

CSI 436/536 (Fall 2024)

# Machine Learning

Lecture 11: Support Vector Machines

Chong Liu

Assistant Professor of Computer Science

Oct 3, 2024

# Announcement

- Homework 1 all submitted on time. Good job!
- Homework 2 has been released.

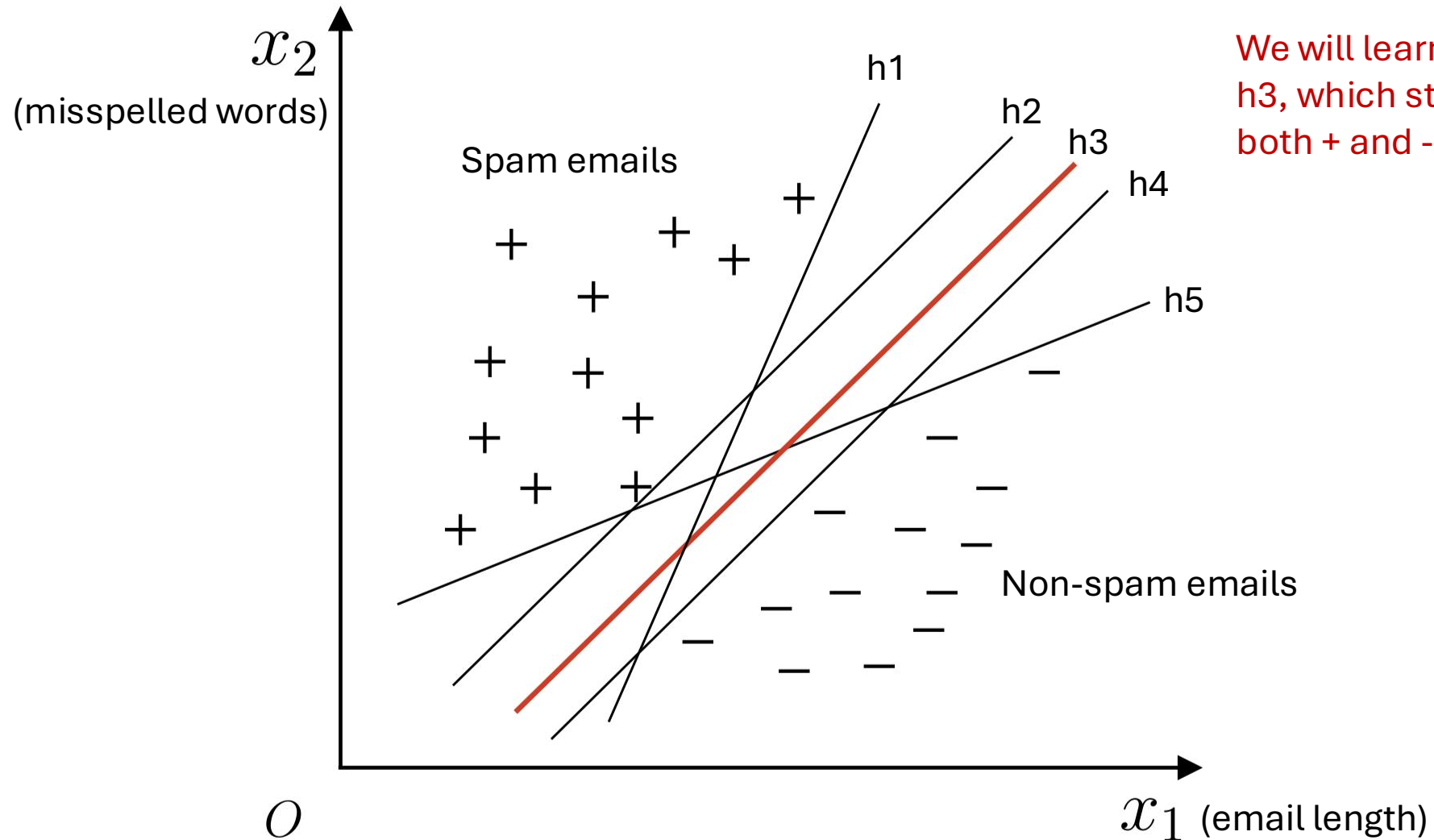
# Recap: Regularization

- Linear regression
  - Solving the Least Square problem {with GD, SGD and direct solver}
- Regularization
  - Controls the parameter complexity of the fitted function
  - Prevents overfitting!
  - Different regularization: L-2 (most popular) and L-1
- Case study: Predict House Price
  - Effect of regularization on training test and test error
  - Regularization path (Effect of regularization on coefficients)

# Today

- Move back to binary classification problem
  - Spam email / non-spam email
- Margin
- Support Vector Machines
- **Warning: While without any proof, today's lecture will be very technical. Feel free to interrupt me at any point to ask questions.**

