: use a **soft margin** instead of a **hard margin**.
- We would like to minimize:

$$\arg \min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + c \sum_{i=1}^N \xi_i \quad (9)$$

$$\text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \quad (10)$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N \quad (11)$$



The Dual Formulation

- Representation of hypothesis set

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i (\mathbf{x} \cdot \mathbf{x}_i) + b \right) \quad (12)$$

- Evaluation

$$\arg \min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^N \alpha_i \quad (13)$$

$$\text{subject to} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad (14)$$

$$0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N \quad (15)$$



Support Vector Terminology

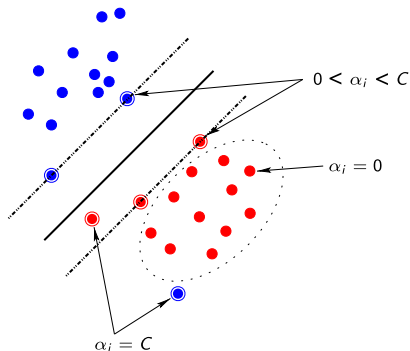
- Training examples \mathbf{x}_i with $\alpha_i > 0$ are **support vectors**.

$$\alpha_i = 0 \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1$$

$$\alpha_i = C \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) < 1$$

$$0 < \alpha_i < C \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$$

(16)



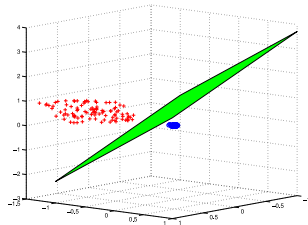
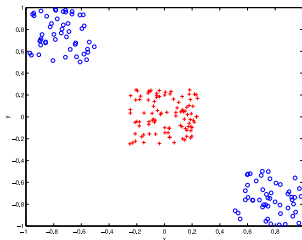


Kernels Support Vector Machines

Non-Linear SVMs



- Project the data into a **higher dimensional space (feature space)**: it should be easier to separate the two classes.
- Given a function $\phi : \mathbb{R}^D \rightarrow \mathcal{F}$, work with $\phi(\mathbf{x}_i)$ instead of working with \mathbf{x}_i .





The Kernel Function

Concept 1

A **kernel** is a function $k(\mathbf{x}, \mathbf{z})$ which represents a dot product in a “hidden” feature space of ϕ .

$$k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) \quad (17)$$

- **Note that:** we have only dot products $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$ to compute; however, this could be very expensive in a high dimensional space.
- **Kernel trick:**

$$\text{instead of } \phi(\mathbf{x}) = \phi \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}, \text{ use } k(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^2$$

Common Kernels



- Polynomial:

$$k(\mathbf{x}, \mathbf{z}) = (u\mathbf{x} \cdot \mathbf{z} + v)^p \quad (u \in \mathbb{R}, v \in \mathbb{R}, p \in \mathbb{N}) \quad (18)$$

- Gaussian:

$$k(\mathbf{x}, \mathbf{z}) = \exp \left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{\sigma^2} \right), \quad \sigma \in \mathbb{R}^+ \quad (19)$$



Techniques for Construction of Kernels

In all the following, k_1, k_2, \dots, k_j are assumed to be valid kernel functions

1. **Scalar multiplication:** The validity of a kernel is conserved after multiplication by a positive scalar, i.e., for any $\alpha > 0$, the function

$$k(\mathbf{x}, \mathbf{z}) = \alpha k_1(\mathbf{x}, \mathbf{z}) \quad (20)$$

2. **Adding a positive constant:** For any positive constant $\alpha > 0$, the function

$$k(\mathbf{x}, \mathbf{z}) = \alpha + k_1(\mathbf{x}, \mathbf{z}) \quad (21)$$

Techniques for Construction of Kernels (cont.)



