

MAC 0460 / 5832

Introduction to Machine Learning

17 — Support Vector Machines (SVM)

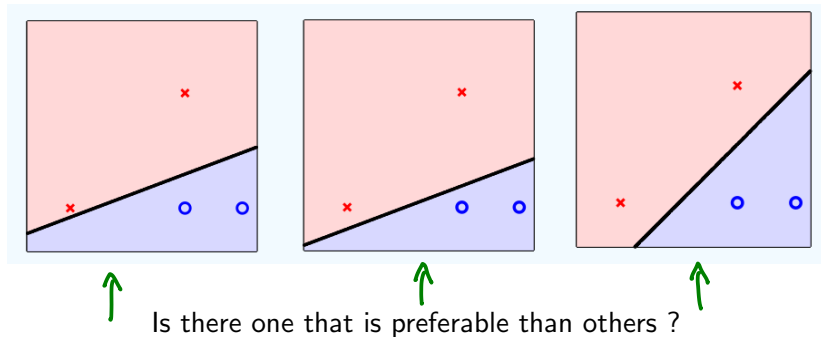
- hyperplane • margin • margin violation •
- QP problems • dual problems • kernel trick •

IME/USP (16/06/2021)

- Binary linear classification: The linearly separable case
- hard-margin SVM: Maximum margin formulation
- Binary linear classification: The non-linearly separable case
- soft-margin SVM: allows margin violation
- hard-margin/soft-margin SVM is a QP problem \Rightarrow
- Dual of hard-margin/soft-margin SVM is also QP \Leftarrow
- How to solve QP problems \Leftarrow
- Non-linear classification: the kernel trick

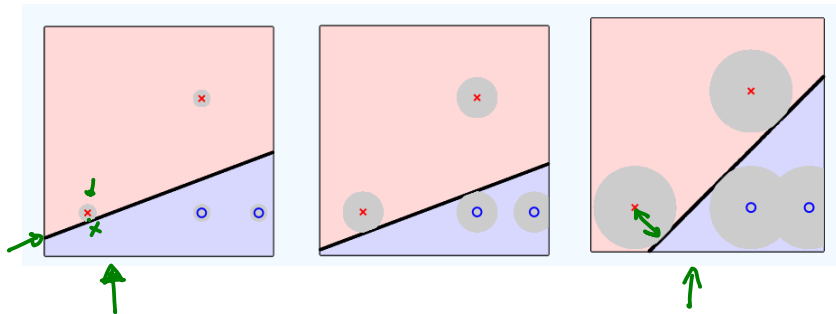
Linearly separable case

Given a linearly separable D , a linear decision boundary separating **negatives** from **positives** can be obtained using, for instance, PLA or logistic regression

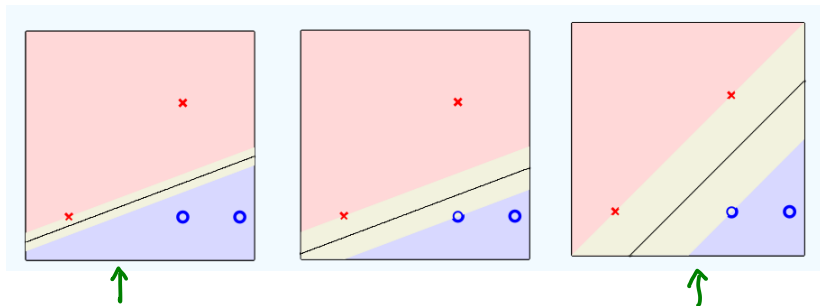


Intuition

Depending on where the separating line is, it is more or less robust to noise



Maximum margin



Any of these lines separate the **negatives** from the **positives**
They have margins of different sizes

How to find the separating hyperplane that maximizes the margin ?

In **SVM**, this is achieved by formulating the problem as a quadratic programming (QP) optimization problem

QP: optimization of quadratic functions with linear constraints on the variables

Notations

Previous Chapters

$$\mathbf{x} \in \{1\} \times \mathbb{R}^d; \mathbf{w} \in \mathbb{R}^{d+1}$$



$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}.$$

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

This Chapter

$$\mathbf{x} \in \mathbb{R}^d; b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d$$

b = bias

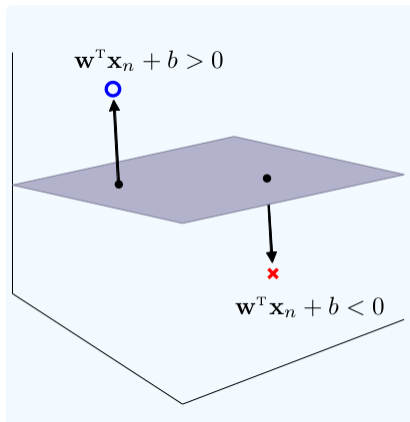
$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}.$$