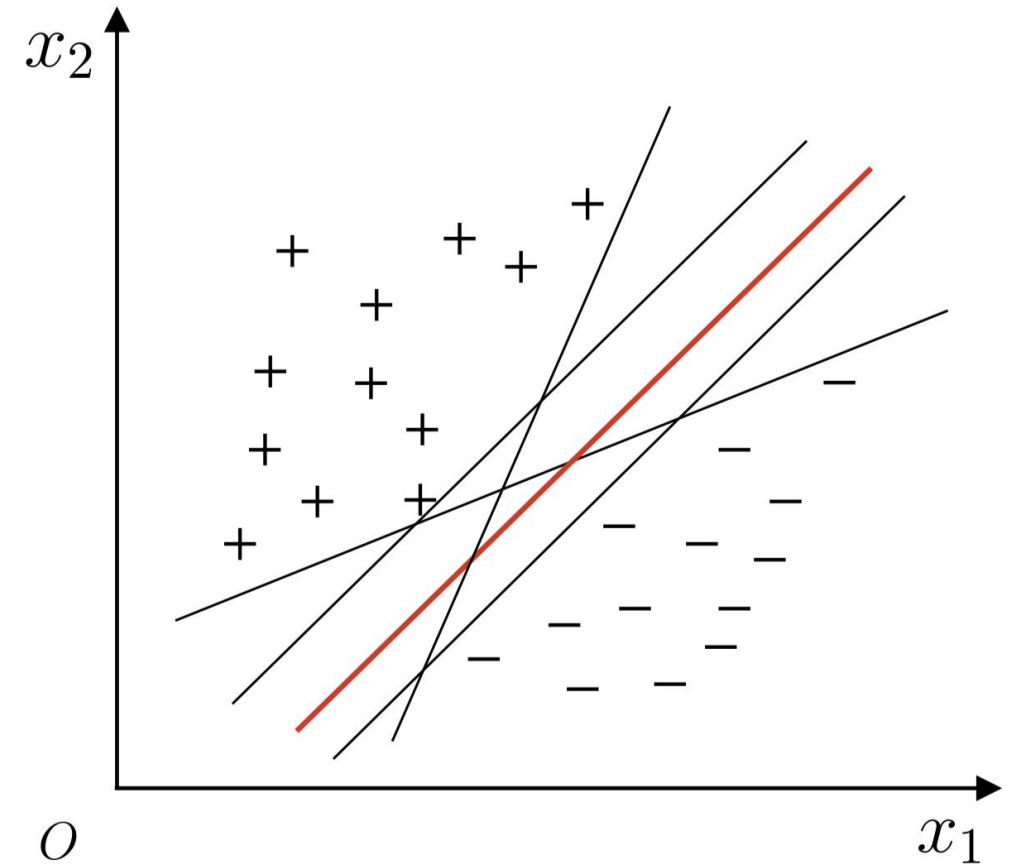
# Discussion: which is the best classifier?



We will learn how to train  $h_3$ , which stays away from both + and -.

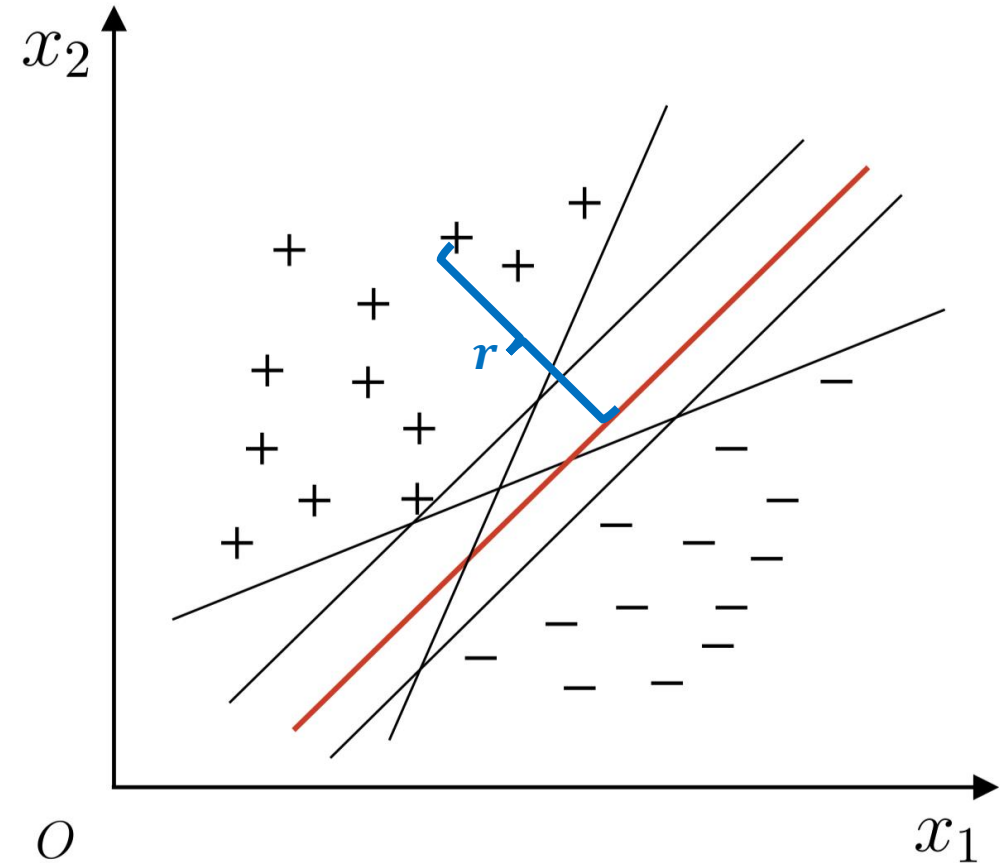
# Linear classification

- Input:  $x = [x_1, x_2] \in R^2$
- Output:  $y \in \{1, -1\}$
- Data:  $n$  data points
- Decision line:
  - $w^T x + b = 0$
  - $w \in R^2, b \in R$  are parameters
  - In-class exercise: Rewrite  $x_2 = x_1 - 5$  in  $w^T x + b = 0$  form.



