# Support vector machines

### Victor Kitov

v.v.kitov@yandex.ru
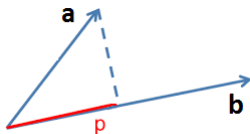
**Yandex**
School of Data Analysis.

# Table of Contents

# Reminder

1. $a = [a^1, \ldots a^D]^T$, $b = [b^1, \ldots b^D]^T$
2. Scalar product $\langle a, b \rangle = a^T b = \sum_{d=1}^{D} a_d b_b$
3. $a \perp b$ means that $\langle a, b \rangle = 0$
4. Norm $\|a\| = \sqrt{\langle a, a \rangle}$
5. Distance $\rho(a, b) = \|a - b\| = \sqrt{\langle a - b, a - b \rangle}$



- $p = \langle a, \frac{b}{\|b\|} \rangle$
- $|p| = \left| a, \frac{b}{\|b\|} \right|$ - unsigned projection length

# Orthogonal vector to hyperplane

### Theorem 1

Vector $w$ is orthogonal to hyperplane $w^T x + w_0 = 0$

*Proof.* Consider arbitrary $x_A, x_B \in \{x : w^T x + w_0 = 0\}$:

$$w^T x_A + w_0 = 0 \tag{1}$$

$$w^T x_B + w_0 = 0 \tag{2}$$

By substracting (2) from (1), obtain $w^T(x_A - x_B) = 0$, so $w$ is orthogonal to hyperplane. $\square$

# Distance from point to hyperplane

### Theorem 2

Distance from point $x$ to hyperplane $w^T x + w_0 = 0$ is equal to $\frac{w^T x + w_0}{\|w_0\|}$ .

Proof. Project $x$ on the hyperplane, let the projection be $p$ and complement $h = x - p$, orthogonal to hyperplane. Then

$$x = p + h$$

Since $p$ lies on the hyperplane,

$$w^T p + w_0 = 0$$

Since $h$ is orthogonal to hyperplane and according to theorem 1

$$h = r \frac{w}{\|w\|}, \ r \in \mathbb{R} \text{ - distance to hyperplane.}$$

# Distance from point to hyperplane

$$x = p + r\frac{w}{\|w\|}$$

After multiplication by $w$ and addition of $w_0$:

$$w^T x + w_0 = w^T p + w_0 + r\frac{w^T w}{\|w\|} = r\|w\|$$

because $w^T p + w_0 = 0$ and $\|w\| = \sqrt{w^T w}$. So we get, that

$$r = \frac{w^T x + w_0}{\|w\|}$$

Comments:

- From one side of hyperplane $r > 0 \Leftrightarrow w^T x + w_0 > 0$
- From the other side $r < 0 \Leftrightarrow w^T x + w_0 < 0$.
- Distance from hyperplane to origin 0 is $\frac{w_0}{\|w\|}$. So $w_0$ accounts for hyperplane offset.

# Binary linear classifier geometric interpretation

Binary linear classifier:

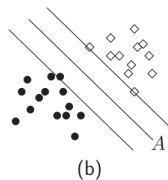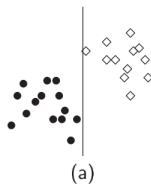$$\widehat{y}(x) = \text{sign}\left(w^T x + w_0\right)$$

divides feature space by hyperplane $w^T x + w_0 = 0$.

- Confidence of decision is proportional to distance to hyperplane $\frac{|w^T x + w_0|}{\|w\|}$.
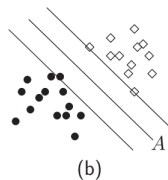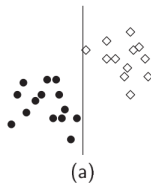- $w^T x + w_0$ is the confidence that class is positive.

# Table of Contents

# Support vector machines



(a)                    (b)

# Support vector machines



### Main idea

Select hyperplane maximizing the spread between classes.

# Support vector machines

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

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

This can be rewritten as

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

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:

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

## 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
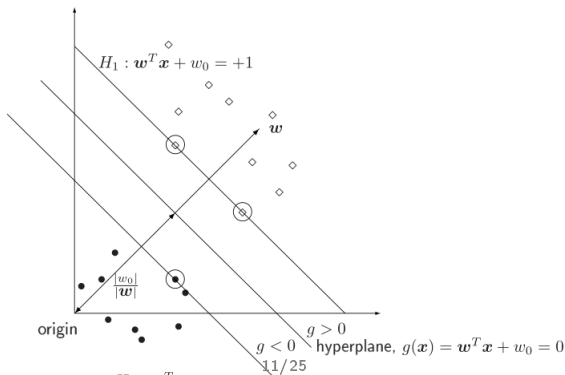- affect the the solution.



