

Advanced Machine Learning and Deep Learning

Dr. Mohammad Pourhomayoun

Assistant Professor

Computer Science Department

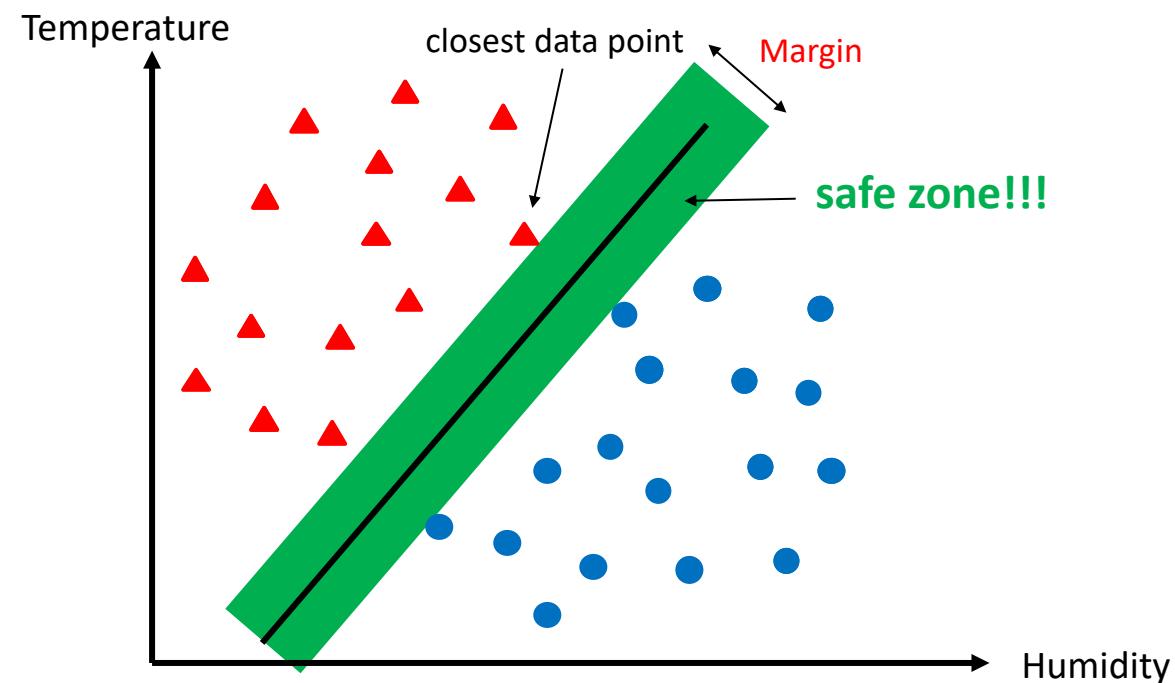
California State University, Los Angeles



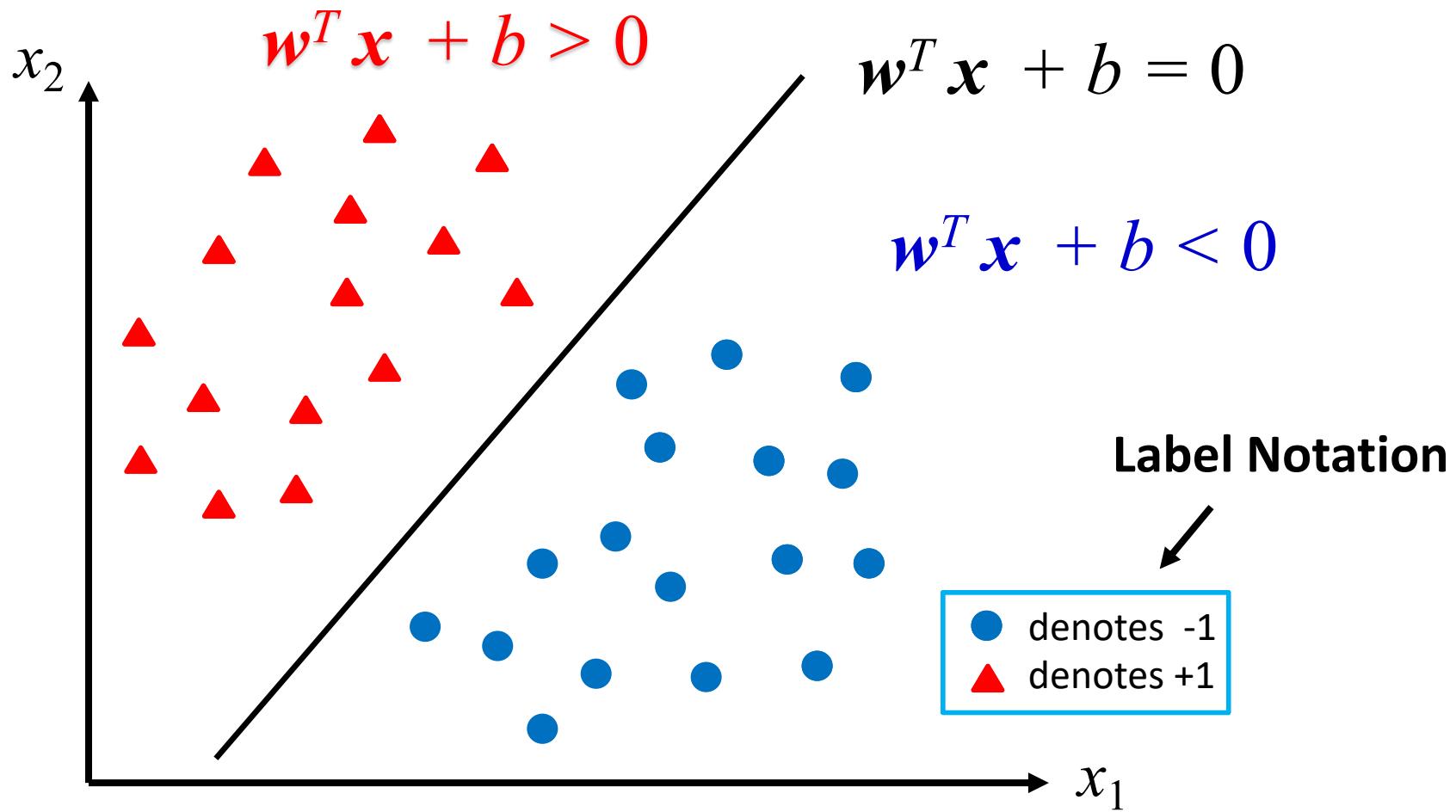
Support Vector Machine (SVM)