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$



$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}_n + b = 0 \text{ defines a hyperplane } H$$

Classification based on H



Output class: $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$

Relate parameters to margin

The classifier has parameters (\mathbf{w}, b) :

$$\underline{h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)}$$

We need to somehow relate $\underline{\mathbf{w}}$ and \underline{b} with the margin

Margin is the distance between H and the closest point among all points in D

\Rightarrow Let us examine $\underline{d(\mathbf{x}, H)}$!

Recap: vector normal to the hyperplane

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

The vector \mathbf{w} is \perp to the plane in the \mathcal{X} space:

Take \mathbf{x}' and \mathbf{x}'' on the plane

??



Recap: vector normal to the hyperplane

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

The vector \mathbf{w} is \perp to the plane in the \mathcal{X} space:

Take \mathbf{x}' and \mathbf{x}'' on the plane

$$\mathbf{w}^T \mathbf{x}' + b = 0 \quad \text{and} \quad \mathbf{w}^T \mathbf{x}'' + b = 0$$

$$\implies \mathbf{w}^T (\mathbf{x}' - \mathbf{x}'') = 0$$



Recap: distance between point and hyperplane

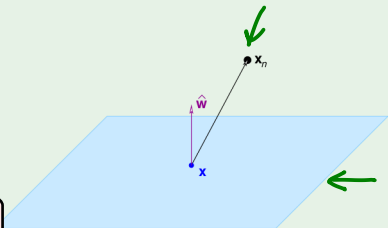
$$d(\mathbf{x}_n, H) = ?$$

Distance between \mathbf{x}_n and the plane:

Take any point \mathbf{x} on the plane

Projection of $\mathbf{x}_n - \mathbf{x}$ on $\hat{\mathbf{w}}$

??



Recap: distance between point and hyperplane

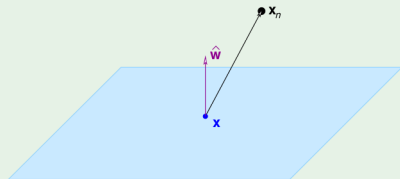
$$d(\mathbf{x}_n, H) = ?$$

Distance between \mathbf{x}_n and the plane:

Take any point \mathbf{x} on the plane

Projection of $\mathbf{x}_n - \mathbf{x}$ on \mathbf{w}

$$\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \implies \text{distance} = \underbrace{|\hat{\mathbf{w}}^\top (\mathbf{x}_n - \mathbf{x})|}_{\substack{\uparrow \\ \text{distance}}}$$



The points in D may be in one or in the other side of H

Thus, distance is given by the absolute value $|\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})|$

The need to treat the two cases (if-else situation) is not convenient

Remember logistic regression? There we used a trick to avoid if-else:


$$\begin{aligned} P(y|\mathbf{x}) &= \theta(y \mathbf{w}^T \mathbf{x}) \\ P(y|\mathbf{x}) &= P(y = 1|\mathbf{x})^y [1 - P(y = 1|\mathbf{x})]^{1-y} \end{aligned}$$

$\frac{p(y|\mathbf{x})}{\mathbf{w}^T \mathbf{x}}$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T (\mathbf{x}_n - \mathbf{x})|$$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$


Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$

$$\mathbf{w}^T(\mathbf{x}_n - \mathbf{x}) = \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}$$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$

$$\mathbf{w}^T(\mathbf{x}_n - \mathbf{x}) = \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_n + b - (\mathbf{w}^T \mathbf{x} + b)$$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$

$$\begin{aligned}\mathbf{w}^T(\mathbf{x}_n - \mathbf{x}) &= \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_n + b - (\mathbf{w}^T \mathbf{x} + b) \\ &= \mathbf{w}^T \mathbf{x}_n + b - 0 = \mathbf{w}^T \mathbf{x}_n + b\end{aligned}$$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} \underbrace{|\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|}$$


$$\begin{aligned} \underbrace{\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})} &= \underbrace{\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}} = \underbrace{\mathbf{w}^T \mathbf{x}_n + b}_{\downarrow} - \underbrace{(\mathbf{w}^T \mathbf{x} + b)}_{\downarrow} \\ &= \underbrace{\mathbf{w}^T \mathbf{x}_n + b}_{\downarrow} - \underbrace{0}_{\downarrow} = \underbrace{\mathbf{w}^T \mathbf{x}_n + b}_{\downarrow} \end{aligned}$$

Why $\mathbf{w}^T \mathbf{x} + b = 0$?

