

Support vector machines

Victor Kitov

Table of contents

- 1 Optimization reminder
- 2 Support vector machines
 - Linearly separable case
 - Linearly non-separable case
- 3 Addition

Kuhn-Takker conditions

Consider the optimization task:

$$\begin{cases} f(x) \rightarrow \min_x \\ g_i(x) \leq 0 \quad i = 1, 2, \dots, m \end{cases} \quad (1)$$

Theorem (necessary conditions for optimality):

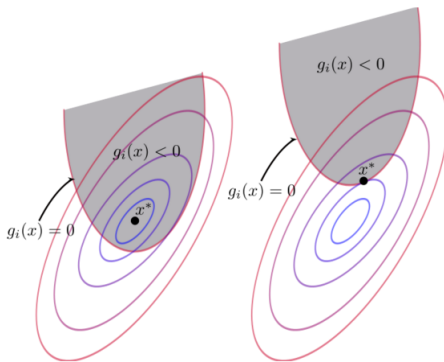
Let

- x^* - be the solution to (1),
- $f(x^*)$ and $g_i(x^*)$, $i = 1, 2, \dots, m$ - continuously differentiable at x^* .
- one of the conditions of regularity is satisfied

Then coefficients $\lambda_1, \lambda_2, \dots, \lambda_m$ exist, such that x^* satisfies the conditions:

$$\begin{cases} \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) = 0 & \text{stationarity} \\ g_i(x^*) \leq 0 & \text{feasibility} \\ \lambda_i \geq 0 & \text{non-negativity} \\ \lambda_i g_i(x^*) = 0 & \text{complementary slackness} \end{cases} \quad (2)$$

Illustration of constrained optimization



Kuhn-Takker conditions

Possible regularity conditions:

- $\{\nabla g_j(x^*), j \in J\}$ - linearly independent, where J are indexes of active constraints $J = \{j : g_j(x^*) = 0\}$.
- Slater condition: $\exists x : g_i(x) < 0 \forall i$ (applicable only when $f(x)$ and $g_i(x), i = 1, 2, \dots, m$ are convex)

Sufficient conditions of optimality:

If $f(x)$ and $g_i(x), i = 1, 2, \dots, m$ are convex, Kuhn-Takker conditions (2) and Slater conditions become sufficient for x^* to be the solution of (1).

Convex optimization

Why convexity of $f(x)$ and $g_i(x)$, $i = 1, 2, \dots, m$ is convenient:

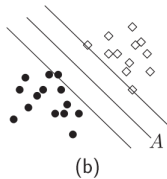
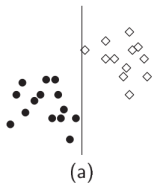
- All local minimums become global minimums
- The set of minimums is convex
- If $f(x)$ is strictly convex and minimum exists, then it is unique.

Table of contents

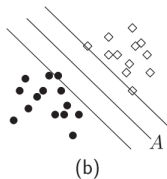
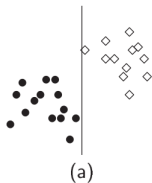
- 1 Optimization reminder
- 2 Support vector machines
 - Linearly separable case
 - Linearly non-separable case
- 3 Addition

- 2 Support vector machines
 - Linearly separable case
 - 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:

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

Problem statement

Problem statement:

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

Lagrangian:

$$L = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}^T \mathbf{x} + w_0) - 1)$$

By Karush-Kuhn-Takker the solution satisfies:

$$\begin{cases} \frac{\partial L}{\partial \mathbf{w}} = \mathbf{0}, \quad \frac{\partial L}{\partial w_0} = 0 \\ y_i(\mathbf{x}_i^T \mathbf{w} + w_0) - 1 \geq 0, \\ \alpha_i (y_i(\mathbf{x}_i^T \mathbf{w} + w_0) - 1) = 0, \\ \alpha_i \geq 0, \quad i = 1, 2, \dots, N \end{cases}$$

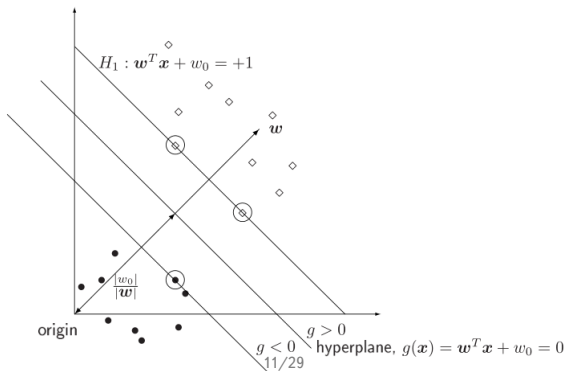
Support vectors

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

- do not affect the solution

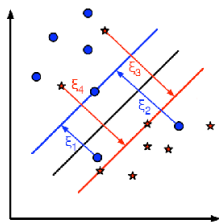
support vectors: $y_i(x_i^T w + w_0) = 1$