Review

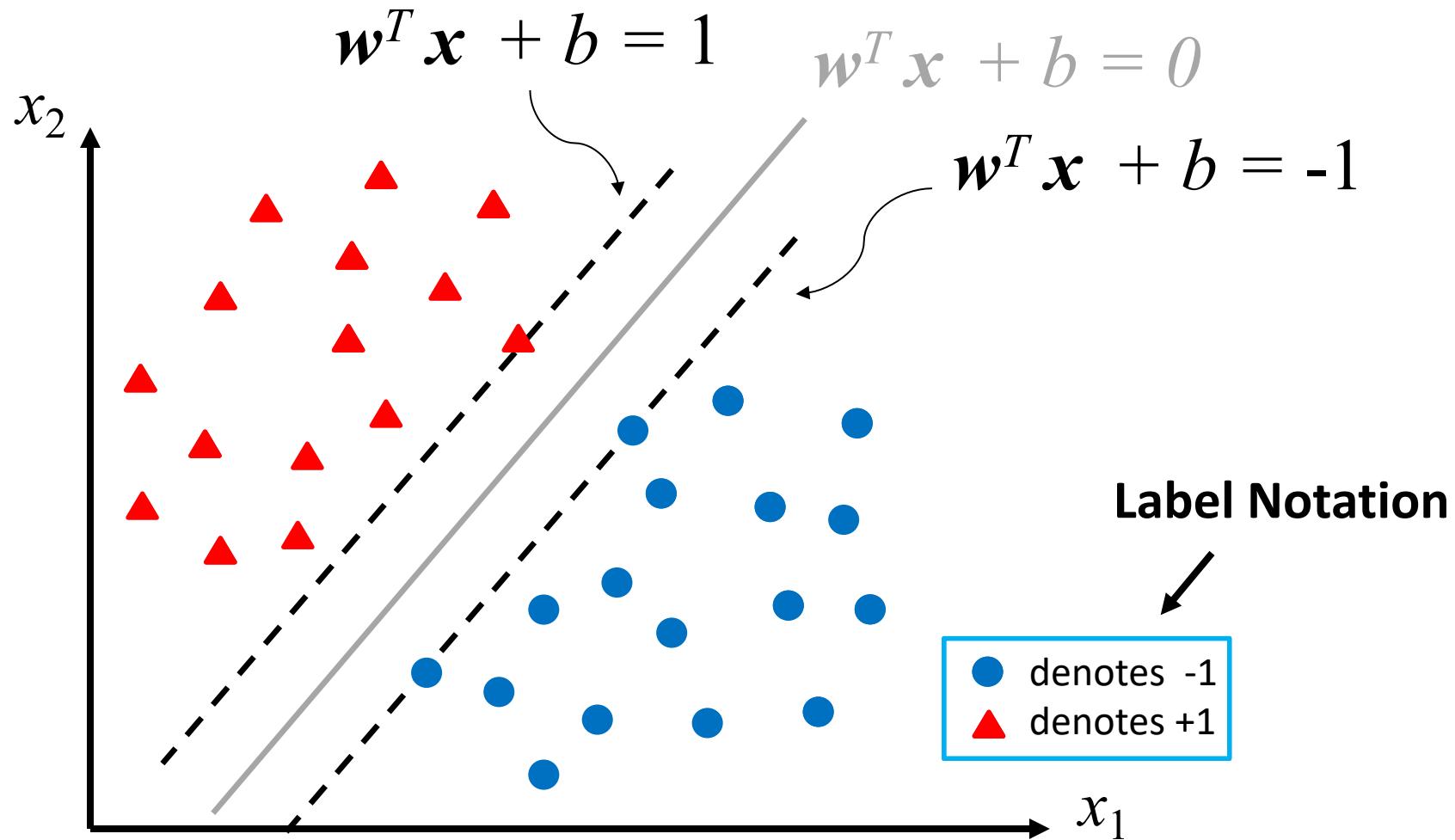
- The best classifier is the one with the **Maximum Margin**.
- **Margin** is widest boundary area before hitting a data point.



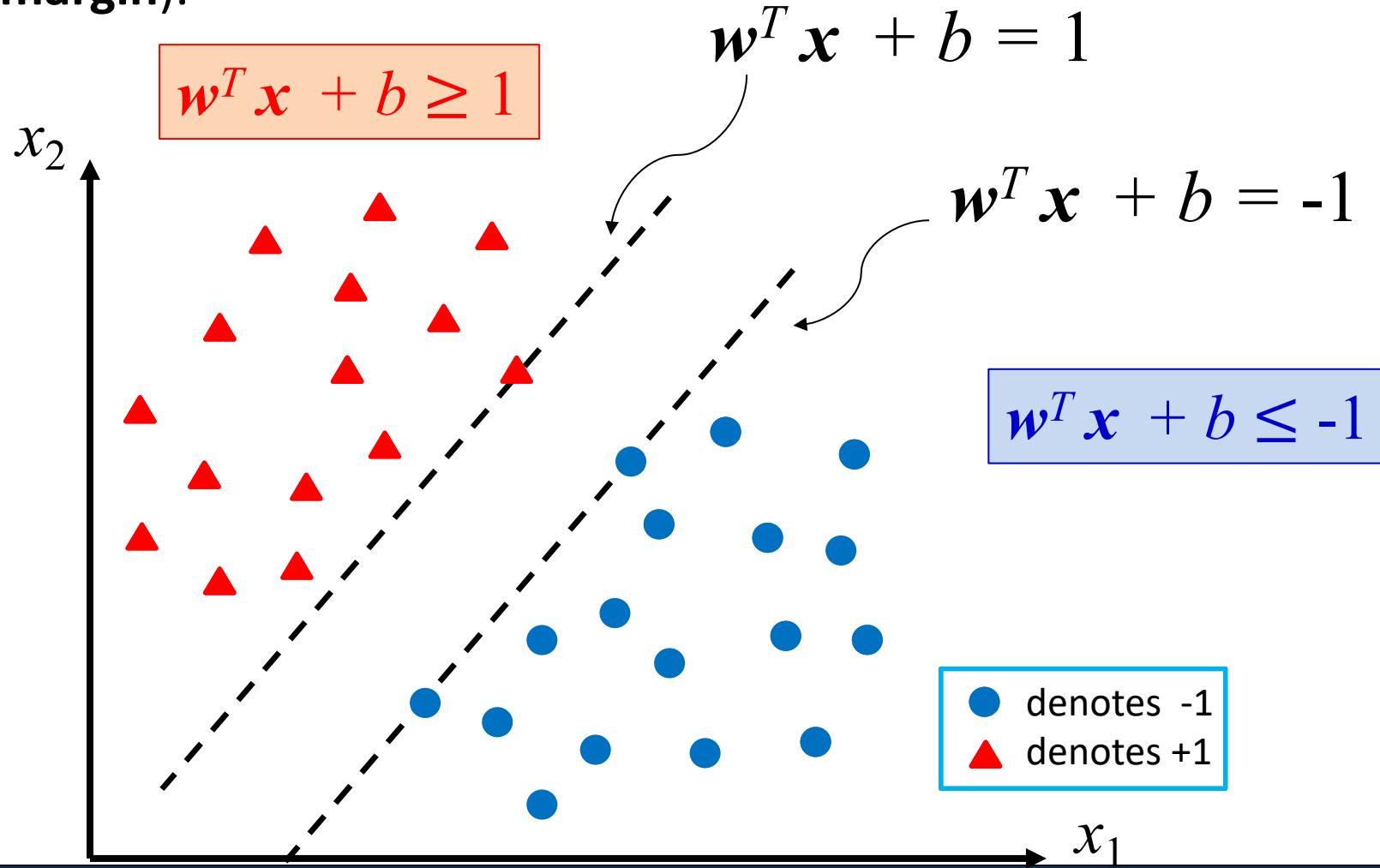
- Given a set of training data points, each point will be considered as a ***m-dimensional vector*** in space. We are looking for ***a line whose value is positive for red samples, and negative for blue samples.***



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:



