

Advanced Machine Learning and Deep Learning

Dr. Mohammad Pourhomayoun

Assistant Professor

Computer Science Department

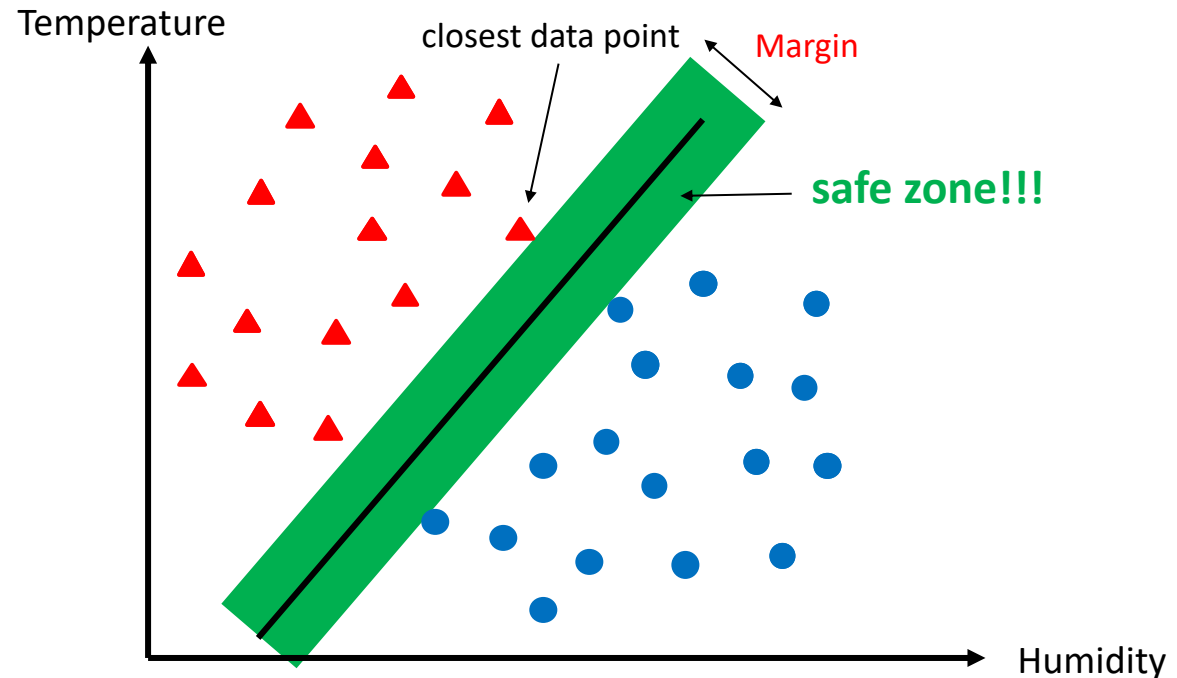
California State University, Los Angeles



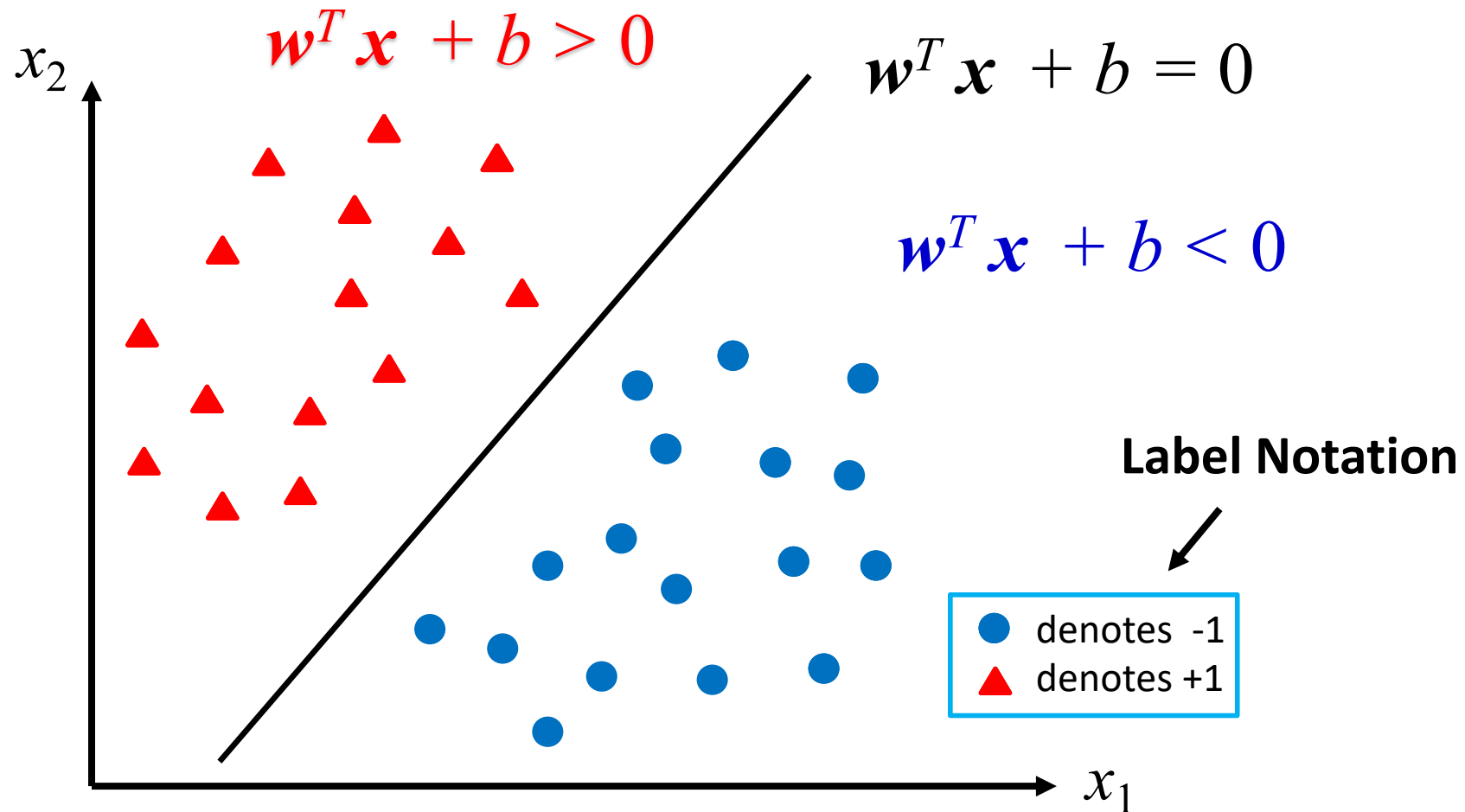
Support Vector Machine (SVM)