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$

$$\begin{aligned}\mathbf{w}^T(\mathbf{x}_n - \mathbf{x}) &= \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_n + b - (\mathbf{w}^T \mathbf{x} + b) \\ &= \mathbf{w}^T \mathbf{x}_n + b - 0 = \mathbf{w}^T \mathbf{x}_n + b\end{aligned}$$

$$\text{dist}(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x}_n + b| = \frac{1}{\|\mathbf{w}\|} y_n (\mathbf{w}^T \mathbf{x}_n + b)$$


Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T (\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T (\mathbf{x}_n - \mathbf{x})|$$

\mathbf{w}^T

$= \mathbf{w}^T$

Why I can do $|\mathbf{w}^T \mathbf{x}_n + b| = y_n (\mathbf{w}^T \mathbf{x}_n + b)$?

$$\text{dist}(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x}_n + b| = \frac{1}{\|\mathbf{w}\|} y_n (\mathbf{w}^T \mathbf{x}_n + b)$$

Rewriting the distance between point and hyperplane

$$\text{dist}(\mathbf{x}_n, H) = |\hat{\mathbf{w}}^T(\mathbf{x}_n - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T(\mathbf{x}_n - \mathbf{x})|$$

$$\begin{aligned}\mathbf{w}^T(\mathbf{x}_n - \mathbf{x}) &= \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_n + b - (\mathbf{w}^T \mathbf{x} + b) \\ &= \mathbf{w}^T \mathbf{x}_n + b - 0 = \mathbf{w}^T \mathbf{x}_n + b\end{aligned}$$

$$\text{dist}(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x}_n + b| = \frac{1}{\|\mathbf{w}\|} \underbrace{y_n(\mathbf{w}^T \mathbf{x}_n + b)}$$

(because if \mathbf{x}_n is at the correct side $\implies \underbrace{y_n(\mathbf{w}^T \mathbf{x}_n)} > 0$)

Choosing a convenient hyperplane representation (weights)

Distance as seen before:

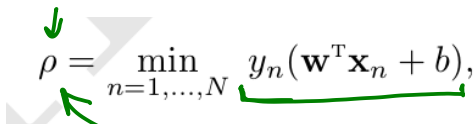
$$\underbrace{dist(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x}_n + b|}_{\text{distance}} = \underbrace{\frac{1}{\|\mathbf{w}\|} y_n (\mathbf{w}^T \mathbf{x}_n + b)}_{\text{margin}} \quad \leftarrow$$

If I manage to make $|\mathbf{w}^T \mathbf{x}_n + b| = 1$, then I will have

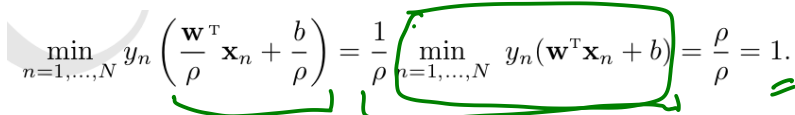
$$dist(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|}$$

We can always rescale (\mathbf{w}, b) so as to have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$.

Let us do that with respect to the closest point to the hyperplane:


$$\rho = \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b),$$

If we divide (\mathbf{w}, b) by ρ , the hyperplane does not change:



$$\min_{n=1, \dots, N} y_n \left(\frac{\mathbf{w}^T}{\rho} \mathbf{x}_n + \frac{b}{\rho} \right) = \frac{1}{\rho} \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = \frac{\rho}{\rho} = 1.$$

Exercise 8.2

Consider the data below and a 'hyperplane' (b, \mathbf{w}) that separates the data.

$$X = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} -1 \\ -1 \\ +1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 1.2 \\ -3.2 \end{bmatrix} \quad b = -0.5$$

- (a) Compute $\rho = \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b)$.
- (b) Compute the weights $\frac{1}{\rho}(b, \mathbf{w})$ and show that they satisfy (8.2).
- (c) Plot both hyperplanes to show that they are the *same* separator.

$$(8.2) \quad \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$$


Wrapping up

Let D be a linearly separable set of points, $\mathbf{x}_n \in D$, and a separating hyperplane H characterized by (\mathbf{w}, b) . Then

$$\text{dist}(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} y_n(\mathbf{w}^T \mathbf{x}_n + b) \quad \swarrow$$