- 3. Linear combination:** A linear combination of kernel functions involving only positive weights, i.e.,

$$k(\mathbf{x}, \mathbf{z}) = \sum_{j=1}^m \alpha_j k_j(\mathbf{x}, \mathbf{z}), \quad \text{with } \alpha_j > 0 \quad (22)$$

is a valid kernel function.

- 4. Product:** The product of two kernel functions, i.e.,

$$k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) k_2(\mathbf{x}, \mathbf{z}) \quad (23)$$

is a valid kernel function.

Techniques for Construction of Kernels (cont.)



- 5. Polynomial functions of a kernel output:** Given a polynomial $f : \mathbb{R} \rightarrow \mathbb{R}$ with positive coefficients, the function

$$k(\mathbf{x}, \mathbf{z}) = f(k_1(\mathbf{x}, \mathbf{z})) \quad (24)$$

is a valid kernel function.

- 6. Exponential function of a kernel output:** The function

$$k(\mathbf{x}, \mathbf{z}) = \exp(k_1(\mathbf{x}, \mathbf{z})) \quad (25)$$

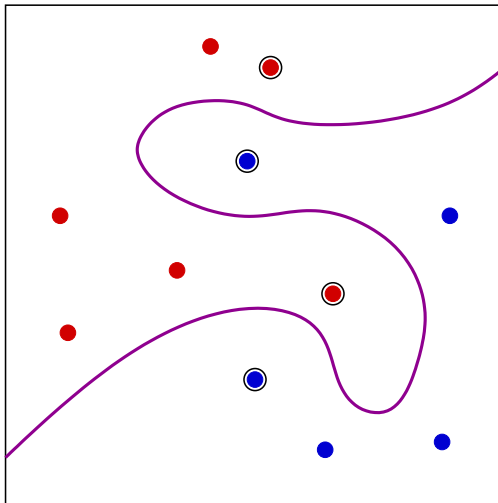
is a valid kernel function.

- 7. Product of matrix and vectors:**

$$k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T A \mathbf{z} \quad (26)$$

where A is a symmetric positive semidefinite matrix.

Decision Boundary and Support Vectors



SVMs in Practice



- In order to tune the **capacity**, the kernel is the most important parameter to choose.
 - Polynomial kernel: increasing the degree will increase the **capacity**.
 - Gaussian kernel: increasing σ will decrease the capacity.
- Tune C , the trade-off between the **margin** and the **errors**.
 - For non-noisy data sets, C usually has not much influence.
 - Carefully choose C for noisy data sets: small values usually give better results.



Multi-class SVM

Multiclass SVM formulations



- There are a few ways of formulating the SVM over multiple classes:
 - One-vs-all
 - One-vs-one
 - Hierarchical
 - **Multiclass**



Score function

The Separable Case
The Non-Separable Case

Kernels
Support Vector
Machines

Multi-class
SVM

Concept 2

The **score function** f that maps the raw features to class scores.

$$z = f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x} \quad (27)$$

input image



\mathbf{W}

1.1	0.2	-0.5	0.1	2.0
3.2	1.5	1.3	2.1	0.0
-1.2	0	0.25	0.2	-0.3

\mathbf{x}

1
56
231
24
2

\times

$=$

-96.8	cat score
437.9	dog score
61.95	ship score

Multiclass SVM loss



- Given the input vector \mathbf{x}_i and the label y_i that specifies the index of the correct class. The multiclass SVM loss (**hinge loss**) for the vector \mathbf{x}_i is then formalized as follows

$$L_i = \sum_{j \neq y_i} \max(0, z_j - z_{y_i} + \Delta) \quad (28)$$

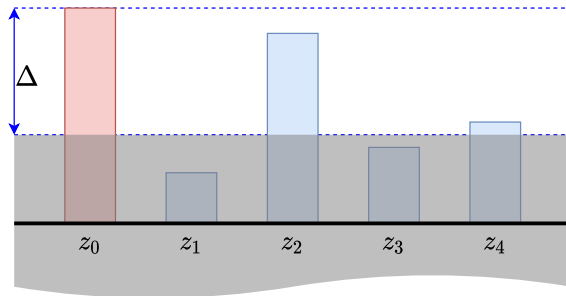
where $\mathbf{z} = f(\mathbf{x}_i; \mathbf{W}) = \mathbf{W}\mathbf{x}_i$