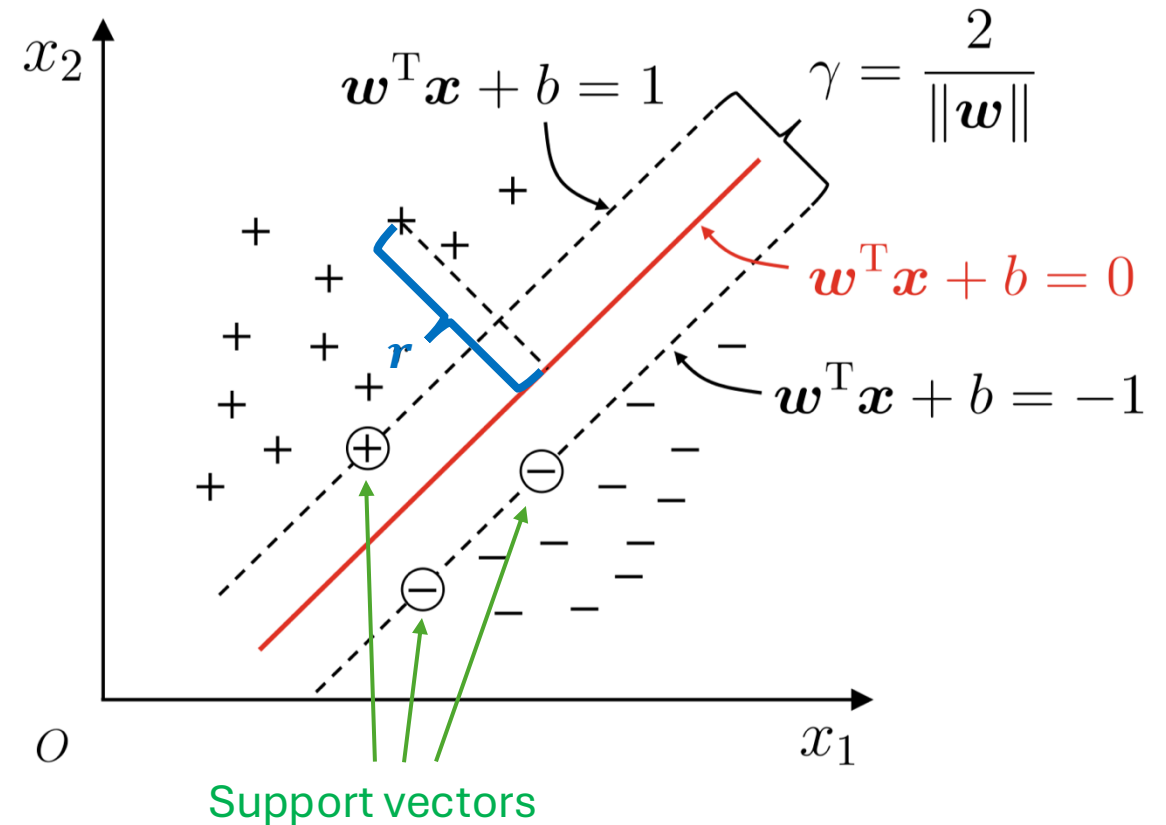
# Margin: min distance of data point to line

- Any data point:
  - $x \in R^2$
- Any line:
  - $w^T x + b = 0$
- Margin:
  - $r = \frac{|w^T x + b|}{\|w\|}$
- **Red** line: We want to learn a max-margin classifier!



# Max-margin classifier

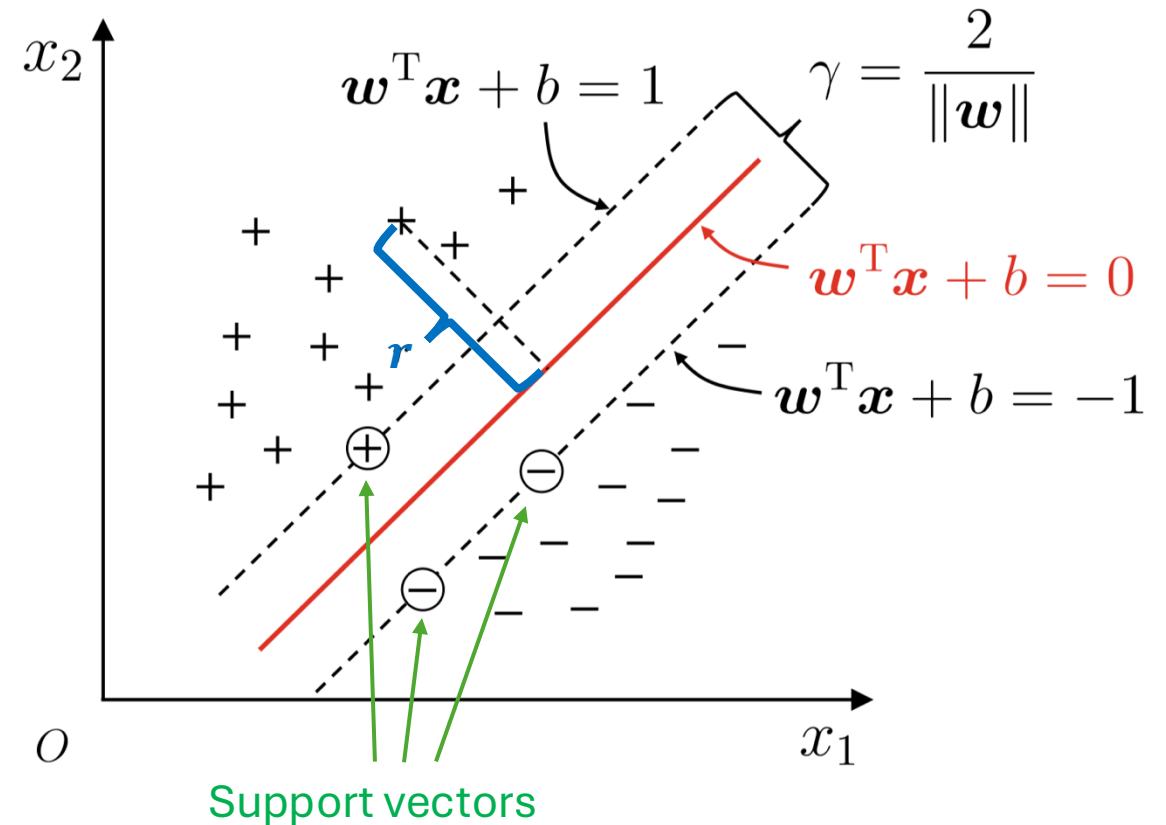
- Discussion: by maximizing margin, which data points are important?
- **Support vectors:**
  - Data points closest to **red** line.
  - Only support vectors affect the training process.
  - Support vector machines (SVM) == max-margin classifier





