

Support vector machines

Victor Kitov
v.v.kitov@yandex.ru

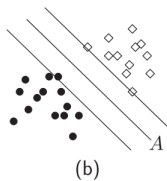
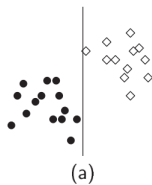
Yandex School of Data Analysis



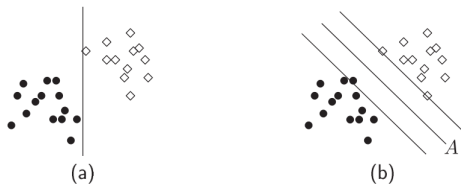
Table of Contents

- 1 Linearly separable case
- 2 Linearly non-separable case

Support vector machines



Support vector machines



Main idea

Select hyperplane maximizing the spread between classes.

Support vector machines

Objects x_i for $i = 1, 2, \dots, n$ lie at distance $b/|w|$ from discriminant hyperplane if

$$\begin{cases} x_i^T w + w_0 \geq b, & y_i = +1 \\ x_i^T w + w_0 \leq -b & y_i = -1 \end{cases} \quad i = 1, 2, \dots, N.$$

This can be rewritten as

$$y_i(x_i^T w + w_0) \geq b, \quad i = 1, 2, \dots, N.$$

The margin is equal to $2b/|w|$. Since w, w_0 and b are defined up to multiplication constant, we can set $b = 1$.

Problem statement

Problem statement:

$$\begin{cases} \frac{1}{2} w^T w \rightarrow \min_{w, w_0} \\ y_i(x_i^T w + w_0) \geq 1, \quad i = 1, 2, \dots, N. \end{cases}$$

Support vectors

non-informative observations: $y_i(x_i^T w + w_0) > 1$

- do not affect the solution

support vectors: $y_i(x_i^T w + w_0) = 1$

- lie at distance $1/|w|$ to separating hyperplane
- affect the the solution.

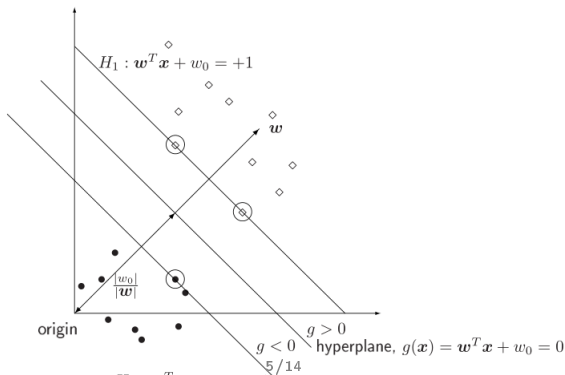
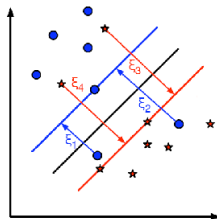


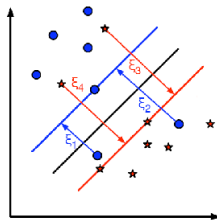
Table of Contents

- 1 Linearly separable case
- 2 Linearly non-separable case

Linearly non-separable case

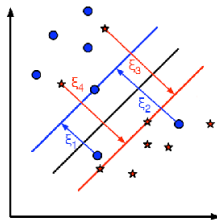


Linearly non-separable case



$$\begin{cases} \frac{1}{2} w^T w \rightarrow \min_{w, w_0} \\ y_i(x_i^T w + w_0) \geq 1, \quad i = 1, 2, \dots, N. \end{cases}$$

Linearly non-separable case



$$\begin{cases} \frac{1}{2} w^T w \rightarrow \min_{w, w_0} \\ y_i(x_i^T w + w_0) \geq 1, \quad i = 1, 2, \dots, N. \end{cases}$$

Problem

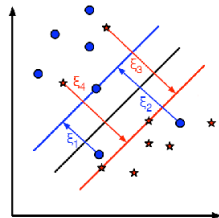
Constraints become incompatible and give empty set!

Linearly non-separable case

No separating hyperplane exists. Errors are permitted by including slack variables ξ_i :

$$\begin{cases} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \rightarrow \min_{w, \xi} \\ y_i(w^T x_i + w_0) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N \\ \xi_i \geq 0, \quad i = 1, 2, \dots, N \end{cases}$$

- Parameter C is the cost for misclassification and controls the bias-variance trade-off.
- It is chosen on validation set.
- Other penalties are possible, e.g. $C \sum_i \xi_i^2$.