Review

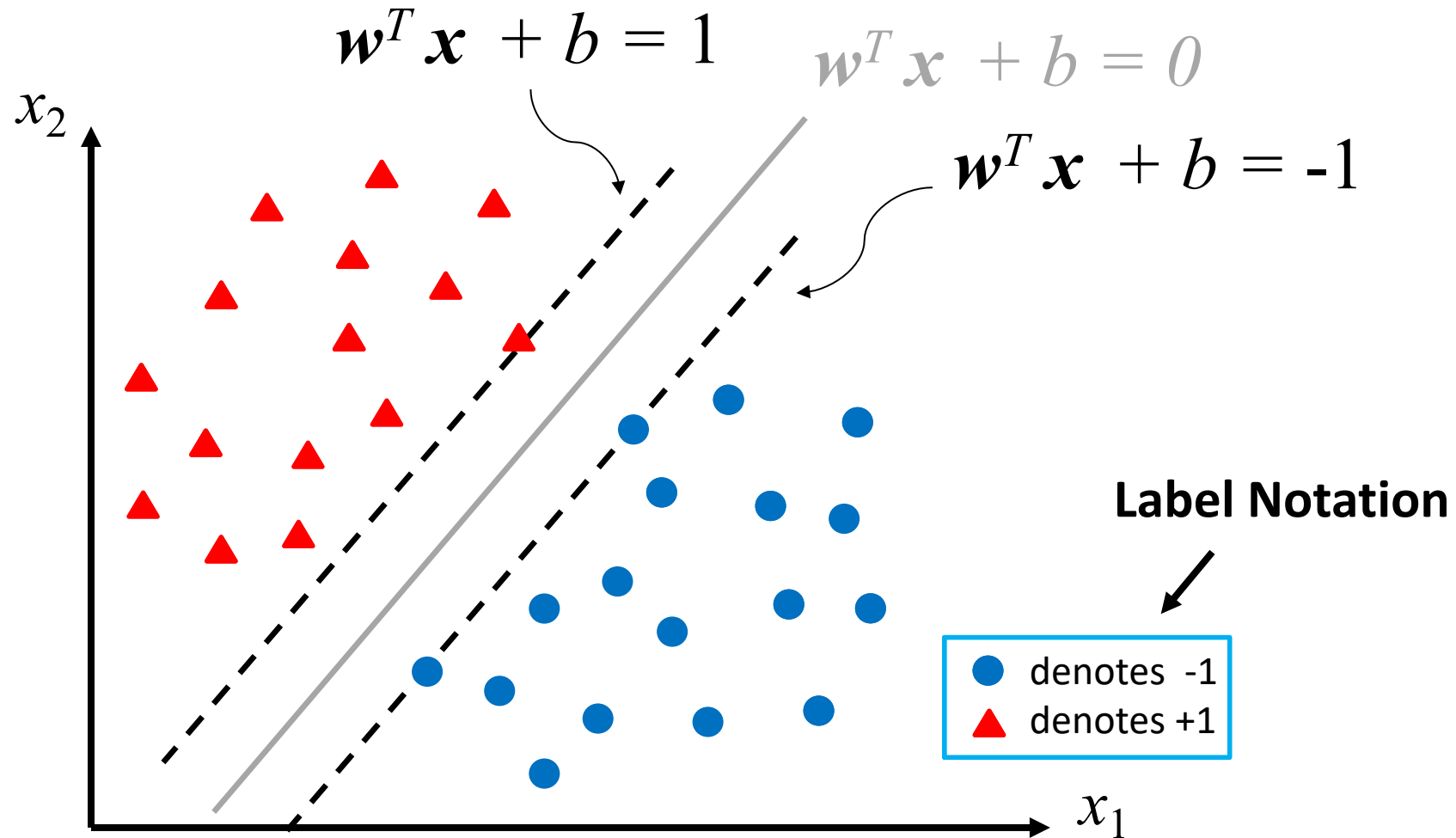
- The best classifier is the one with the **Maximum Margin**.
- **Margin** is widest boundary area before hitting a data point.



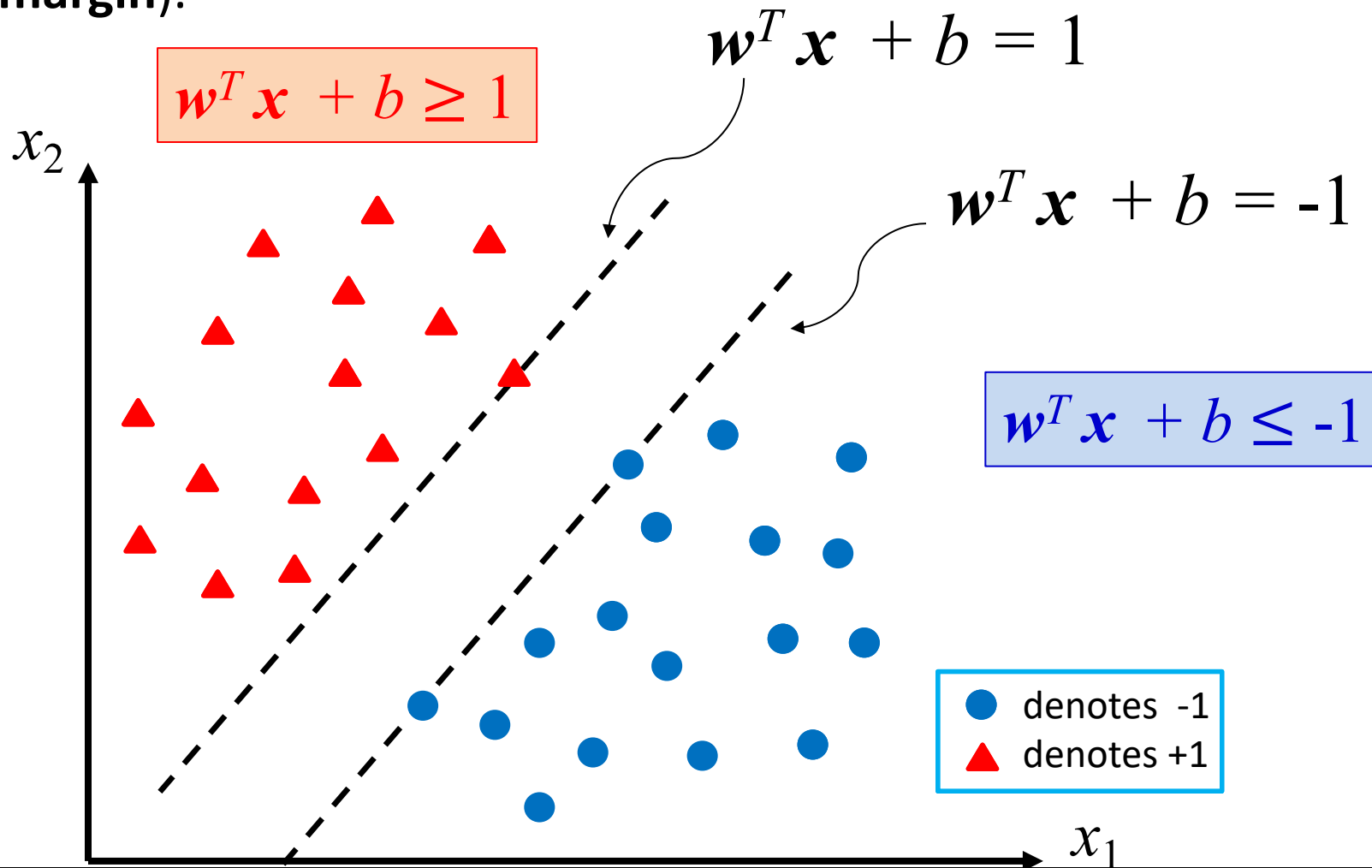
- Given a set of training data points, each point will be considered as a ***m-dimensional vector*** in space. We are looking for ***a line whose value is positive for red samples, and negative for blue samples.***