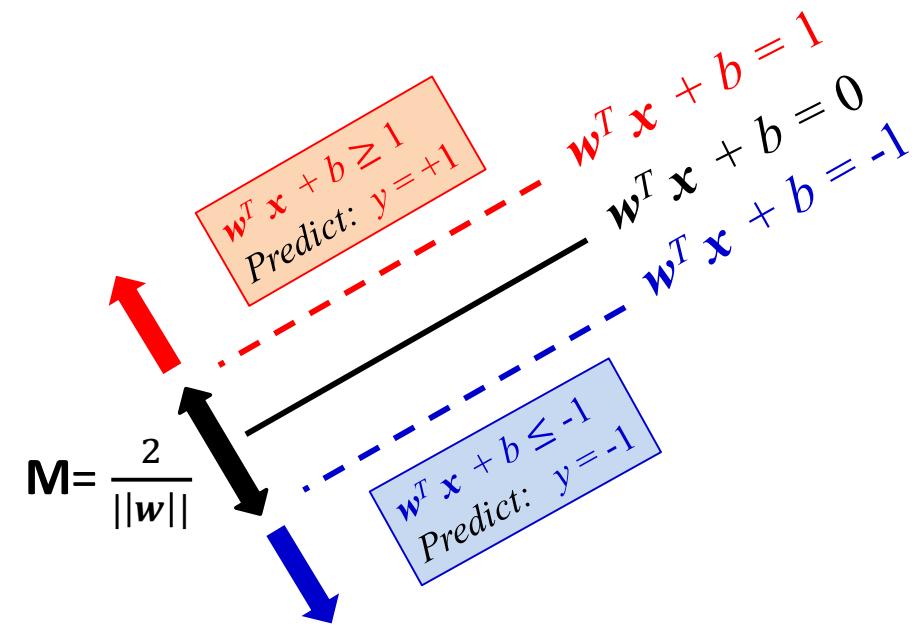
Now, Let's do even better than better: Now, let's find the **two parallel lines** that separate the two classes, **AND** the distance between them is as large as possible (**maximum margin**).