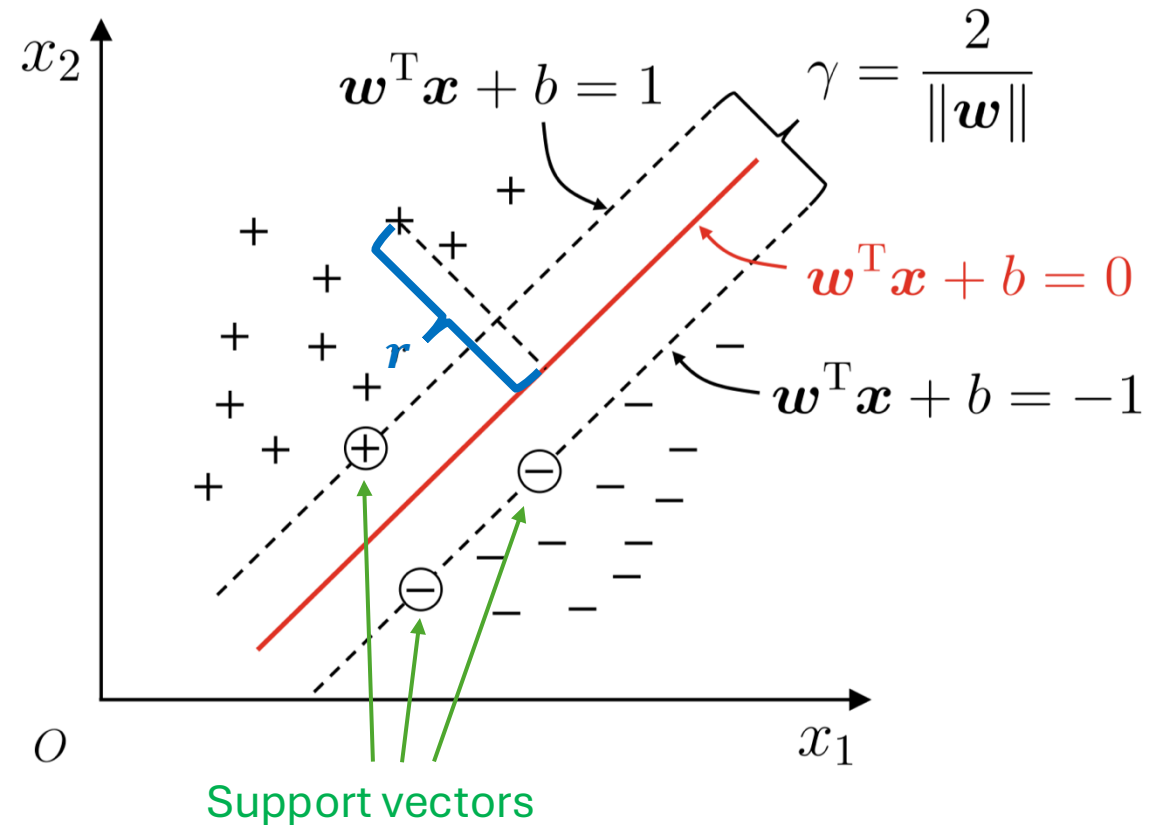
# How to train max-margin classifier?

- Assumption:
  - Linearly separable data points
- Recap: Linear classifier
  - If  $y = 1, w^T x + b > 0$
  - If  $y = -1, w^T x + b < 0$
- Key idea of SVM:
  - If  $y = 1, w^T x + b \geq 1$
  - If  $y = -1, w^T x + b \leq -1$
  - Why? Support vectors are only data points that matter.