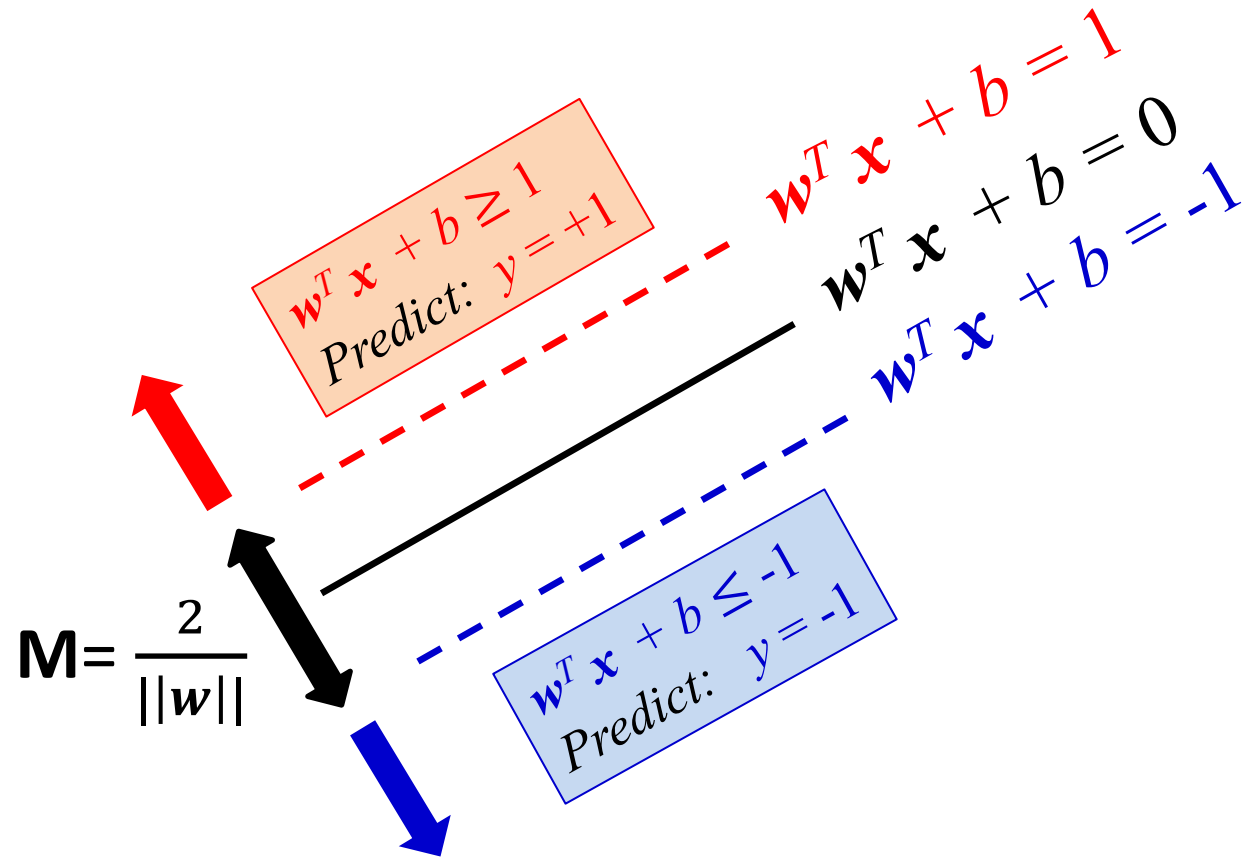


Let's do even better: We can find two parallel lines (rather than just one line) that separate the two classes of data! Red samples are ABOVE the TOP line, and blue samples are BELOW the BOTTOM line. These two hyperplanes can be described by the following equations:



Now, Let's do even better than better: Now, let's find the **two parallel lines** that separate the two classes, **AND** the distance between them is as large as possible (**maximum margin**).