SVM Classifier

- **Formulation:**

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

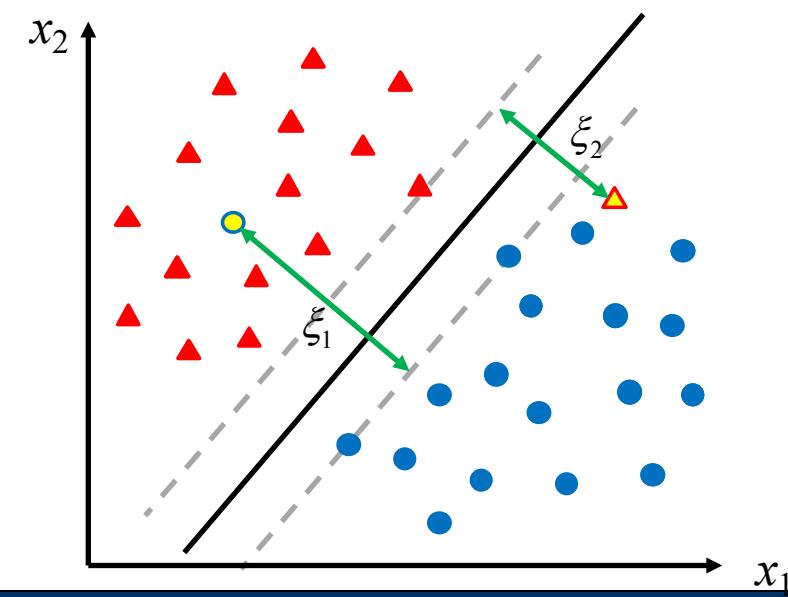


- **Hard Margin (our previous formula):**

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

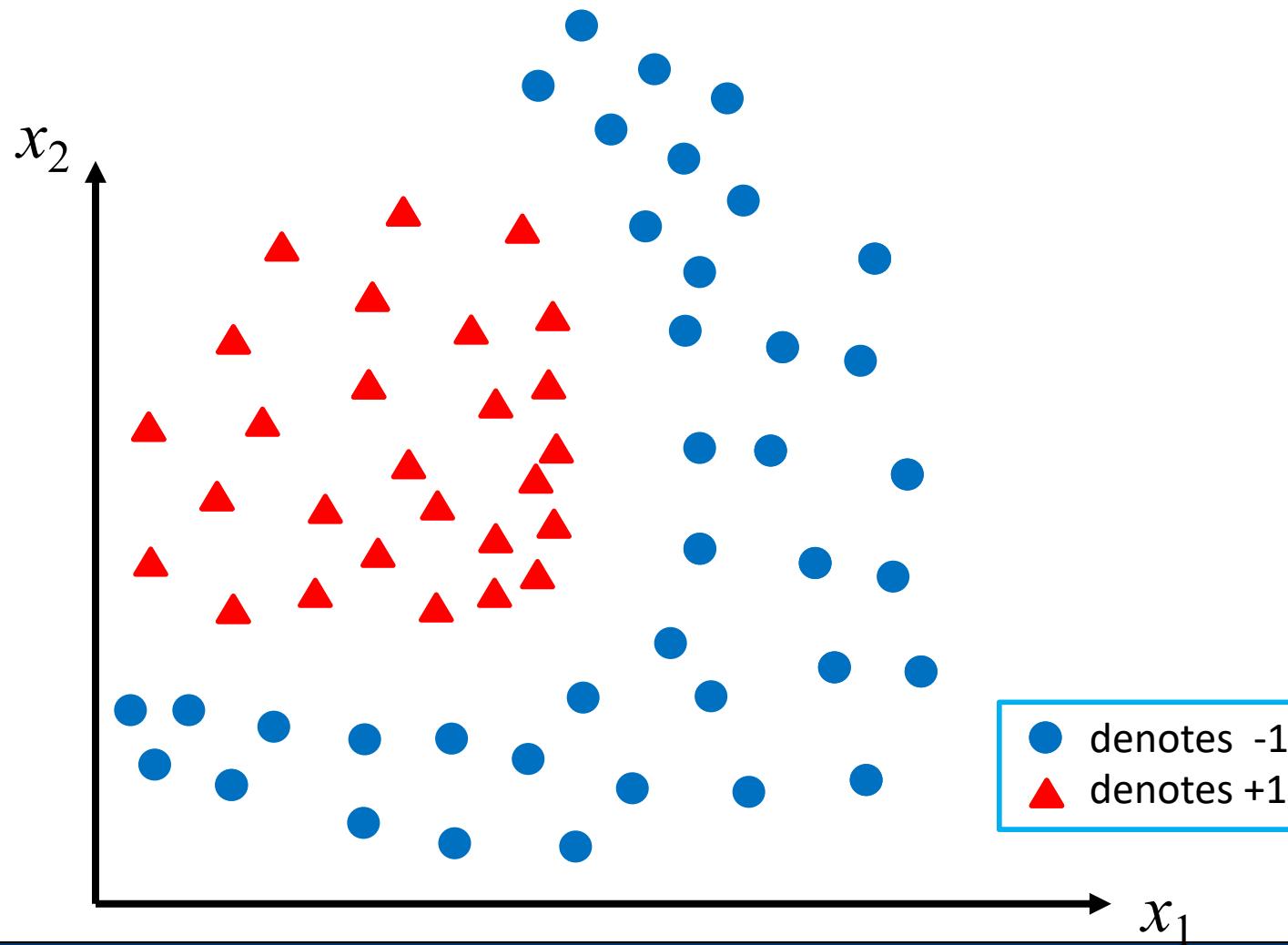
- **Soft Margin:**

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

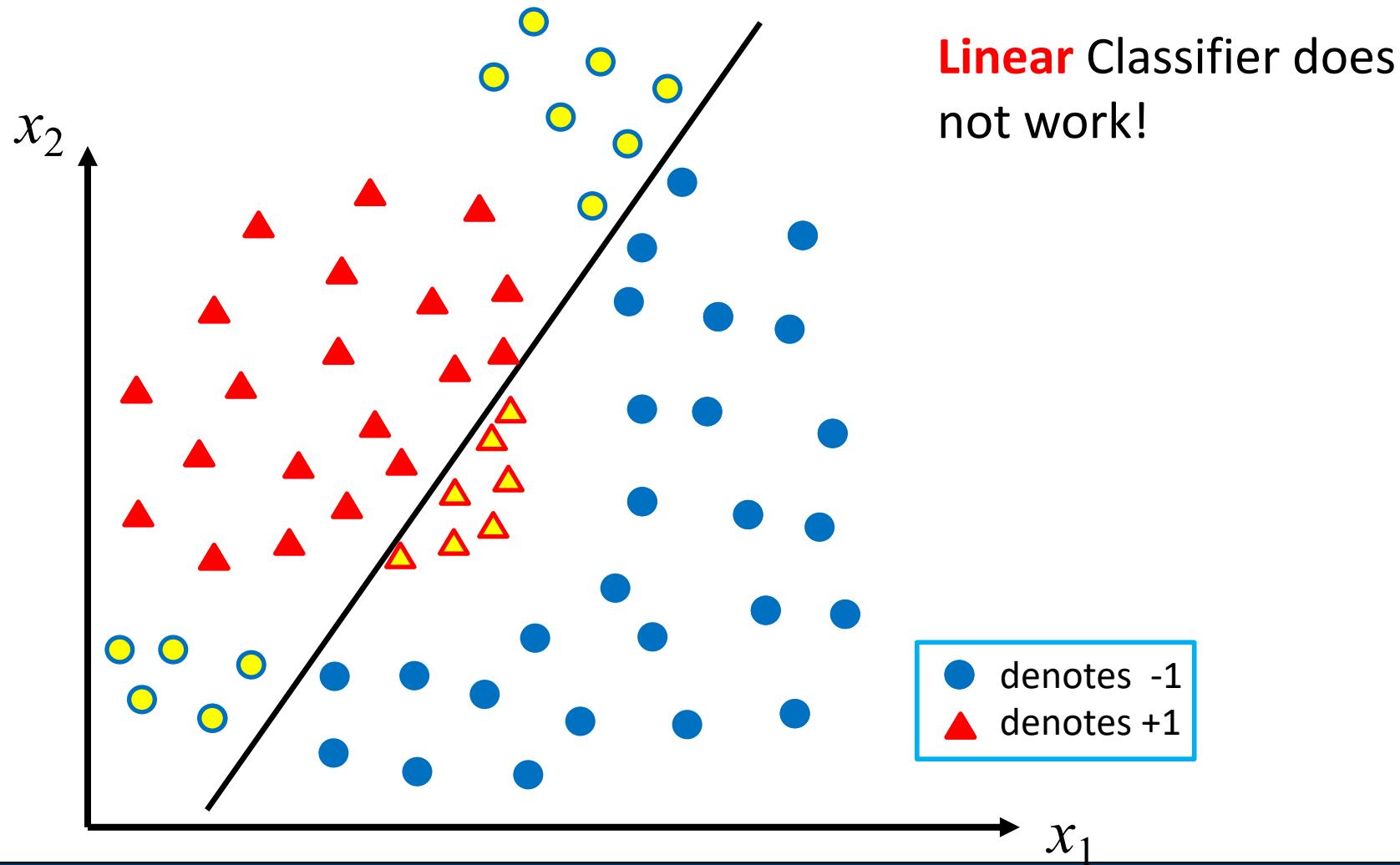


Non-Linear SVM and Kernel Trick

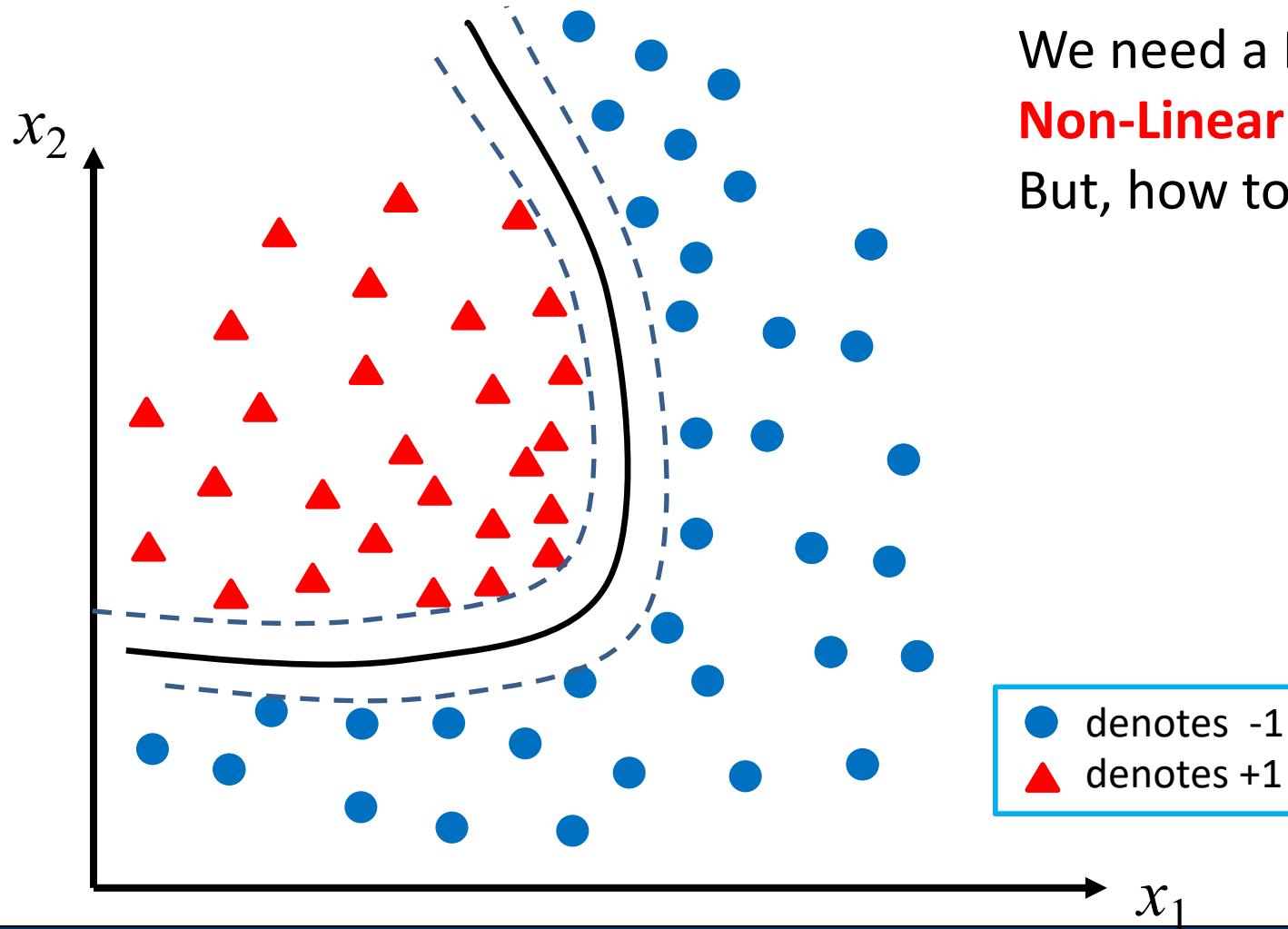
What if the data is not linearly separable!!?