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

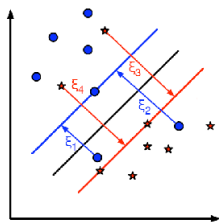


- 2 Support vector machines
 - Linearly separable case
 - Linearly non-separable case

Linearly non-separable case

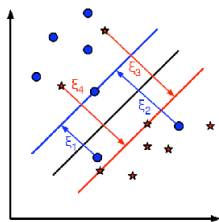


Linearly non-separable case



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

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

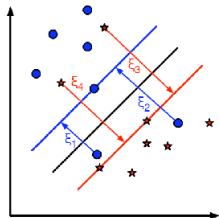
Constraints become incompatible and give empty set!

Linearly non-separable case

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

$$\begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \rightarrow \min_{\mathbf{w}, \xi} \\ y_i (\mathbf{w}^T \mathbf{x}_i + \mathbf{w}_0) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N \\ \xi_i \geq 0, \quad i = 1, 2, \dots, N \end{cases}$$

- 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. $C \sum_i \xi_i^2$.



Linearly non-separable case

Lagrangian:

$$L = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i) - \sum_{i=1}^N r_i \xi_i$$

By Karush-Kuhn-Takker conditions, the solution satisfies constraints:

$$\begin{cases} \frac{\partial L_P}{\partial \mathbf{w}} = \mathbf{0}, \quad \frac{\partial L_P}{\partial w_0} = 0, \quad \frac{\partial L_P}{\partial \xi_i} = 0 \\ \xi_i \geq 0, \quad \alpha_i \geq 0, \quad r_i \geq 0 \\ y_i (\mathbf{x}_i^T \mathbf{w} + w_0) \geq 1 - \xi_i, \\ \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i) = 0 \\ r_i \xi_i = 0, \quad i = 1, 2, \dots, N \end{cases}$$

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.

Solving Karush-Kuhn-Takker conditions

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{0} : \mathbf{w} = \sum_{i=1}^N \alpha_i \mathbf{y}_i \mathbf{x}_i \quad (3)$$

$$\frac{\partial L}{\partial \mathbf{w}_0} = 0 : \sum_{i=1}^N \alpha_i \mathbf{y}_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = 0 : C - \alpha_i - r_i = 0 \quad (4)$$

Substituting these constraints into L , we obtain the *dual problem*¹:

$$\begin{cases} L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i \mathbf{y}_i = 0 \\ 0 \leq \alpha_i \leq C \quad (\text{using (4) and that } \alpha_i \geq 0, r_i \geq 0) \end{cases} \quad (5)$$

¹Dual Lagrangian is maximized because original Lagrangian has saddlepoint in optimum, min for $\mathbf{w}, \mathbf{w}_0, \xi_i$ and max for α_i, r_i .

Comments on support vectors

- non support vectors: $y_i(w^T x_i + w_0) > 1 \Leftrightarrow \xi_i = 0$,
 $y_i(w^T x_i + w_0) - 1 + \xi_i > 0 \Rightarrow \alpha_i = 0$
 - support vectors SV will have $\alpha_i > 0$.
- non-boundary support vectors $SV \setminus \widetilde{SV}$: $y_i(w^T x_i + w_0) < 1$
 $\Leftrightarrow \xi_i > 0 \Rightarrow r_i = 0 \Leftrightarrow \alpha_i = C$.
- boundary support vectors \widetilde{SV} : $y_i(w^T x_i + w_0) = 1 \Rightarrow \xi_i = 0$
 - since $\alpha_i \in [0, C]$, $\alpha_i \in (0, C)$ for boundary support vectors.