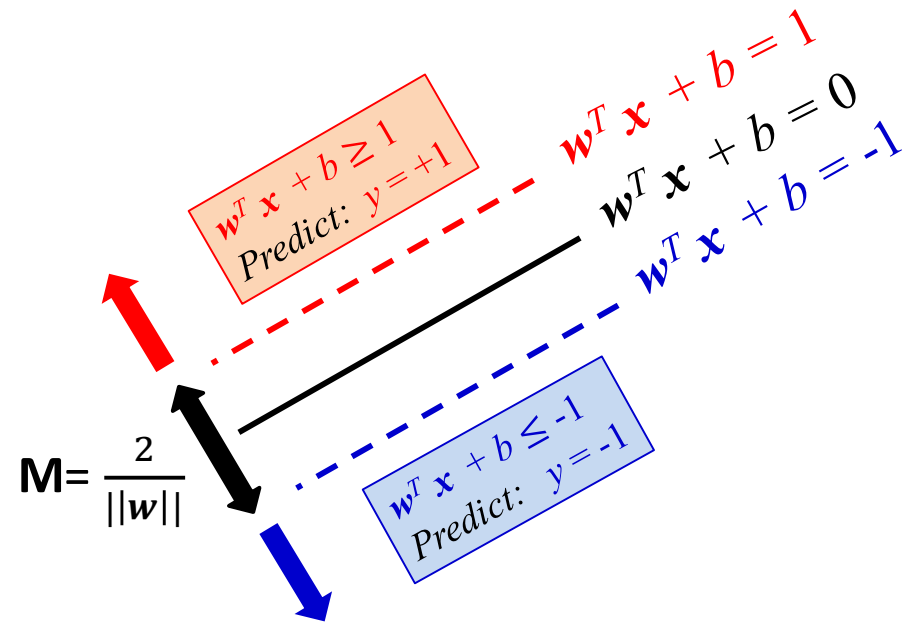


- Given a Training set, the goal is to find the best classifier that classifies the samples correctly and has the maximum margin (maximum $M = \frac{2}{\|w\|}$).

SVM Classifier

- **Formulation:**

$$\left\{ \begin{array}{l} \text{Maximize } \frac{2}{\|w\|} \\ \text{Such that:} \\ \text{for } y = +1, w^T x + b \geq 1 \\ \text{for } y = -1, w^T x + b \leq -1 \end{array} \right.$$



The Optimization Problem

- Thus, given a training set $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$, the new formulation is:

$$\begin{cases} \text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{Subject to: } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ \text{for } i = 1, 2, \dots, n \end{cases}$$

The Optimization Problem (Optional)

Quadratic
programming
with linear
constraints

$$\left\{ \begin{array}{l} \text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{Subject to: } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ \text{for } i = 1, 2, \dots, n \end{array} \right.$$

This is a “***quadratic programming with linear constraints***”. It can be solved using Lagrange Multipliers:

$$\min L_p(\mathbf{w}, b, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1), \alpha_i \geq 0$$

The Optimization Problem (Optional)

This is a “*quadratic programming with linear constraints*”. It can be solved using Lagrange Multipliers:

$$\min L_p(\mathbf{w}, b, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1), \alpha_i \geq 0$$

$$\left\{ \begin{array}{ll} \frac{\partial L_p}{\partial \mathbf{w}} = 0 & \longrightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial L_p}{\partial b} = 0 & \longrightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{array} \right.$$

The Optimization Problem (Optional)

- We solve it by constructing a *Lagrange dual problem*:

Find $\alpha_1 \dots \alpha_N$:

$$\left\{ \begin{array}{ll} \text{maximize} & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} & \alpha_i \geq 0, \text{ and } \sum_{i=1}^n \alpha_i y_i = 0 \end{array} \right.$$

The Optimization Problem (Optional)

- By solving this problem, we will find the best " w " and " b " corresponding to the maximum margin:

$$w = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$
$$b = y_k - w^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ such that } \alpha_k \neq 0$$

- It turns out that only Support Vectors have $\alpha_i \neq 0$
- Support Vectors are the samples which lie nearest to classifier (**on the border**)

SVM Classifier

- By solving this problem, we will find the best $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$ and b corresponding to the maximum margin.
- Then, the ***middle line*** " $(\mathbf{w}^T \mathbf{x} + b) = 0$ " will be our best SVM classifier.
- **SVM Classifier:**

If $\mathbf{w}^T \mathbf{x} + b \geq 0 \rightarrow \text{Predict: } y = +1$
If $\mathbf{w}^T \mathbf{x} + b < 0 \rightarrow \text{Predict: } y = -1$

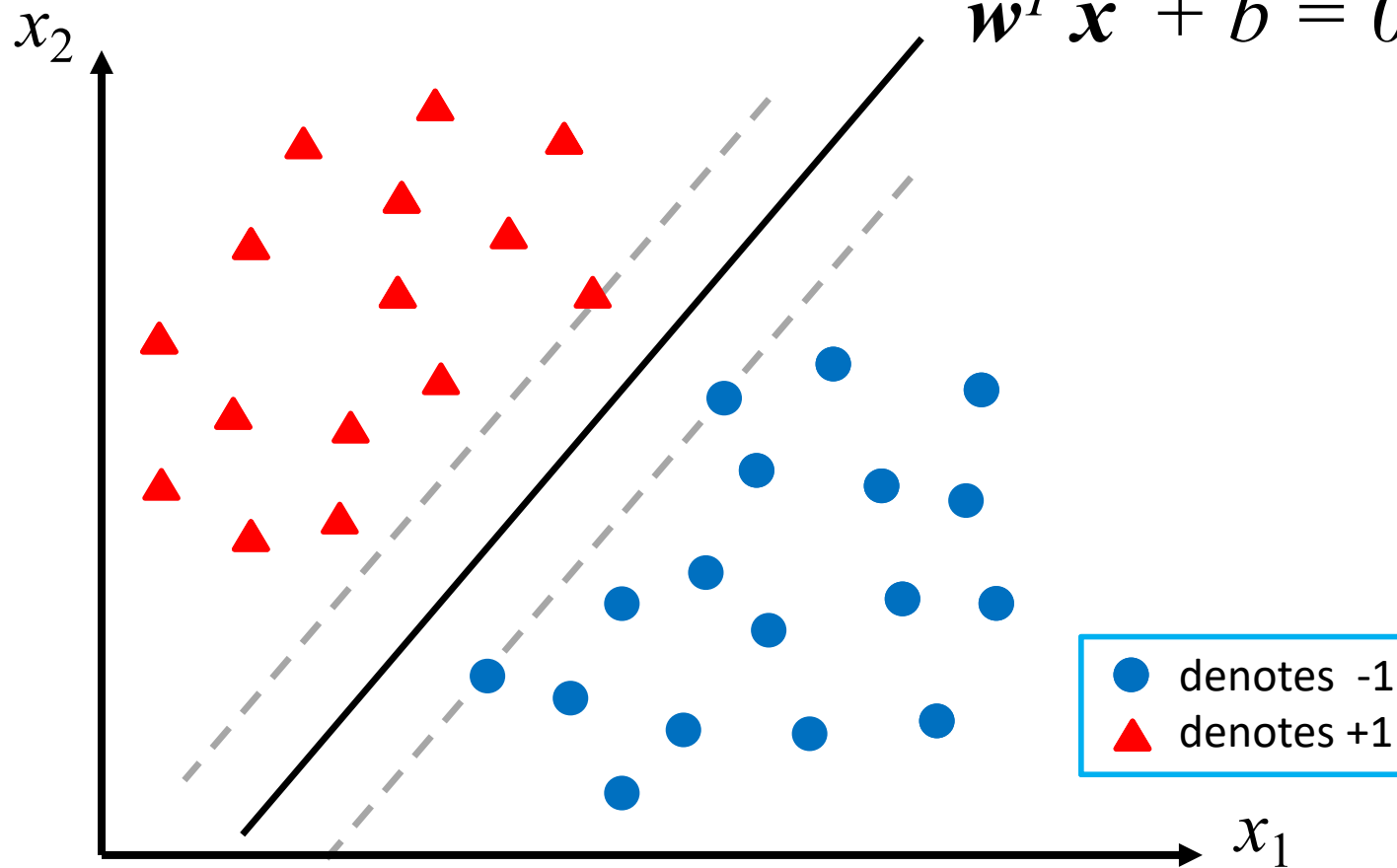
Or Equivalently:

$$y = \text{Sign}(\mathbf{w}^T \mathbf{x} + b)$$

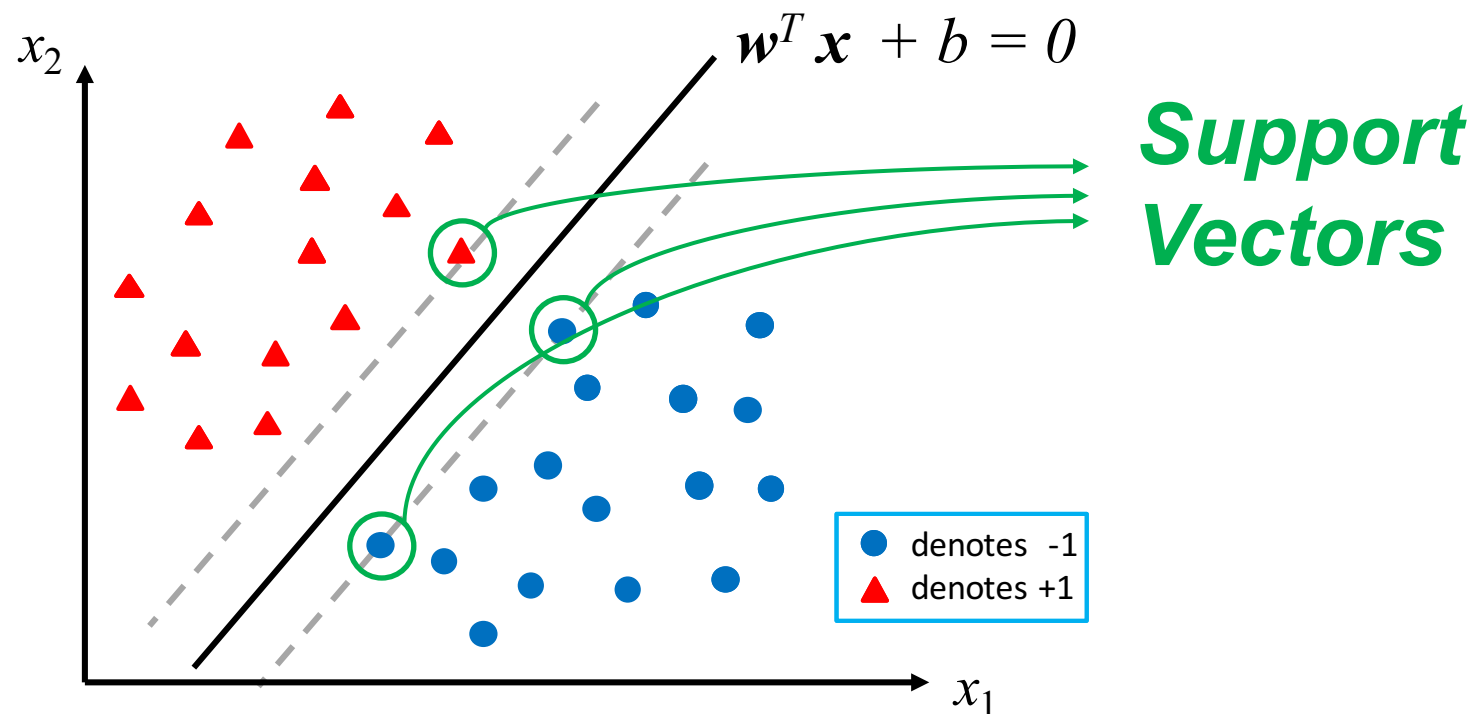
$$, \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

SVM Classifier:

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



Important Note: The max-margin hyperplane is completely determined by those samples \mathbf{x}_i which lie nearest to it (***on the border***). These \mathbf{x}_i 's are called **support vectors**.