Classification of training objects

- **Non-informative objects:**

- $y_i(w^T x_i + w_0) > 1$

- **Support vectors SV :**

- $y_i(w^T x_i + w_0) \leq 1$

- **boundary support vectors \widetilde{SV} :**

- $y_i(w^T x_i + w_0) = 1$

- **violating support vectors:**

- $y_i(w^T x_i + w_0) > 0$: violating support vector is correctly classified.

- $y_i(w^T x_i + w_0) < 0$: violating support vector is misclassified.

Solution of linearly non-separable case

- 1 Solve dual task to find α_i^* , $i = 1, 2, \dots, N$

$$\begin{cases} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \end{cases}$$

- 2 Find optimal w_0 :

$$w_0 = \frac{1}{n_{\tilde{S}V}} \left(\sum_{j \in \tilde{S}V} y_j - \sum_{j \in \tilde{S}V} \sum_{i \in SV} \alpha_i^* y_i \langle x_i, x_j \rangle \right)$$

- 3 Make prediction for new x :

$$\hat{y} = \text{sign}[w^T x + w_0] = \text{sign} \left[\sum_{i \in SV} \alpha_i^* y_i \langle x_i, x \rangle + w_0 \right]$$

Making predictions

- 1 Solve dual task to find α_i^* , $i = 1, 2, \dots, N$

$$\begin{cases} L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \end{cases}$$

- 2 Find optimal w_0 :

$$w_0 = \frac{1}{n_{\tilde{S}V}} \left(\sum_{j \in \tilde{S}V} y_j - \sum_{j \in \tilde{S}V} \sum_{i \in SV} \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

- 3 Make prediction for new \mathbf{x} :

$$\hat{y} = \text{sign}[w^T \mathbf{x} + w_0] = \text{sign} \left[\sum_{i \in SV} \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + w_0 \right]$$

- On all steps we don't need exact feature representations, only scalar products $\langle \mathbf{x}, \mathbf{x}' \rangle$!

Kernel trick generalization

- 1 Solve dual task to find α_i^* , $i = 1, 2, \dots, N$

$$\begin{cases} L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j) \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \end{cases}$$

- 2 Find optimal w_0 :

$$w_0 = \frac{1}{n_{\tilde{S}V}} \left(\sum_{j \in \tilde{S}V} y_j - \sum_{j \in \tilde{S}V} \sum_{i \in SV} \alpha_i^* y_i K(x_i, x_j) \right)$$

- 3 Make prediction for new x :

$$\hat{y} = \text{sign}[w^T x + w_0] = \text{sign} \left[\sum_{i \in SV} \alpha_i^* y_i K(x_i, x) + w_0 \right]$$

- We replaced $\langle x, x' \rangle \rightarrow K(x, x')$ for $K(x, x') = \langle \phi(x), \phi(x') \rangle$ for some feature transformation $\phi(\cdot)$.

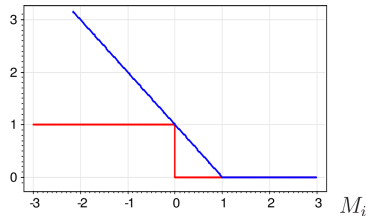
Another view on SVM

Optimization problem:

$$\begin{cases} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \rightarrow \min_{w, \xi} \\ y_i(w^T x_i + w_0) = M_i(w, w_0) \geq 1 - \xi_i, \\ \xi_i \geq 0, i = 1, 2, \dots, N \end{cases}$$

can be rewritten as

$$\frac{1}{2C} |w|^2 + \sum_{i=1}^N [1 - M_i(w, w_0)]_+ \rightarrow \min_{w, \xi}$$



Thus SVM is linear discriminant function with cost approximated with $\mathcal{L}(M) = [1 - M]_+$ and L_2 regularization.

Sparsity of solution

- SVM solution depends only on support vectors
- This is also clear from loss function, satisfying $\mathcal{L}(M) = 0$ for $M \geq 1$.
 - objects with $\text{margin} \geq 1$ don't affect solution!
- Sparsity causes SVM to be less robust to outliers
 - because outliers are always support vectors