What if the data is not linearly separable!!?



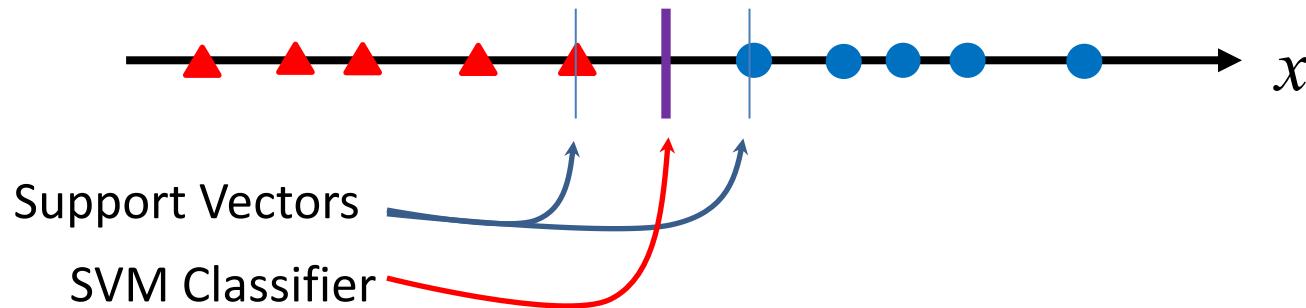
What if the data is not linearly separable!!?



We need a Maximum-Margin
Non-Linear Classifier!
But, how to do that?

A Simple Example in 1-D

- In this example, the Data points are linearly separable:



- In this example, the Data points are **NOT** linearly separable:



- How to classify this dataset!!?**

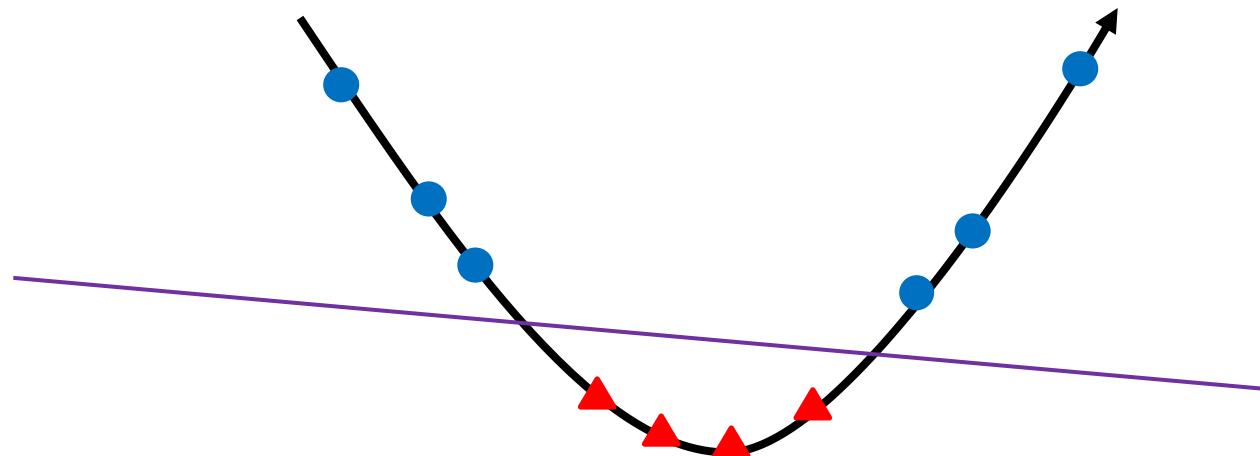
Ref: Douglas Eck, Music and Machine Learning, University of Montreal

A Simple Example in 1-D

- Data points are NOT linearly separable:



- What if we bend the space!!?
- We can **move the data into a higher-dimensional space** to make it linearly separable:

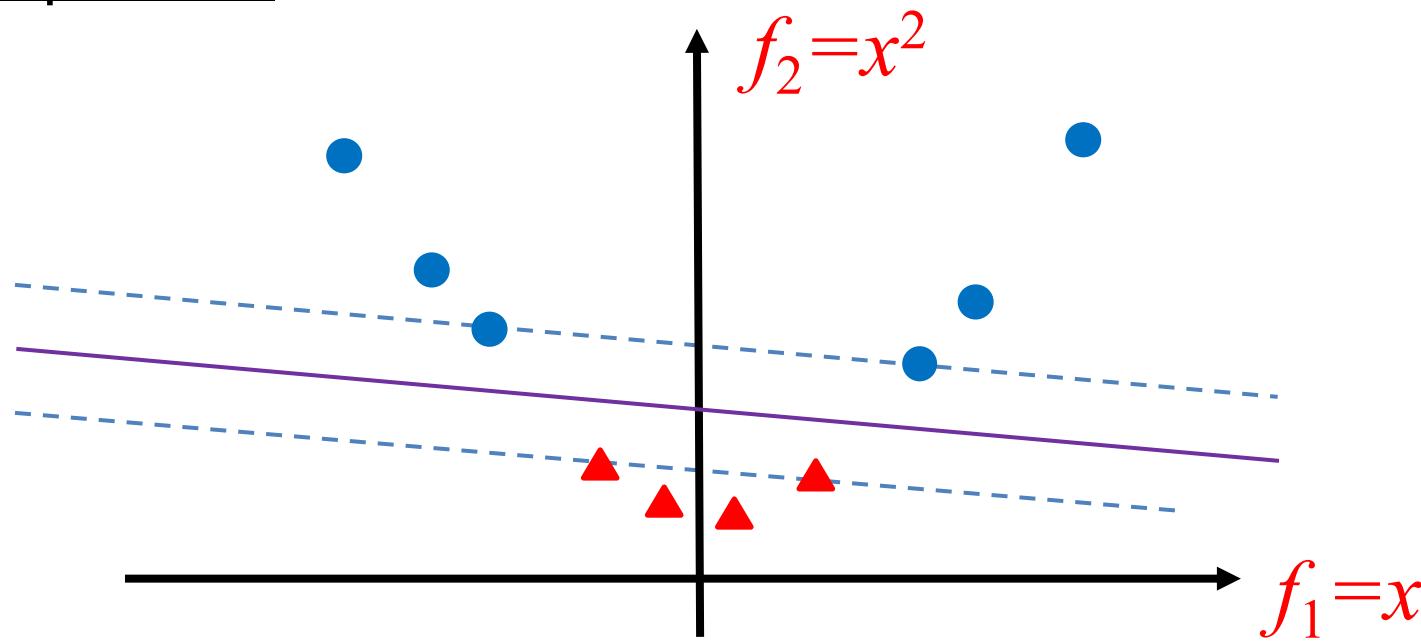


A Simple Example in 1-D

- Data points are NOT linearly separable:



- We can move the data into a higher-dimensional space to make it linearly separable:

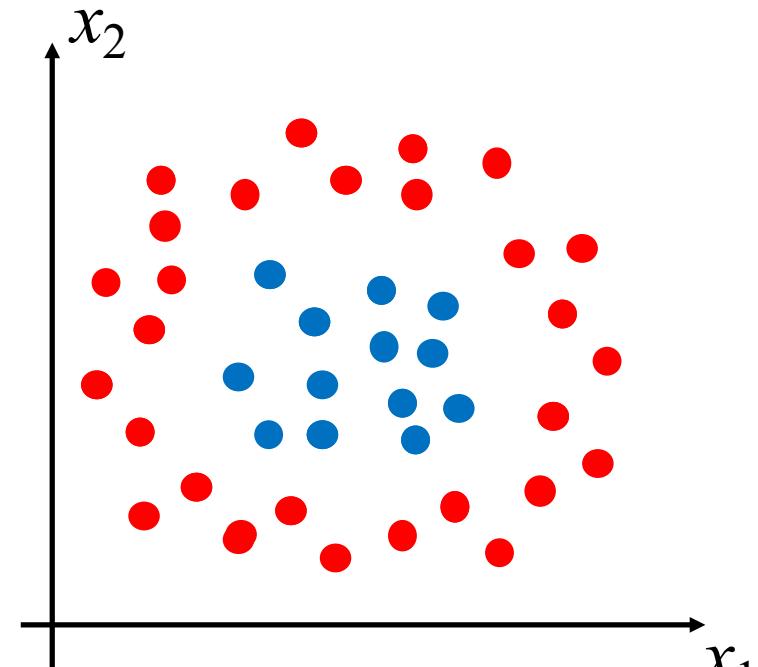


A Simple Example in 2-D

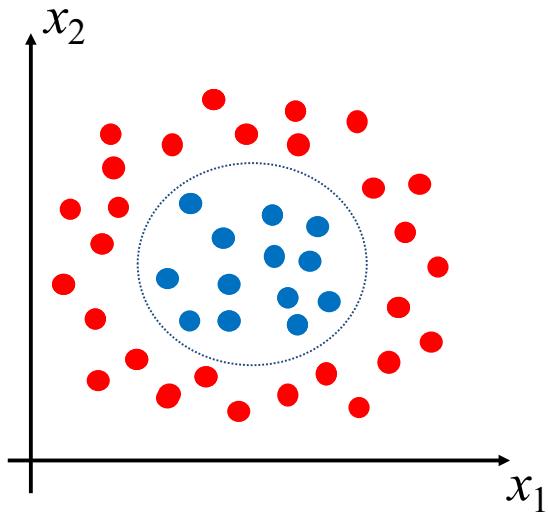
- In this example, the Data points are **NOT** linearly separable: you can **not** find a linear hyperplane (a line in 2D) to separate them:
- Can we **move the data into a higher-dimensional space** to make it linearly separable?
- We need a function Φ (/fai/) that moves data samples x to a higher-dimensional space:

$$\Phi: x \rightarrow \varphi(x)$$

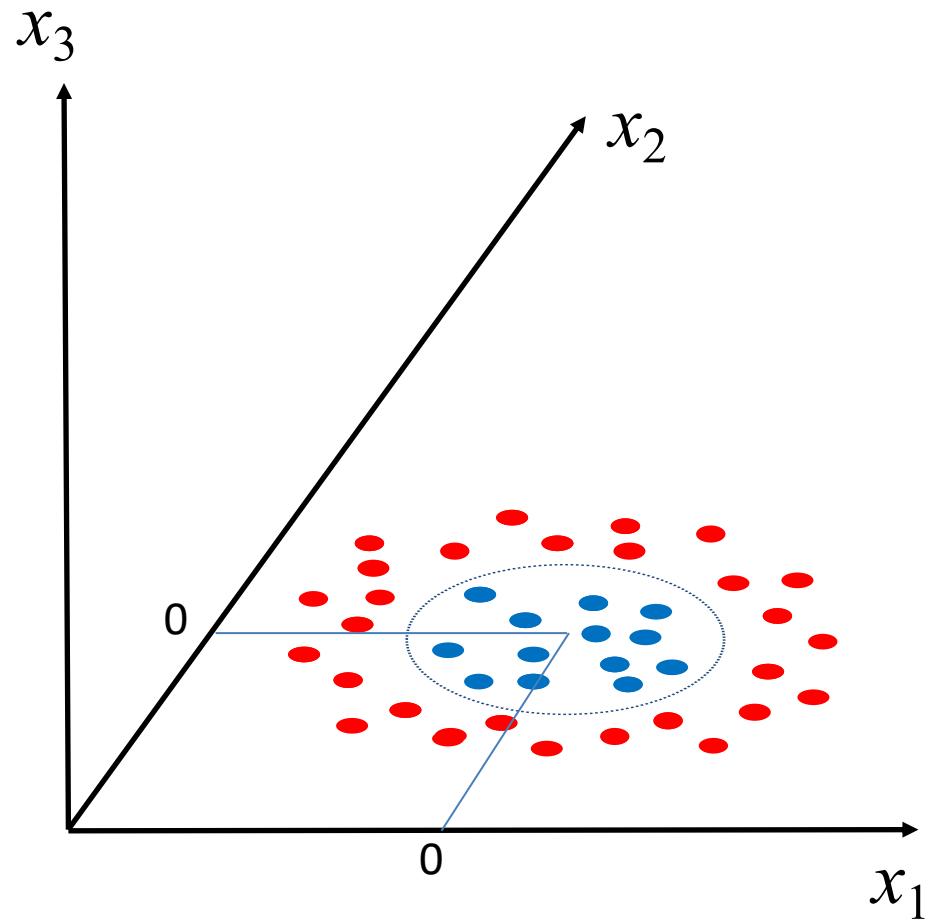
Low-dimensional space High-dimensional space