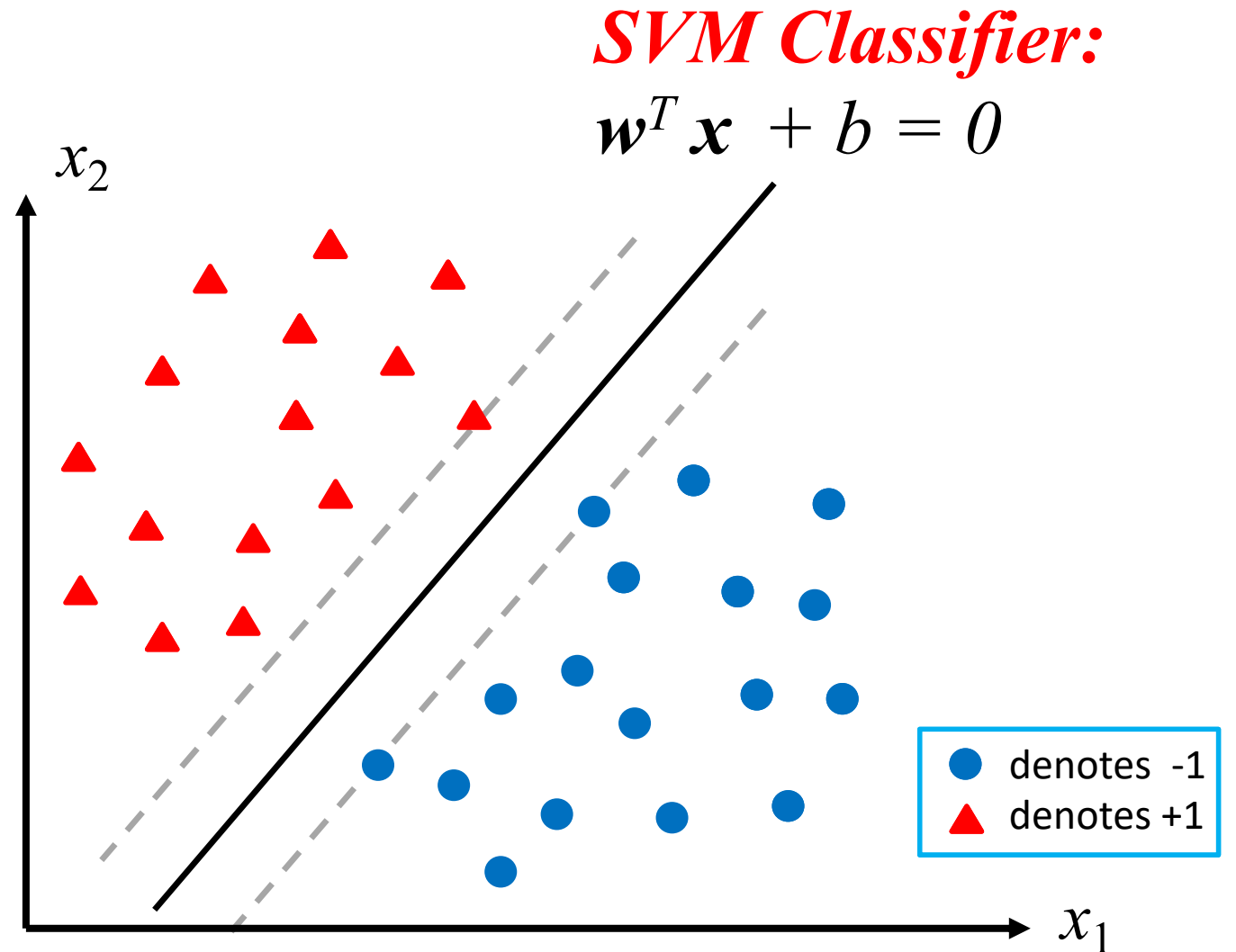
SVM Classifier

- By replacing $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$ in $y = \text{Sign}(\mathbf{w}^T \mathbf{x} + b)$, we will have:
- **SVM Classifier:**
$$y = \text{Sign} \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \right)$$
- **Note1 (Prediction):** The prediction is made by computing the *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i (because only for the support vectors $\alpha_i \neq 0$).
- **Note2 (Training):** To Solve the optimization problem (find the best line), we need to calculate the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points.

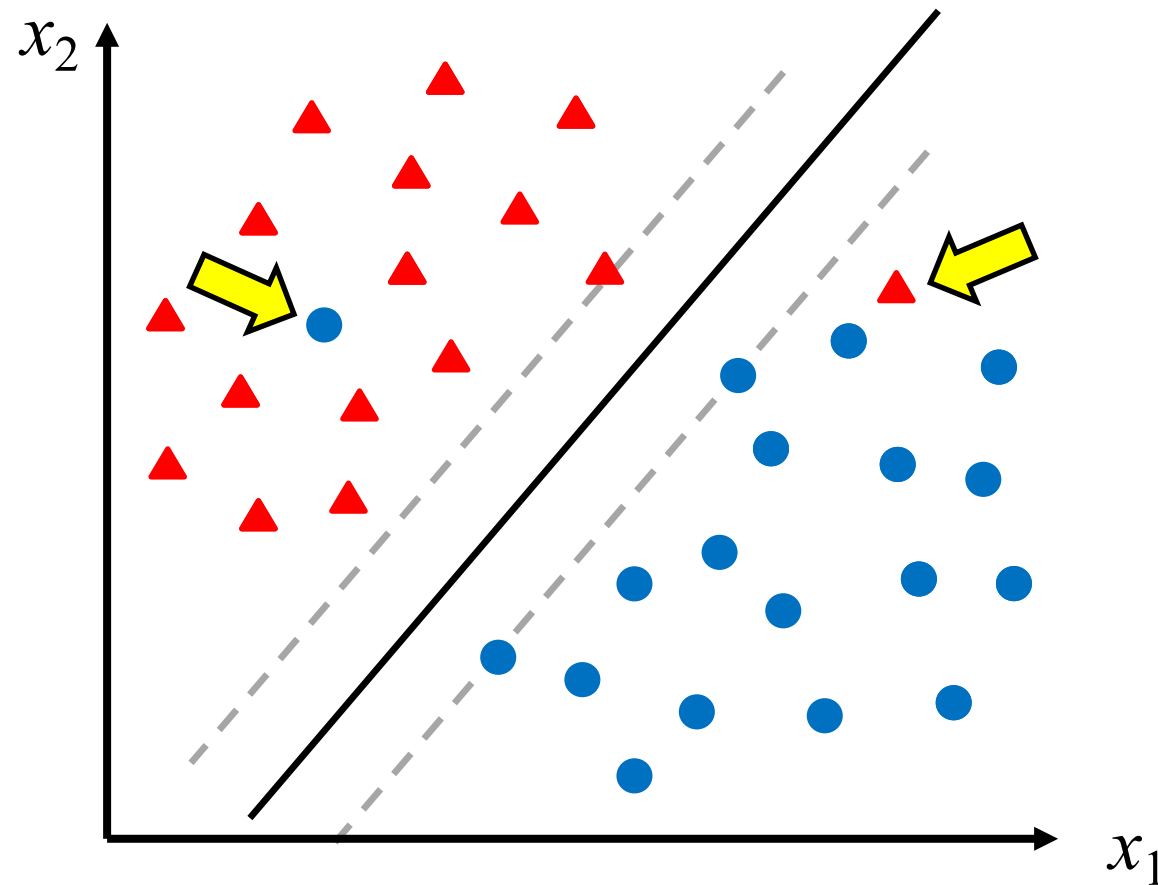
Reminder: Inner product between vector \mathbf{a} and \mathbf{b} : $\mathbf{a}^T \mathbf{b} = \mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^N a_i b_i$