We can always choose (\mathbf{w}, b) such that the closest point \mathbf{x}_n to H satisfies

$$\underline{y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1}$$

In such case

$$\text{dist}(\mathbf{x}_n, H) = \frac{1}{\|\mathbf{w}\|} \quad \leftarrow$$

The problem we want to solve

$$\begin{array}{ll}\text{maximize}_{\mathbf{w}, b} & \frac{1}{\|\mathbf{w}\|} \quad \leftarrow \\ \text{subject to} & \min_{i=1, \dots, N} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \quad \leftarrow\end{array}$$

The problem we want to solve



$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{maximize}} && \frac{1}{\|\mathbf{w}\|} \\ & \text{subject to} && \min_{i=1, \dots, N} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \quad \leftarrow \end{aligned}$$

- The constraint $\min_{i=1, \dots, N} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$ implies $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$ which has the effect of forcing all examples to be classified correctly
- The equality $\min_{i=1, \dots, N} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$ implies that the distance of the closest point to the hyperplane is $\frac{1}{\|\mathbf{w}\|}$ (a nice objective function!)

The problem we want to solve

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{maximize}} && \frac{1}{\|\mathbf{w}\|} \\ & \text{subject to} && \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 \end{aligned}$$

Equivalent formulation


$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} && \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 \end{aligned}$$

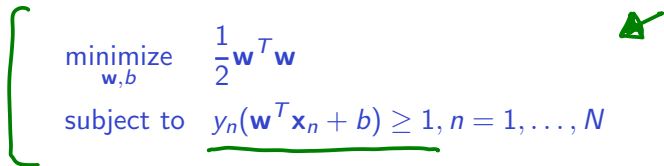
Quadratic function

Relaxed formulation

Original minimization formulation:

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} && \min_{n=1, \dots, N} y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1 \end{aligned}$$

Equivalent relaxed formulation:


$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} && \underline{y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, n = 1, \dots, N} \end{aligned}$$

The equivalence can be proved by contradiction (see Chapter on SVM, page 7)

A toy example

$N=4$

Constraints: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$? 

$$X = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$$

A toy example

$$X = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad y = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$$

Constraints: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$?



$$\begin{aligned} -b &\geq 1 \quad (1) \\ -(2w_1 + 2w_2 + b) &\geq 1 \quad (2) \\ 2w_1 + b &\geq 1 \quad (3) \\ 3w_1 + b &\geq 1 \quad (4) \end{aligned}$$

$$\frac{1}{2} \mathbf{w}^T \mathbf{w}$$

Solving it by hand

$$-b \geq 1 \quad (1)$$

$$-(2w_1 + 2w_2 + b) \geq 1 \quad (2)$$

$$2w_1 + b \geq 1 \quad (3)$$

$$3w_1 + b \geq 1 \quad (4)$$

Solving it by hand

$$-b \geq 1 \quad (1)$$

$$-(2w_1 + 2w_2 + b) \geq 1 \quad (2)$$

$$\underline{2w_1 + b \geq 1} \quad (3)$$

$$\underline{3w_1 + b \geq 1} \quad (4)$$

- From (3) and (1)

$$\underline{2w_1 + b \geq 1} \rightsquigarrow 2w_1 \geq 1 - b \rightsquigarrow \underline{w_1 \geq \frac{1}{2}(1 - b)} \ \&\& \ \underline{b \leq -1}$$
$$\Rightarrow \underline{\underline{w_1 \geq 1}}$$

Solving it by hand

$$-b \geq 1 \quad (1)$$

$$-(2w_1 + 2w_2 + b) \geq 1 \quad (2)$$

$$2w_1 + b \geq 1 \quad (3) \quad \checkmark$$

$$3w_1 + b \geq 1 \quad (4)$$

- From (3) and (1)

$$2w_1 + b \geq 1 \rightsquigarrow 2w_1 \geq 1 - b \rightsquigarrow w_1 \geq \frac{1}{2}(1 - b) \text{ \&\& } b \leq -1$$

$$\implies w_1 \geq 1$$

- From (2) and (3):

$$-(2w_1 + 2w_2 + b) \geq 1 \rightsquigarrow -2w_1 - 2w_2 - b \geq 1 \rightsquigarrow$$

$$2w_2 \leq -2w_1 - b - 1 \text{ \&\& } 2w_1 + b \geq 1 \implies w_2 \leq -1$$

Solving it by hand

$$-b \geq 1 \quad (1)$$

$$-(2w_1 + 2w_2 + b) \geq 1 \quad (2)$$

$$2w_1 + b \geq 1 \quad (3)$$

$$3w_1 + b \geq 1 \quad (4)$$

- From (3) and (1)

$$2w_1 + b \geq 1 \rightsquigarrow 2w_1 \geq 1 - b \rightsquigarrow w_1 \geq \frac{1}{2}(1 - b) \text{ \&\& } b \leq -1 \\ \implies w_1 \geq 1$$