Solution

- 1 Solve (5) to find optimal dual variables α_i^*
- 2 Using (3) and that $\alpha_i^* = 0$ for non support vectors, find optimal w

$$w = \sum_{i \in \mathcal{SV}} \alpha_i^* y_i x_i$$

- 3 w_0 can be found from any edge equality for boundary support vector:

$$y_i(x_i^T w + w_0) = 1, \forall i \in \widetilde{\mathcal{SV}} \quad (6)$$

Solution for w_0

By multiplying (6) by y_i obtain

$$x_i^T w + w_0 = y_i \quad \forall i \in \widetilde{S\mathcal{V}}$$

By summing over all $i \in \widetilde{S\mathcal{V}}$ for more robust solution we obtain

$$n_{\widetilde{S\mathcal{V}}} w_0 = \sum_{j \in \widetilde{S\mathcal{V}}} (y_j - x_j^T w) = \sum_{j \in \widetilde{S\mathcal{V}}} y_j - \sum_{j \in \widetilde{S\mathcal{V}}} x_j^T \sum_{i \in S\mathcal{V}} \alpha_i^* y_i x_i$$

where $n_{\widetilde{S\mathcal{V}}}$ is the number of boundary support vectors.

Final solution for w_0 :

$$w_0 = \frac{1}{n_{\widetilde{S\mathcal{V}}}} \left(\sum_{j \in \widetilde{S\mathcal{V}}} y_j - \sum_{j \in \widetilde{S\mathcal{V}}} \sum_{i \in S\mathcal{V}} \alpha_i^* y_i x_j^T x_i \right)$$

Making predictions

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

$$\begin{cases} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \quad (\text{using (4) and that } \alpha_i \geq 0, r_i \geq 0) \end{cases}$$

- 2 Find optimal w_0 :

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

- 3 Make prediction for new x :

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

Making predictions

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

$$\begin{cases} L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \quad (\text{using (4) and that } \alpha_i \geq 0, r_i \geq 0) \end{cases}$$

- 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^* , $i = 1, 2, \dots, N$

$$\begin{cases} L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \rightarrow \max_{\alpha} \\ \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \quad (\text{using (4) and that } \alpha_i \geq 0, r_i \geq 0) \end{cases}$$

- 2 Find optimal w_0 :

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

- 3 Make prediction for new x :

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

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

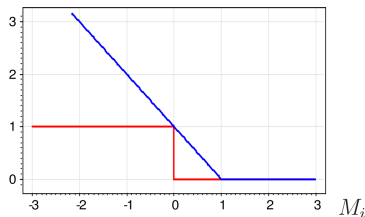
Another view on SVM

Optimization problem:

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

can be rewritten as

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



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 margin ≥ 1 don't affect solution!
- Sparsity causes SVM to be less robust to outliers
 - because outliers are always support vectors

Multiclass classification

C classes $\omega_1, \omega_2, \dots, \omega_C$.

- One-against-all:
 - build C binary classifiers, classifying class ω_i against other classes
 - select the class with highest margin
- One-against-one:
 - build $C(C-1)/2$ classifiers, classifying class ω_i against ω_j .
 - select the class having maximum votes
- Multiclass variant of initial algorithm

Table of contents

- 1 Optimization reminder
- 2 Support vector machines
 - Linearly separable case
 - Linearly non-separable case
- 3 Addition

SVM regression

Predict real-valued output with

$$\hat{y}(x) = w^T x + w_0$$

where parameters w, w_0 are found from

$$\begin{cases} (x^T x_n + w_0) - y_n \leq \varepsilon + \xi_n \\ y_n - (x^T x_n + w_0) \leq \varepsilon + \tilde{\xi}_n \\ \xi_n, \tilde{\xi}_n \geq 0, \quad n = 1, 2, \dots, N. \\ \frac{1}{2} w^T w + C \sum_{n=1}^N (\xi_n + \tilde{\xi}_n) \rightarrow \min_{w, w_0, \xi_n, \tilde{\xi}_n} \end{cases}$$

Gives ε -insensitive loss!

Multiclass SVM

C discriminant functions are built simultaneously:

$$g_k(x) = (w^k)^T x + w_0^k$$

Linearly separable case:

$$\begin{cases} \sum_{k=1}^C (w^k)^T w^k \rightarrow \min_w \\ (w^{y(i)})^T x + w_0^{y(i)} - (w^k)^T x - w_0^k \geq 1 \quad \forall k \neq y(i), \\ i = 1, 2, \dots, N. \end{cases}$$

Linearly non-separable case:

$$\begin{cases} \sum_{k=1}^C (w^k)^T w^k + C \sum_{i=1}^N \xi_i \rightarrow \min_w \\ (w^{y(i)})^T x + w_0^{y(i)} - (w^k)^T x - w_0^k \geq 1 - \xi_i \quad \forall k \neq y(i), \\ \xi_i \geq 0, \quad i = 1, 2, \dots, N. \end{cases}$$

Is slower, but shows similar accuracy to usual SVM.