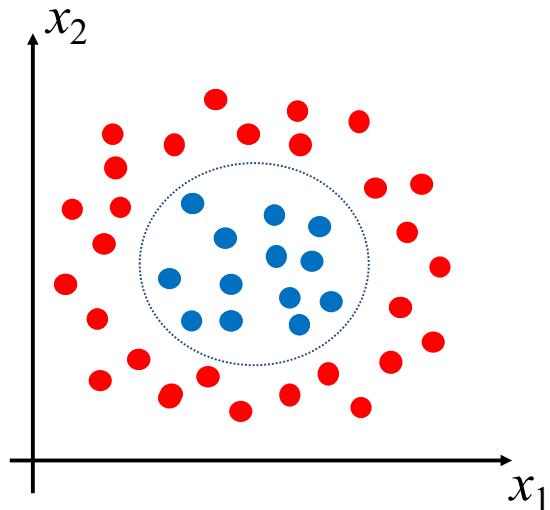
A Simple Example in 2-D



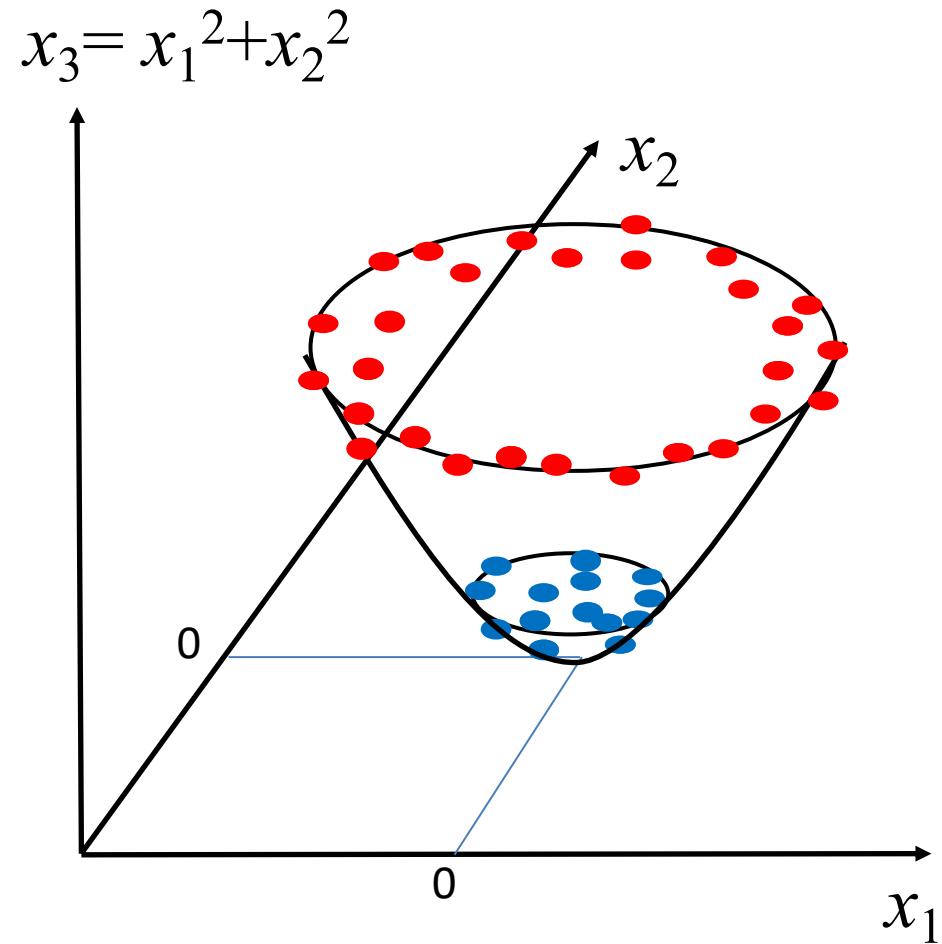
$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



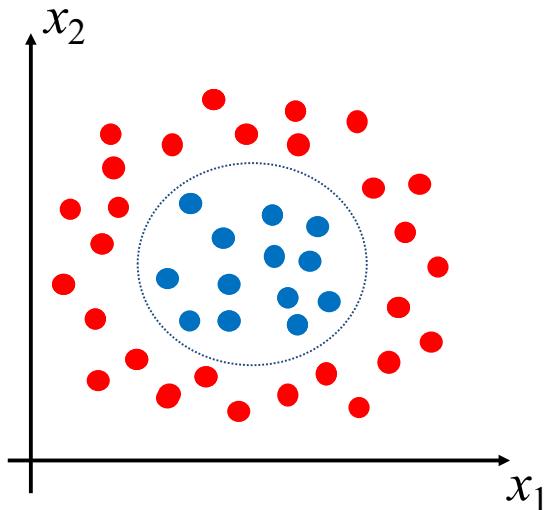
A Simple Example in 2-D



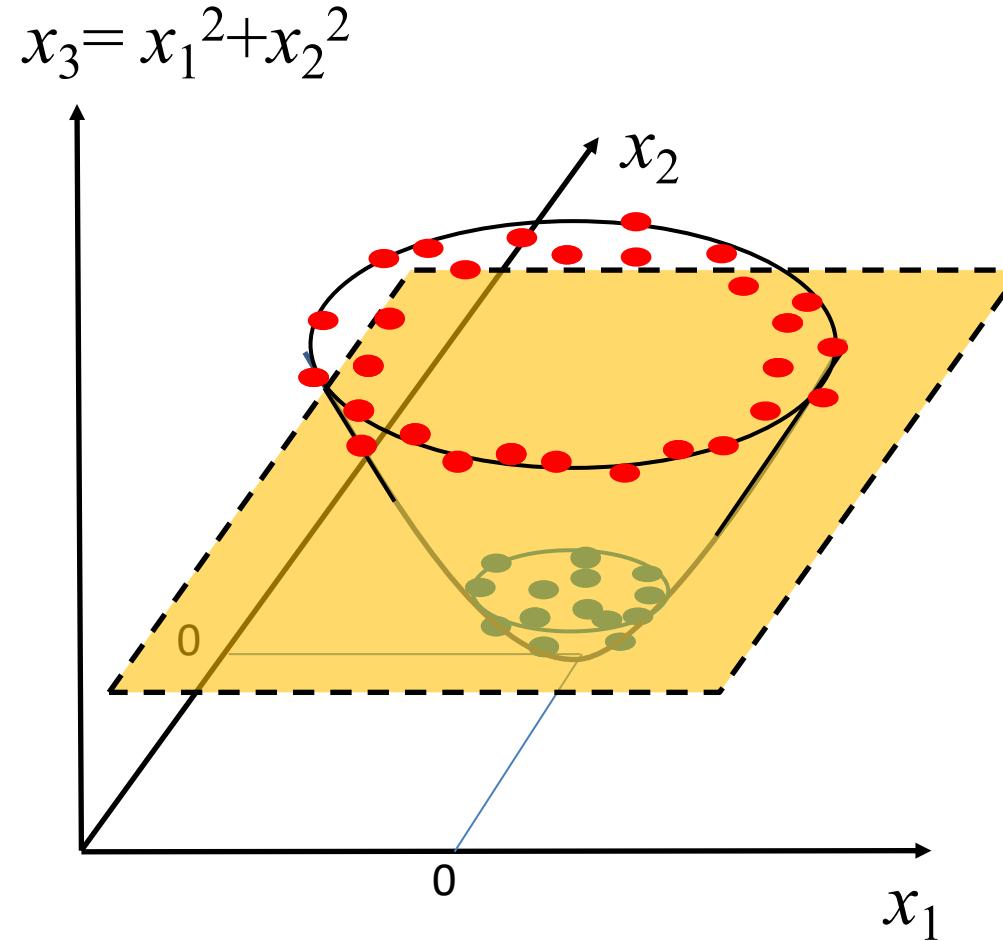
$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



A Simple Example in 2-D



$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



Nonlinear SVM

- **The main idea for non-linear SVM:** The original data can be transformed to some **higher-dimensional feature space where the data is linearly separable.**
- **Any Problem with this idea?**
 - What is the best new Space?
 - In general, we don't know the best function Φ that maps the data samples to a new space.
 - Even if we know it, the computation in the new space can be costly because it is high dimensional.
- **Solution?**
 - The kernel Trick comes to help!

Review: SVM Classifier

- By putting $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$, in $y = \text{Sign}(\mathbf{w}^T \mathbf{x} + b)$, we will have:

$$y = \text{Sign}\left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b\right)$$

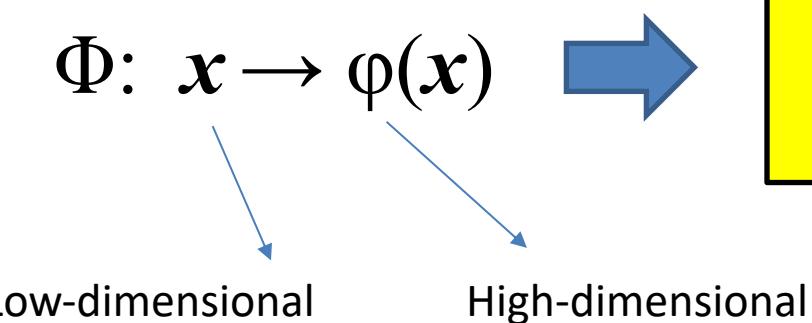
- Reminder: Inner product between vector \mathbf{a} and \mathbf{b} : $\mathbf{a}^T \mathbf{b} = \mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^N a_i b_i$
- **Note1:** The prediction is made by computing the *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i (where $\alpha_i \neq 0$).
- **Note2:** To Solve the optimization problem we need to calculate the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points.
- **THUS, everything appears in the form of inner products!**

Nonlinear SVM

- Original Space:

$$y = \text{Sign} \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \right)$$

- Now, If we map the data into a higher-dimensional space using a function Φ , then our prediction function is:



$$y = \text{Sign} \left(\sum_{i=1}^n \alpha_i y_i \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}) + b \right)$$

Nonlinear SVM

- High-Dimensional Space ($\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$):

$$y = \text{Sign} \left(\sum_{i=1}^n \alpha_i y_i \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}) + b \right)$$

- Again, everything appears in the form of inner products $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x})$.
So, everything will become a scalar number again!
- Thus, we really do NOT need to know $\varphi(\mathbf{x})$. We just need to find $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x})$.
- This is called Kernel Function: $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$

Nonlinear SVM: Kernel Trick!