- From (2) and (3):

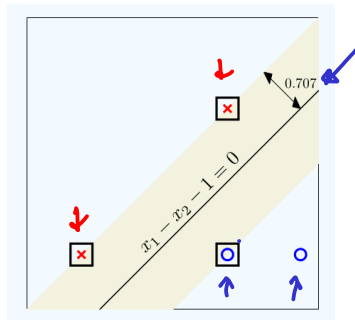
$$-(2w_1 + 2w_2 + b) \geq 1 \rightsquigarrow -2w_1 - 2w_2 - b \geq 1 \rightsquigarrow \\ 2w_2 \leq -2w_1 - b - 1 \text{ \&\& } 2w_1 + b \geq 1 \implies w_2 \leq -1$$

Thus, $\frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2}(w_1^2 + w_2^2) \geq 1$ and the minimum is at $\mathbf{w} = (1, -1)$;
 $(b = -1, w_1 = 1, w_2 = -1)$ satisfies the 4 constraints

Solution (by hand) of the toy example

The separating hyperplane H with maximum margin is given by $x_1 - x_2 - 1 = 0$.

$$X = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad y = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$$



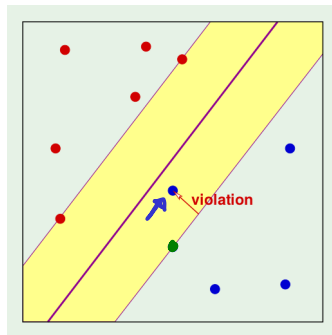
The margin is $\frac{1}{\|w\|} = \frac{1}{\sqrt{2}} \approx 0.707$

Summary (linearly separable case)

- The goal is to find a hyperplane that maximizes the margin
- We examined the formulation of the hard margin SVM
- It can be written as a QP optimization (quadratic objective function with linear inequality constraints)
- We solved a toy example by hand
- We still do not know how to solve QP problems

Non-linearly separable case

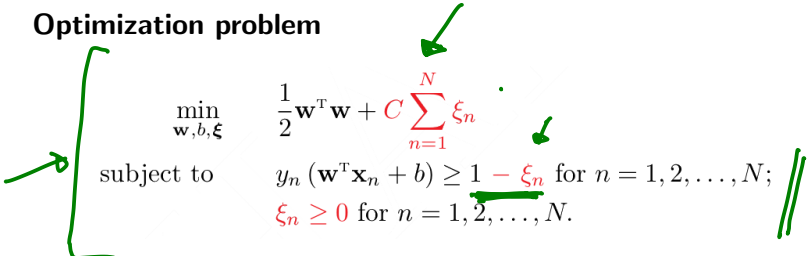
This case is dealt by considering a **soft margin** formulation as opposed to the (previous) **hard margin** formulation:



Soft margin: $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq \underline{1 - \xi_n}$ ✓

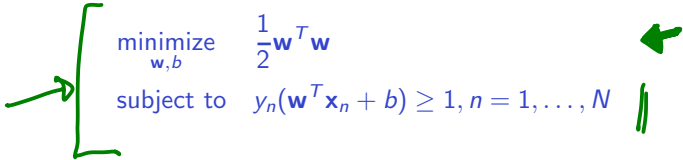
(Hard margin: $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq \underline{1}$) =

Optimization problem

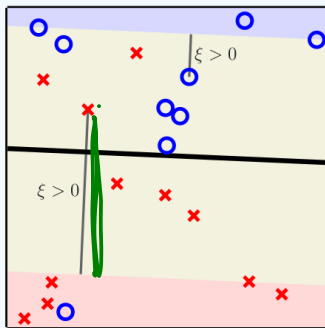

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \text{ for } n = 1, 2, \dots, N; \\ & \xi_n \geq 0 \text{ for } n = 1, 2, \dots, N. \end{aligned}$$

$C \geq 0$ is an user-specified parameter; the larger it is, the smaller the allowed margin violation

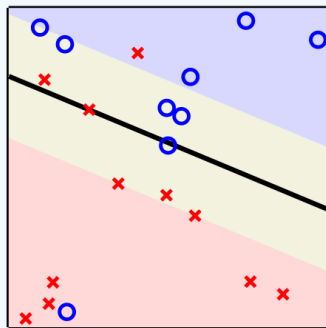
Compare to the hard-margin formulation:


$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, n = 1, \dots, N \end{aligned}$$

Intuition on constant C



(a) $C = 1$



(b) $C = 500$



$$\underset{\mathbf{w}, b, \xi}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i$$

How to solve QP optimization problems ?