SVM Classifier

- Everything looks good so far!
- But, there is still a **problem** that should be addressed!
- So far, for sake of simplicity, we have assumed that all samples are separable!
- This approach called **Hard Margin**.
- **How about Non-Separable samples???**

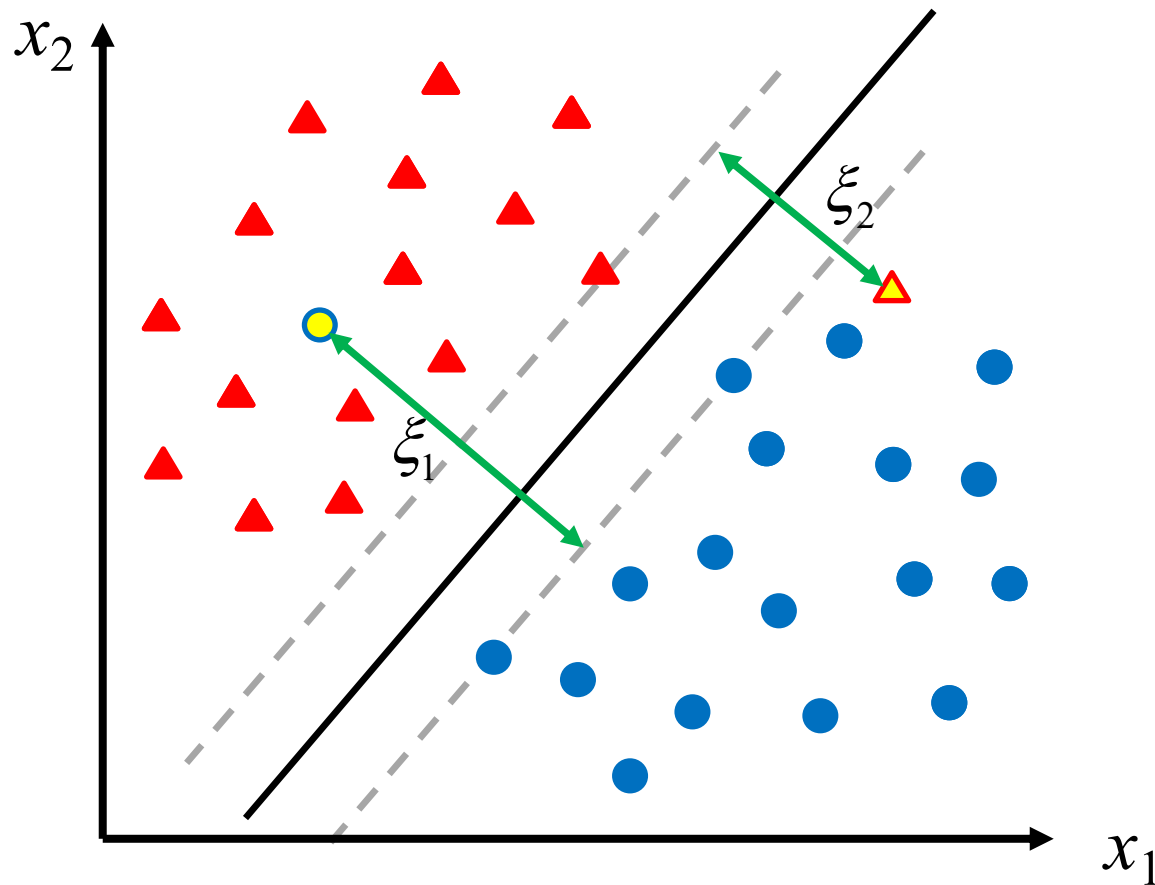


Non-Separable Samples



- So far, we assumed that all data points are separable, and will be classified correctly (no training error)!
- How about the samples that are not separable?

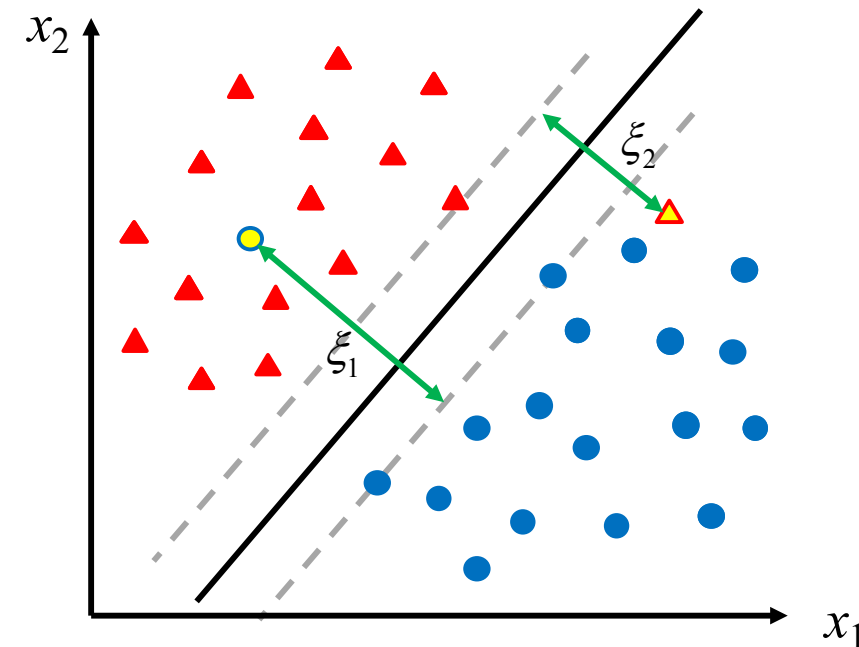
Non-Separable Samples



- How about the samples that are not separable?
- **Solution:** Don't ruin your model because of some bad data samples.
- Instead, we define **Slack Variables** ξ_i to allow mis-classification of a few difficult or bad/noisy data points. (ξ_i is pronounced /zai/)