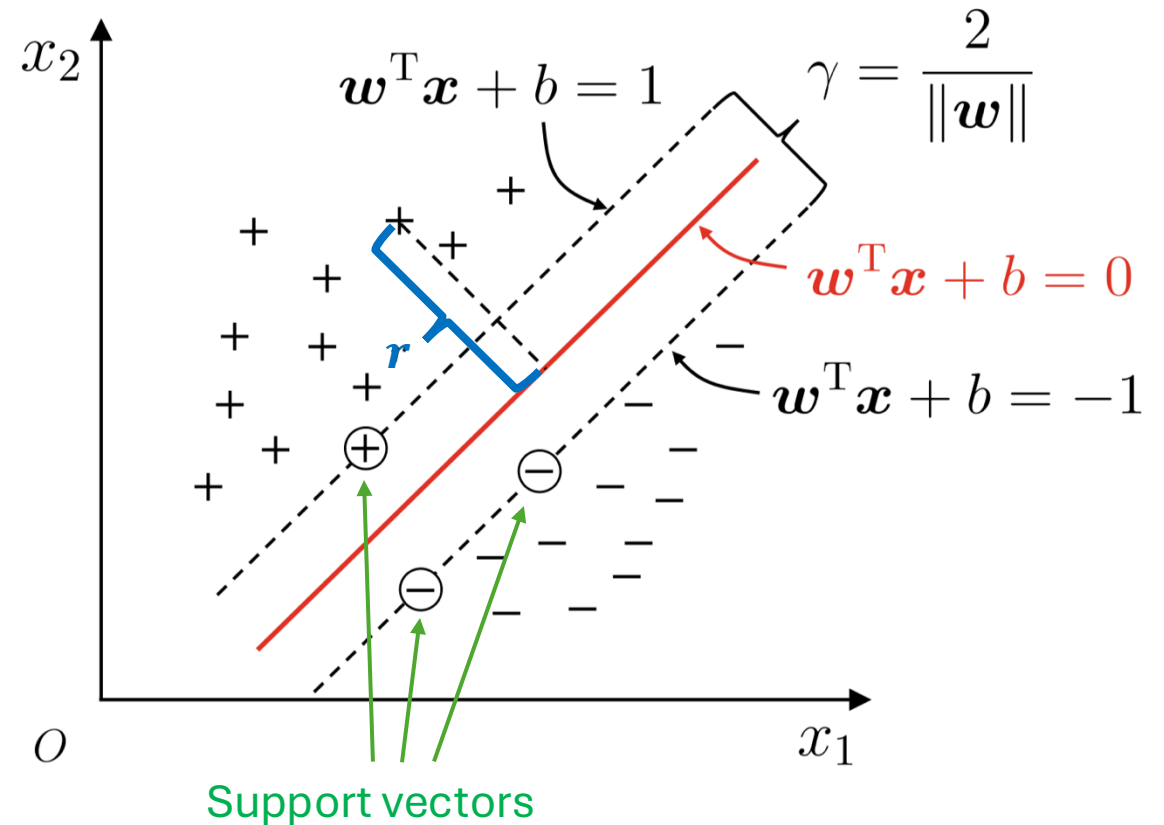
# How to train max-margin classifier?

- Key idea of SVM:
  - If  $y = 1, w^T x + b \geq 1$
  - If  $y = -1, w^T x + b \leq -1$
- Recap: Margin for any data point  $x$ 
  - $r = \frac{|w^T x + b|}{\|w\|}$
- Total margin between support vectors:
  - $\gamma = \frac{2}{\|w\|}$



# How to train max-margin classifier?

- Key idea of SVM:
  - If  $y = 1, w^T x + b \geq 1$
  - If  $y = -1, w^T x + b \leq -1$
- Total margin between support vectors:
  - $\gamma = \frac{2}{\|w\|}$
- **Optimization problem of SVM:**
  - $\max_{w,b} \frac{2}{\|w\|}$
  - s. t.  $y_i(w^T x_i + b) \geq 1, i = 1, \dots, n$