Example

- Suppose that we have five classes $\{0, 1, 2, 3, 4\}$ that receive the scores $\mathbf{z} = [17, 4, 15, 6, 8]$ and the true class $y_i = 0$
- Also assume that $\Delta = 10$

$$L_i = \max(0, 4 - 17 + 10) + \max(0, 15 - 17 + 10) + \max(0, 6 - 17 + 10) + \max(0, 8 - 17 + 10) = 9$$





Regularization loss

- The most common regularization penalty is the L_2 norm that discourages large weights through an elementwise quadratic penalty over all parameters:

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (29)$$

- The data loss (which is the average loss L_i over all examples) and the regularization loss. That is, the full multiclass SVM loss becomes:

$$\mathcal{L} = \underbrace{\mathcal{L}_{\text{data}}}_{\text{data loss}} + \underbrace{\lambda R(W)}_{\text{regularization loss}} = \frac{1}{N} \sum_i L_i + \lambda R(W) \quad (30)$$

- Learning goal:** Find W that minimize

$$\arg \min_W \mathcal{L} \quad (31)$$

Practical considerations



- **Setting Delta:** It can safely be set to $\Delta = 1.0$ in all cases
- **Relation to Binary Support Vector Machine:** The loss for the i -th example (\mathbf{x}_i, y_i) can be written as

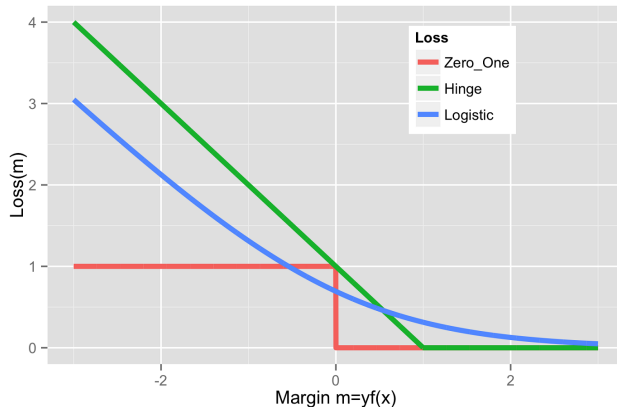
$$L_i = C \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i) + R(\mathbf{w}) \quad (32)$$

where C is a hyperparameter, and $y_i \in \{-1, 1\}$

Binary classification losses



- Perceptron (zero-one)
- SVM (hinge)
- Logistic





SGD for hinge loss

- Consider linear hypothesis space:

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \quad (33)$$

- Hinge loss of (\mathbf{x}, y)

$$\mathcal{L}(\mathbf{x}) = \max(0, 1 - y \mathbf{w}^T \mathbf{x}) \quad (34)$$

- Gradient of hinge loss (\mathbf{x}, y) :

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{x}) = \begin{cases} -y \mathbf{x} & \text{if } y h_{\mathbf{w}}(\mathbf{x}) < 1 \\ 0 & \text{if } y h_{\mathbf{w}}(\mathbf{x}) > 1 \\ \text{undefined} & \text{if } y h_{\mathbf{w}}(\mathbf{x}) = 1 \end{cases} \quad (35)$$

- A point with margin $m = y h_{\mathbf{w}}(\mathbf{x}) = 1$ is correctly classified \rightarrow we can skip SGD update for these points.

References



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

Deep learning.

MIT press.



Lê, B. and Tô, V. (2014).

Cở sở trí tuệ nhân tạo.

Nhà xuất bản Khoa học và Kỹ thuật.



Russell, S. and Norvig, P. (2021).

Artificial intelligence: a modern approach.

Pearson Education Limited.