Both cases, hard and soft margin SVM, can be formulated as a QP optimization problem

Primal formulation: Standard QP optimization

Dual formulation: based on Lagrange formulation, dual QP

Standard QP optimization

Standard form of QP problems

M inequality constraints and Q positive semi-definite

$$\begin{array}{ll} \underset{\mathbf{u}}{\text{minimize}} & \frac{1}{2}\mathbf{u}^T Q \mathbf{u} + \mathbf{p}^T \mathbf{u} \\ \text{subject to:} & \mathbf{a}_m^T \mathbf{u} \geq c_m \quad (m = 1, \dots, M) \end{array}$$

In matrix form

$$\begin{array}{ll} \underset{\mathbf{u}}{\text{minimize}} & \frac{1}{2}\mathbf{u}^T Q \mathbf{u} + \mathbf{p}^T \mathbf{u} \\ \text{subject to:} & \mathbf{A} \mathbf{u} \geq \mathbf{c} \end{array}$$

QP solvers can be used to compute the optimal solution \mathbf{u}^* :

$$\mathbf{u}^* \leftarrow \text{QP}(Q, \mathbf{p}, \mathbf{A}, \mathbf{c})$$

SVM – standard QP formulation

QP problem formulation

$$\begin{aligned} & \underset{\mathbf{u}}{\text{minimize}} && \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \underline{\mathbf{p}}^T \mathbf{u} \\ & \text{subject to:} && \mathbf{a}_m^T \mathbf{u} \geq c_m \\ & && i = m, \dots, M \end{aligned}$$

QP of hard-margin SVM ↗

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad \text{---} \\ & \text{subject to:} && y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \uparrow \\ & && i = 1, \dots, N \end{aligned}$$

Denoting $\mathbf{u} = \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$, we have

$$\begin{aligned} \mathbf{w}^T \mathbf{w} &= \begin{bmatrix} b & \mathbf{w}^T \end{bmatrix} \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix} \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} = \mathbf{u}^T \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix} \mathbf{u}, \\ \underline{\mathbf{a}}_n^T &= y_n \begin{bmatrix} 1 & \mathbf{x}_n^T \end{bmatrix} \quad \text{and} \quad c_n = 1 \end{aligned}$$

Linear Hard-Margin SVM with QP

- 1: Let $\mathbf{p} = \mathbf{0}_{d+1}$ ($(d+1)$ -dimensional zero vector) and $\mathbf{c} = \mathbf{1}_N$ (N -dimensional vector of ones). Construct matrices Q and A , where

$$Q = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & I_d \end{bmatrix}, \quad A = \underbrace{\begin{bmatrix} y_1 & -y_1 \mathbf{x}_1^T \\ \vdots & \vdots \\ y_N & -y_N \mathbf{x}_N^T \end{bmatrix}}_{\text{signed data matrix}}.$$

Handwritten annotations:
 - Green arrow points to the top-left element of Q .
 - Green arrow points to the top-right block of Q ($\mathbf{0}_d^T$), labeled $1 \times d$.
 - Green arrow points to the bottom-left block of Q ($\mathbf{0}_d$), labeled $d \times 1$.
 - Green arrow points to the top-right element of A .
 - Green arrow points to the bottom-right element of A .
 - Green arrow points to the entire A matrix, labeled "signed data matrix".

- 2: Calculate $\begin{bmatrix} b^* \\ \mathbf{w}^* \end{bmatrix} = \mathbf{u}^* \leftarrow \text{QP}(Q, \mathbf{p}, A, \mathbf{c})$.

- 3: Return the hypothesis $g(\mathbf{x}) = \text{sign}(\mathbf{w}^{*T} \mathbf{x} + b^*)$.

$$W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$v = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

b

$$Q = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$v^T Q v = [b \ w_1 \ w_2] \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

$$= [0 \ w_1 \ w_2] \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

$$= W^T W$$