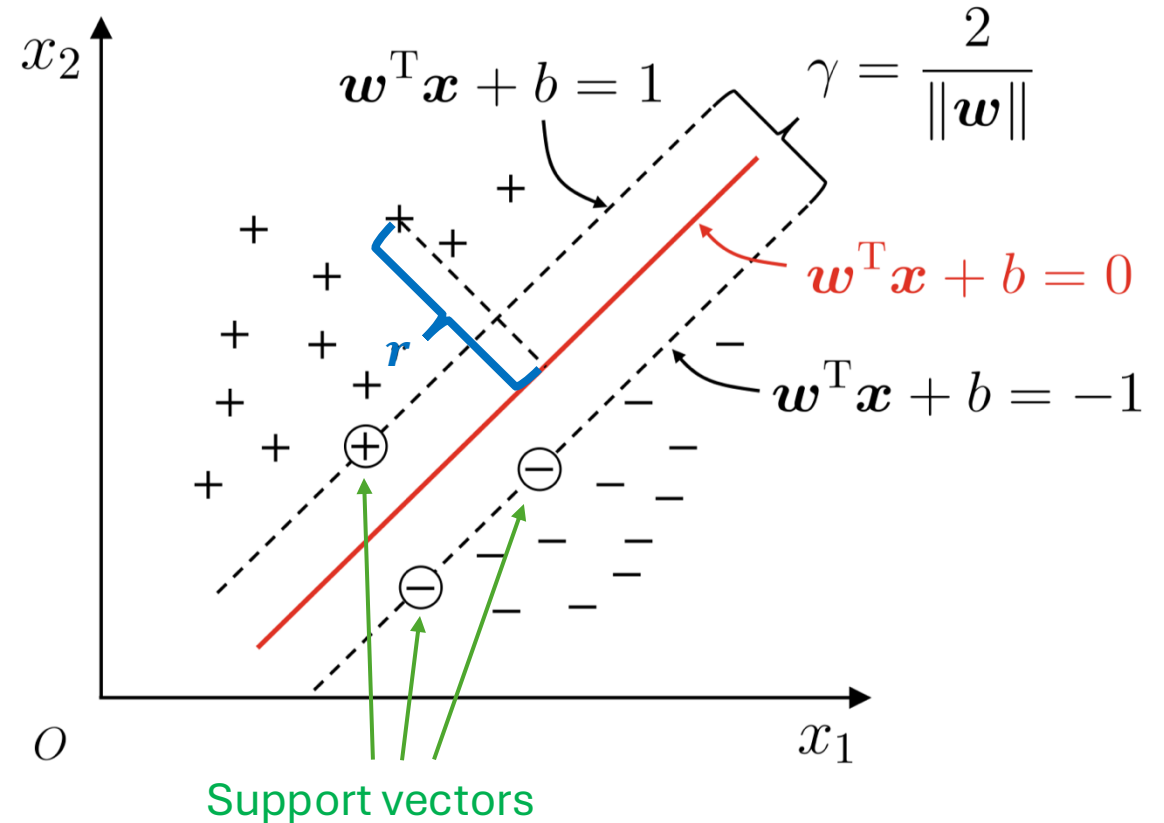
# How to train max-margin classifier?

- Optimization problem of SVM:

- $\max_{w,b} \frac{2}{\|w\|}$
- s. t.  $y_i(w^T x_i + b) \geq 1, i = 1, \dots, n$

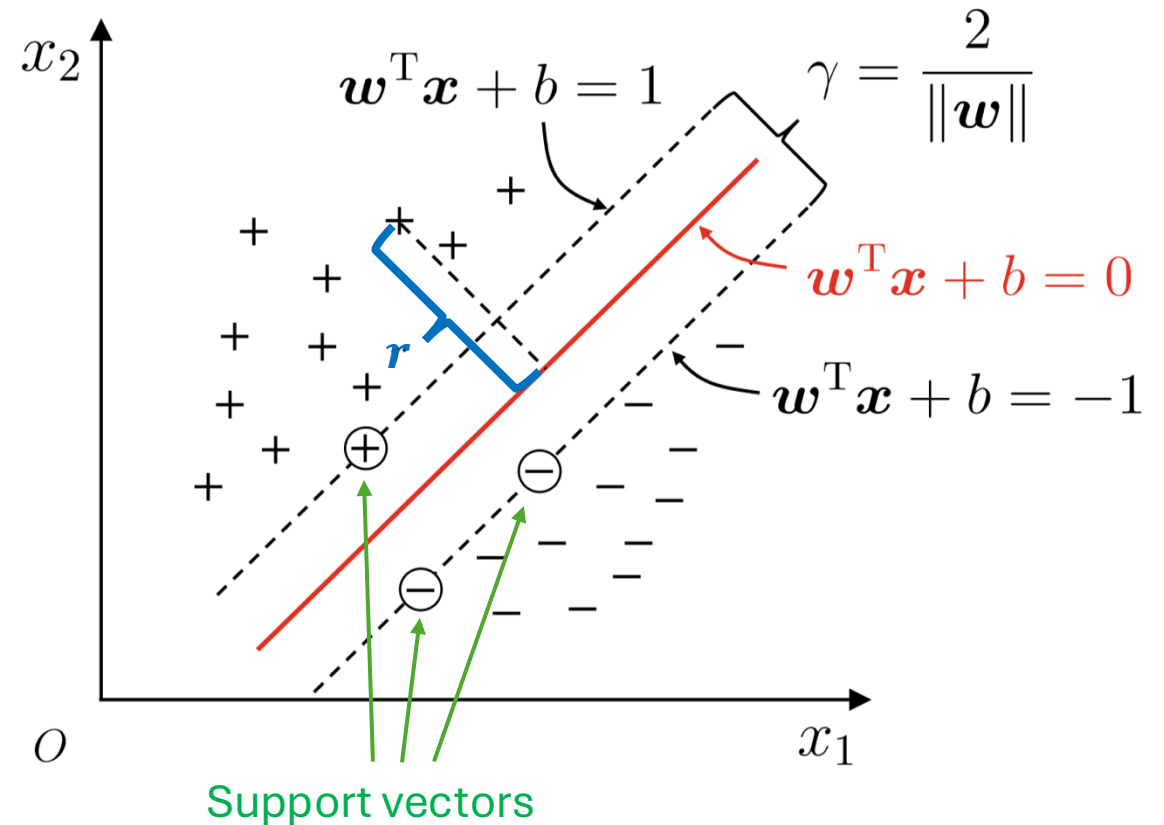
- Equivalent optimization problem:

- $\min_{w,b} \frac{1}{2} \|w\|$
- s. t.  $y_i(w^T x_i + b) \geq 1, i = 1, \dots, n$
- Quadratic programming problem
  - Can be solved using some optimization tools, e.g., CPLEX.