Non-Separable Samples: Soft Margin

- **Solution:** Don't ruin your model because of some bad data samples.
 - Instead, we define **Slack Variables** ξ_i to allow mis-classification of a few difficult or bad/noisy data points.
 - In fact, we cheat a little by relaxing our criteria, and allow some of the data samples to be on the **wrong side** of the line.
 - This way, we still keep our previous model and just **let it tolerate some errors.**
 - However, remember that we still need to **minimize** this amount of error. So, we add some penalty on **Slack Variables!**
- This is called **Soft Margin.**

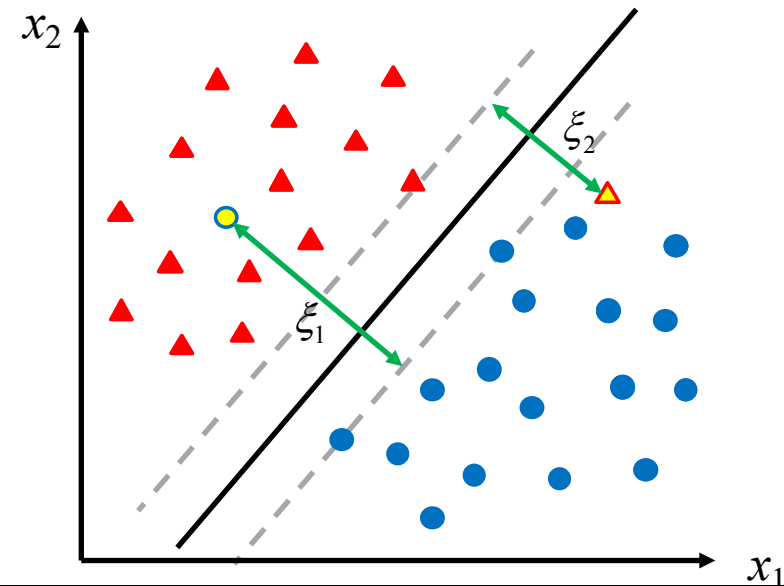


- Hard Margin (our previous formula):

$$\begin{cases} \text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{Subject to: } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ \text{for } i = 1, 2, \dots, n \end{cases}$$

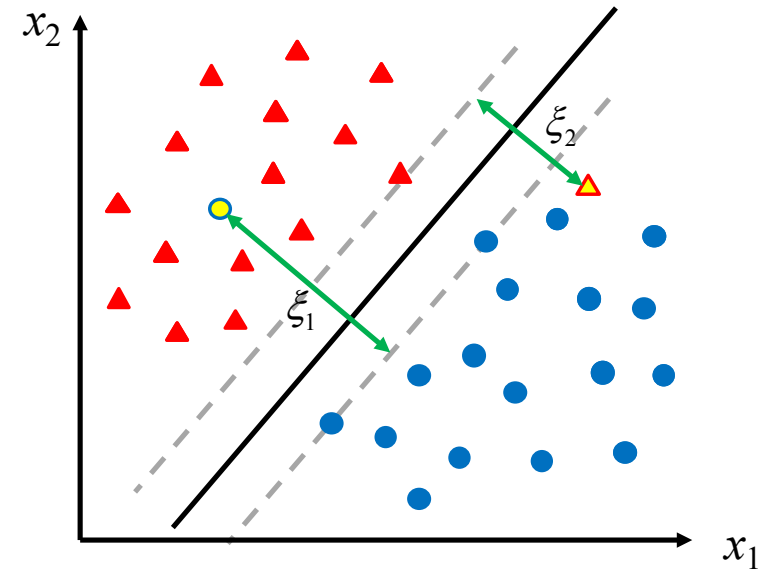
- Soft Margin:

$$\begin{cases} \text{Minimize } \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_i \right) \\ \text{Subject to: } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ \text{for } i = 1, 2, \dots, n \end{cases}$$



Soft Margin:

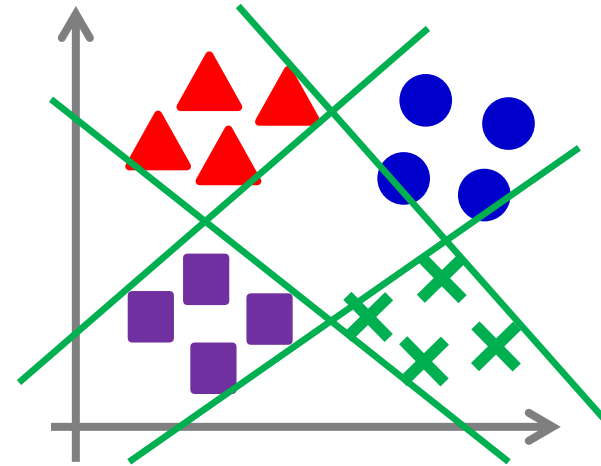
$$\left\{ \begin{array}{l} \text{Minimize } \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_i \right) \\ \text{Subject to: } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ \text{for } i = 1, 2, \dots, n \end{array} \right.$$



- “C” is a very important parameter to adjust SVM results and control Overfitting:
 - Large C: Penalize $\xi \rightarrow$ No tolerance for misclassification \rightarrow Low bias, high variance \rightarrow prone to overfitting.
 - Small C: Let larger $\xi \rightarrow$ allow misclassification \rightarrow High bias, low variance \rightarrow prone to underfitting.

Multi-Class Classification

- SVM is originally a binary classifier (2 labels). However, it is possible to perform multiclass classification using **One-vs-All method**.



- **One-vs-All method**: In this approach, assuming that we have K different labels, we train K SVM-classifiers, one to distinguish class k from the rest (for $k=1,2,...,K$), and learn the weights $\mathbf{w}^{(k)}$ for each classifier ($\mathbf{w}^{(1)}$ for classifying class1 vs. the rest, $\mathbf{w}^{(2)}$ for classifying class2 vs. the rest, ...). Then, in testing stage, for a new sample \mathbf{x} , pick the class with the largest $(\mathbf{w}^T \mathbf{x}_i + b)$.

SVM in Python

- Linear SVM:

```
from sklearn.svm import LinearSVC  
my_SVM = LinearSVC(parameters)  
my_SVM.fit (X_train, y_train)  
y_predict = my_SVM.predict (X_test)
```

- Non-Linear SVM (will cover it next):

```
from sklearn.svm import SVC  
my_SVM = SVC(parameters)  
my_SVM.fit (X_train, y_train)  
y_predict = my_SVM.predict (X_test)
```

Thank You!

Questions?