- By putting $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$, the new classification rule based on **Kernel Trick** is:

$$y = \text{Sign} \left(\sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \right)$$

- A **kernel function** is a function that corresponds to an inner product in some high-dimensional expanded feature space.
- In other word, The kernel function plays the role of the inner product in high dimensional space.

Nonlinear SVM: Kernel Trick!

- **Summary:**
 - I. The main idea: The original data can be mapped to some higher-dimensional feature space, where the data is linearly separable.
 - II. Assume that $\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ is the function that can map the data to some higher-dimensional space.
 - III. It turns out that in every calculation in both testing and training stages, every time, $\varphi(\mathbf{x})$ only appears in the form of inner product $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$.
 - IV. Thus, we really do not need to know the mapping function φ explicitly!
 - V. Instead, we define and use a *kernel function* as a function that corresponds to the inner product of two feature vectors in some expanded high-dimensional feature space: $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$.

Nonlinear SVM: Kernel Trick

- Examples of most commonly-used kernel functions:

- Polynomial kernel:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{r} + \gamma \mathbf{x}_i^T \mathbf{x}_j)^d$$

- Gaussian Radial-Basis Function (RBF):

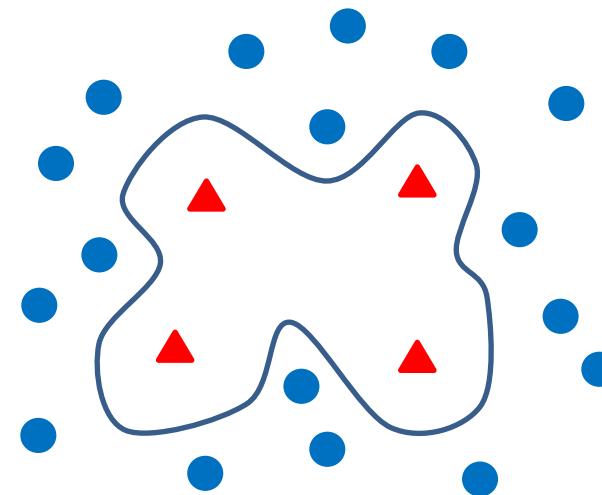
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

- Hyperbolic Tangent:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\mathbf{r} + \gamma \mathbf{x}_i^T \mathbf{x}_j)$$

Another Interpretation (Optional)

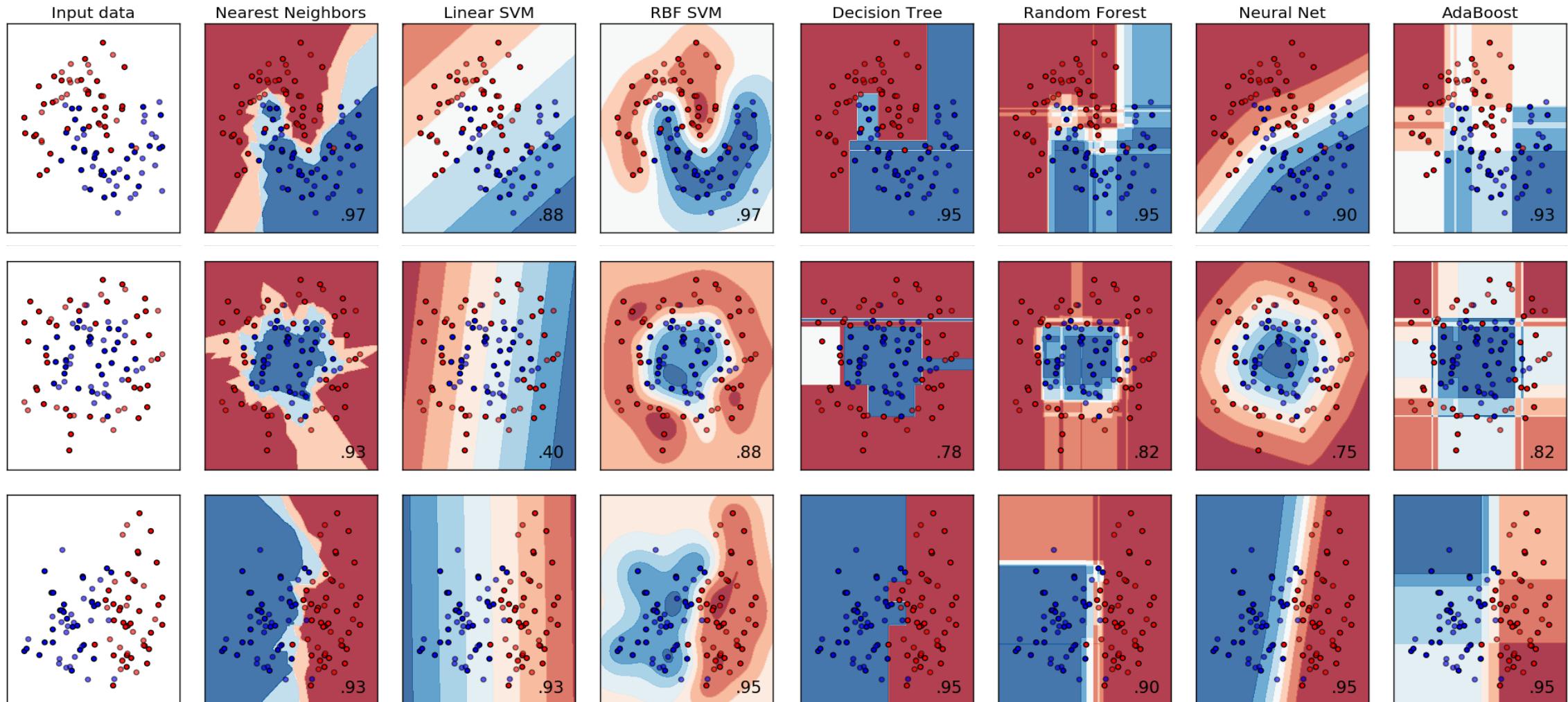
- Kernel function represents some kind of **similarity** (or closeness) between two vectors in the space. That is why sometimes it is called ***Similarity Function***.
- Thus, the Kernel between two vectors in training and testing stages, can measure the closeness of a testing vector to the training vectors.
- This closeness, somehow defines a **non-linear** hidden border that can be our decision boundary.



- Now, after CS4661 and CS4662, you know how the **most popular** and the **most effective** data analytics and machine learning algorithms work! More importantly, you know how to use them! ☺



Classifier Comparison



Ref: Sklearn Documentation

Project Progress Report

- **Due Date for Project Progress Report: Sunday, March 28.**
- Please submit your progress report on Canvas by March 28.
- **Your report should include:**
 - Team members information and their responsibilities.
 - Project description/details.
 - Your results, methods, progress so far.
- **Note: One submission per group (the team leader submits on behalf of everyone).**

Thank You!

Questions?