# How to train max-margin classifier?

- Equivalent optimization problem:
  - $\min_{w,b} \frac{1}{2} \|w\|^2$
  - s. t.  $y_i(w^T x_i + b) \geq 1, i = 1, \dots, n$
- In-class exercise:
  - Write the optimization problem with three support vectors:  $(6, 2) +, (7, 1) -, (8, 2) -$ .



# Take a deeper look at optimization problem

- Equivalent optimization problem:

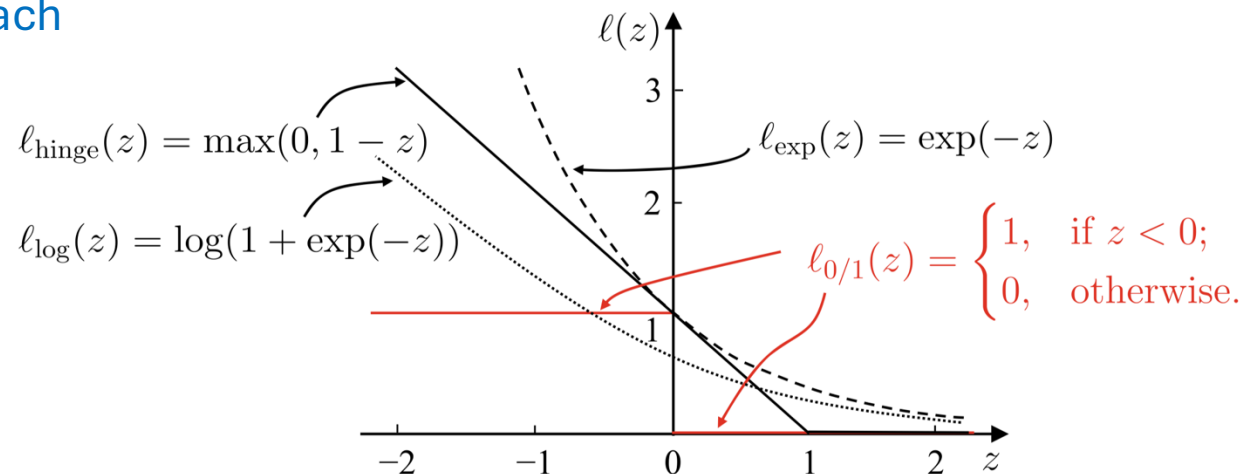
- $\min_{w,b} \frac{1}{2} \|w\|^2$
- s. t.  $y_i(w^T x_i + b) \geq 1, i = 1, \dots, n$

- In-class exercise: Write the Lagrange function

- $\min_{w,b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \lambda_i (1 - y_i(w^T x_i + b))$  Related to a New surrogate loss function - Hinge loss!

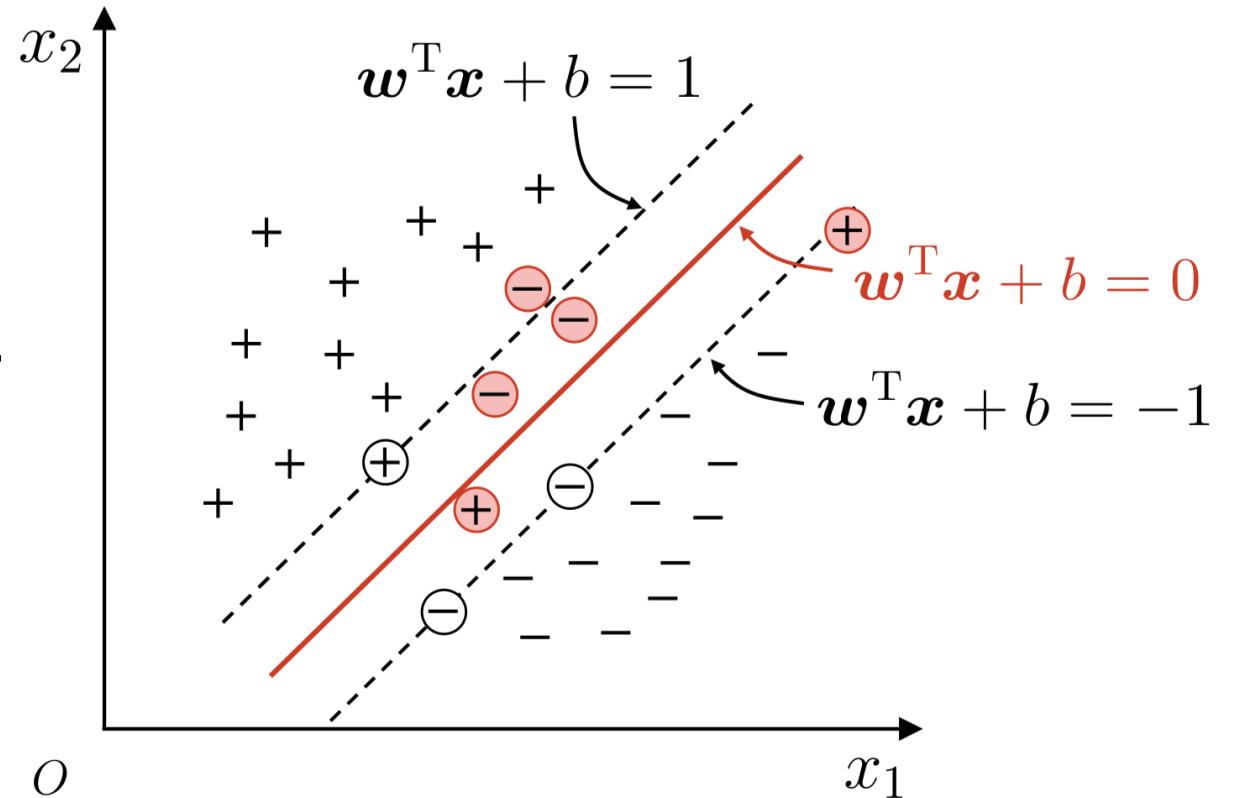
L-2 regularization on parameter!