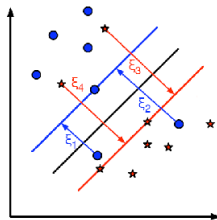$H_1 : \boldsymbol{w}^T \boldsymbol{x} + w_0 = +1$

$\boldsymbol{w}$

$\frac{|w_0|}{|\boldsymbol{w}|}$

origin

$g > 0$

$g < 0$ hyperplane, $g(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + w_0 = 0$

# Table of Contents

# Linearly non-separable case

# Linearly non-separable case



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

# Linearly non-separable case



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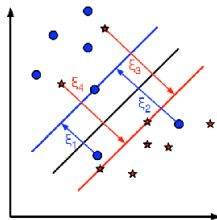y_i(x_i^T w + w_0) \geq 1, \quad i = 1, 2, ...N.
\end{cases}
$$

## Problem

Constraints become incompatible and give empty set!

# Linearly non-separable case

No separating hyperplane exists. Errors are permitted by including slack variables $\xi_i$:

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

- 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.

# SVM with unconstrained optimization

Optimization problem:

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

can be rewritten as



$$\frac{1}{2C} \|w\|_2^2 + \sum_{i=1}^{N} [1 - M_i(w, w_0)]_+ \to \min_{w, w_0, \xi}$$

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

# Table of Contents

# Dual problem

Solving Karush-Kuhn-Takker conditions, get **dual optimization problem**:

$$\begin{cases} L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j x_i^T x_j \to \max_\alpha \\ \sum_{n=1}^{N} \alpha_n y_n = 0 \\ 0 \le \alpha_n \le C, \quad n = \overline{1, N} \end{cases} \quad (3)$$

It is standard quadratic programming task.

# Comments on support vectors

- **non-informative vectors**: $y_i(w^T x_i + w_0) > 1$ have $\alpha_i = 0$
- **non-boundary support vectors** $SV \setminus \tilde{SV}$:
  $y_i(w^T x_i + w_0) < 1$ have $\alpha_i = C$.
- **boundary support vectors** $\widetilde{\mathcal{SV}}$: $y_i(w^T x_i + w_0) = 1$
  Typically $\alpha_i \in (0, C)$, though $\alpha_i = 0, C$ are possible as special cases.

# Solution

1. Solve (3) to find optimal dual variables $\alpha_i^*$
2. Find optimal $w$ ($\alpha_i^* \neq 0$ only for support vectors):

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

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

$$y_i(x_i^T w + w_0) = 1, \ \forall i \in \widetilde{\mathcal{SV}} \qquad (4)$$

# Solution for $w_0$

By multiplyting (4) by $y_i$ obtain

$$x_i^T w + w_0 = y_i \quad \forall i \in \widetilde{\mathcal{SV}} \tag{5}$$

Get more numerically stable from summing 5 over all $i \in \widetilde{\mathcal{SV}}$:

$$n_{\widetilde{SV}} w_0 = \sum_{j \in \widetilde{SV}} \left( y_j - x_j^T w \right) = \sum_{j \in \widetilde{SV}} y_j - \sum_{j \in \widetilde{SV}} x_j^T w, \quad n_{\widetilde{SV}} = \left| \widetilde{SV} \right|$$

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

If there exist no boundary support vectors (only violating SV), then find $w_0$ by grid search.

# Making predictions

1. Solve dual task to find $\alpha_i^*$, $i = 1, 2, ...N$

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

2. Find optimal $w_0$:

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

3. Make prediction for new $x$:

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

# Making predictions

1. Solve dual task to find $\alpha_i^*$, $i = 1, 2, ...N$

$$\begin{cases} L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \to \max_\alpha \\ \sum_{i=1}^{N} \alpha_i y_i = 0 \\ 0 \le \alpha_i \le C \quad \text{(using (\textbf{??}) and that } \alpha_i \ge 0, \ r_i \ge 0\text{)} \end{cases}$$

2. Find optimal $w_0$:

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

3. Make prediction for new $x$:

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

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

# Kernel trick generalization

1. Solve dual task to find $\alpha_i^*$, $i = 1, 2, ... N$

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

2. Find optimal $w_0$:

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

3. Make prediction for new $x$:

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

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

# Summary

- SVM - linear classifier with $L_2$ regularization and hinge loss.
- Geometrically SVM maximizes border between classes.
- Solution depends only on support vectors, having margin$\leq 1$.
- Solution depends on $x$ only through $\langle x_i, x_j \rangle$
  - may generalize $\langle x_i, x_j \rangle$ to $K(x_i, x_j)$.