Weight of each data point



# What if data points are not linearly separable?

- $y_i(w^T x_i + b) \geq 1$  is violated.
- What can we do?
  - Key idea: we give some tolerance.
- New constraint:
  - $y_i(w^T x_i + b) \geq 1 - \xi$
  - $\xi > 0$
- Discussion: what happens when  $\xi$  is very large / small?



# Soft-Margin Support Vector Machines

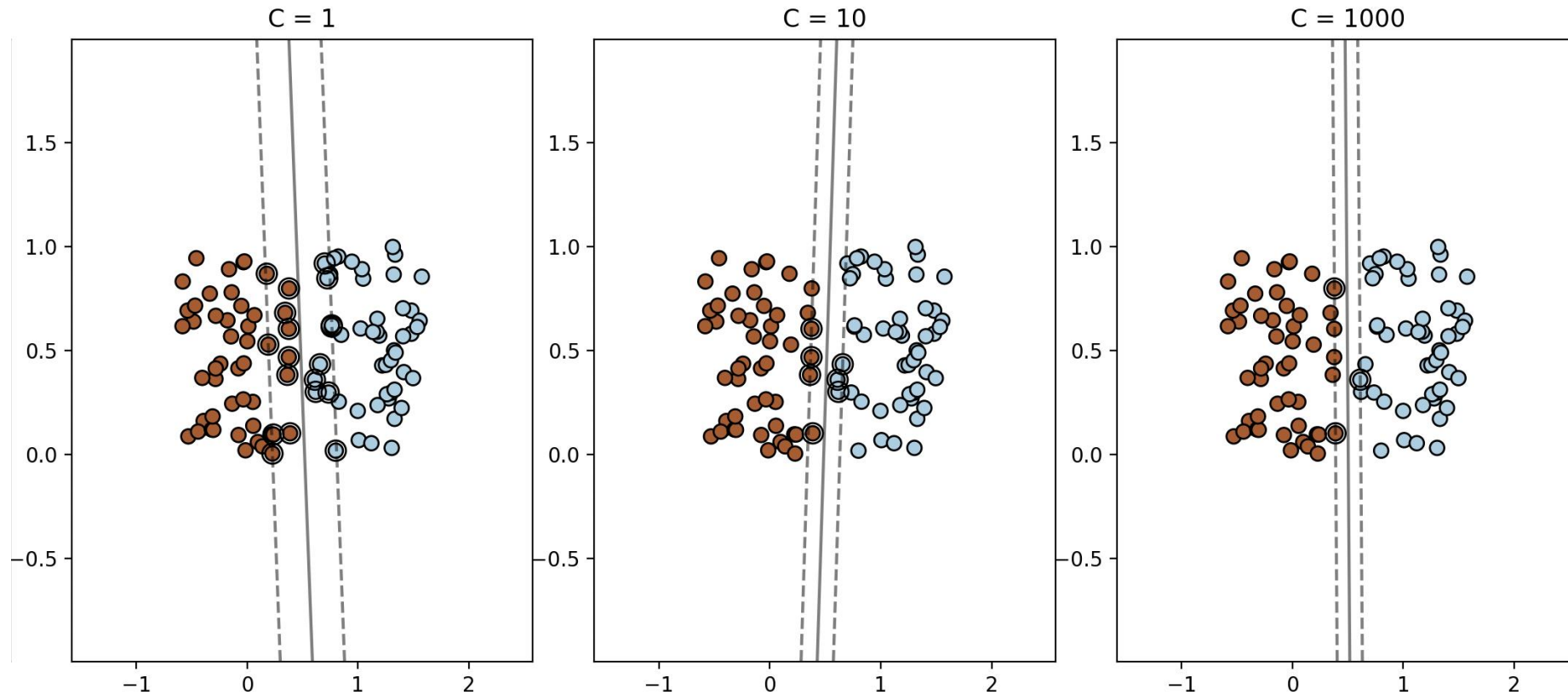
$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i \in [n] \\ & \xi_i \geq 0 \quad \forall i \in [n] \end{aligned}$$

Equivalent to minimizing **Hinge losses**:

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max [1 - y_i(\mathbf{w}^T \mathbf{x} + b), 0]$$



# Hyperparameter C in **soft-margin SVM** and how they affect the margins and “support vectors”.



As C increases, smaller tolerance and fewer soft-margin support vectors.

# Checkpoint of Lecture 1-11

- Tasks of ML:
  - Classification (spam / non-spam email) and regression (house price)
- Philosophy of designing ML algorithms:
  - Regularization: Control the complexity of parameters
    - Prevent overfitting
    - Fun fact: L-2 regularization is associated with max margin classifier
  - Optimization: Toolbox of ML
    - ML problem => optimization problem
      - Direct solver, GD, SGD, and much more!
    - Minimize the loss / parameter complexity
    - Maximize the margin