# Support Vector Machine (SVM) 1

| المحاضرة 6 | كلية الهندسة المعلوماتية | مقرر تعلم الآلة |
|---|---|---|

**SYRIAN PRIVATE UNIVERSITY**

الجامعة السورية الخاصة

د. رياض سنبل

# ML Pipeline



Data
مجموعة معطيات

مجموعة معطيات الاختبار          مجموعة معطيات التعلم

Feature Engineering
هندسة السمات

Training Approach
منهجية التعلم

New input

model
النموذج المتعلَّم

Predicted output
الخرج المتوقع
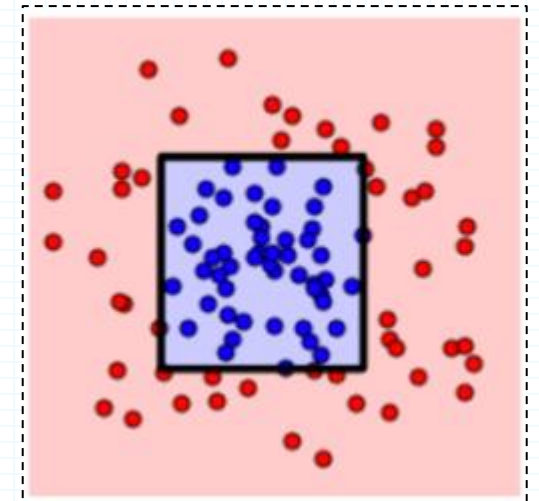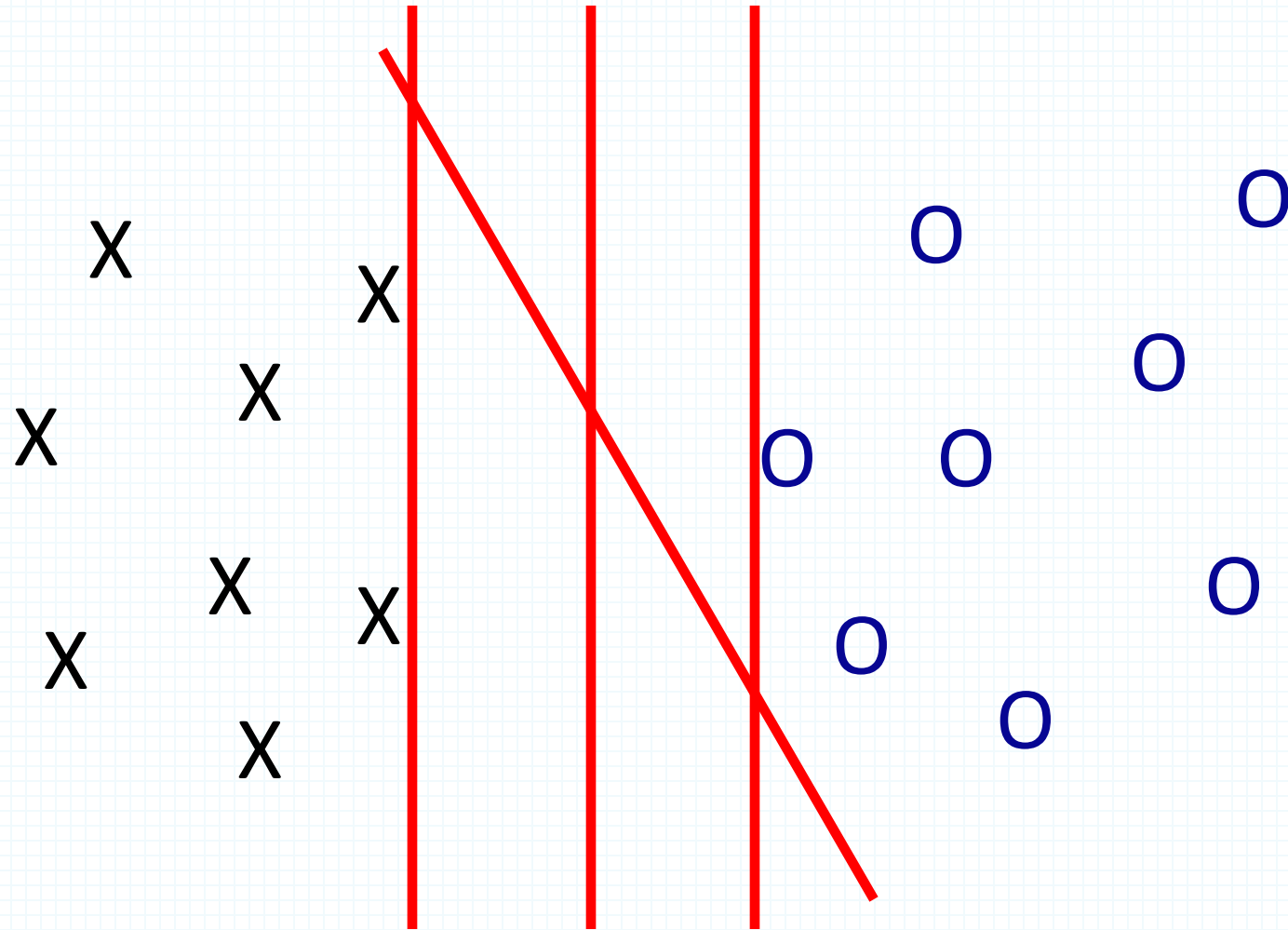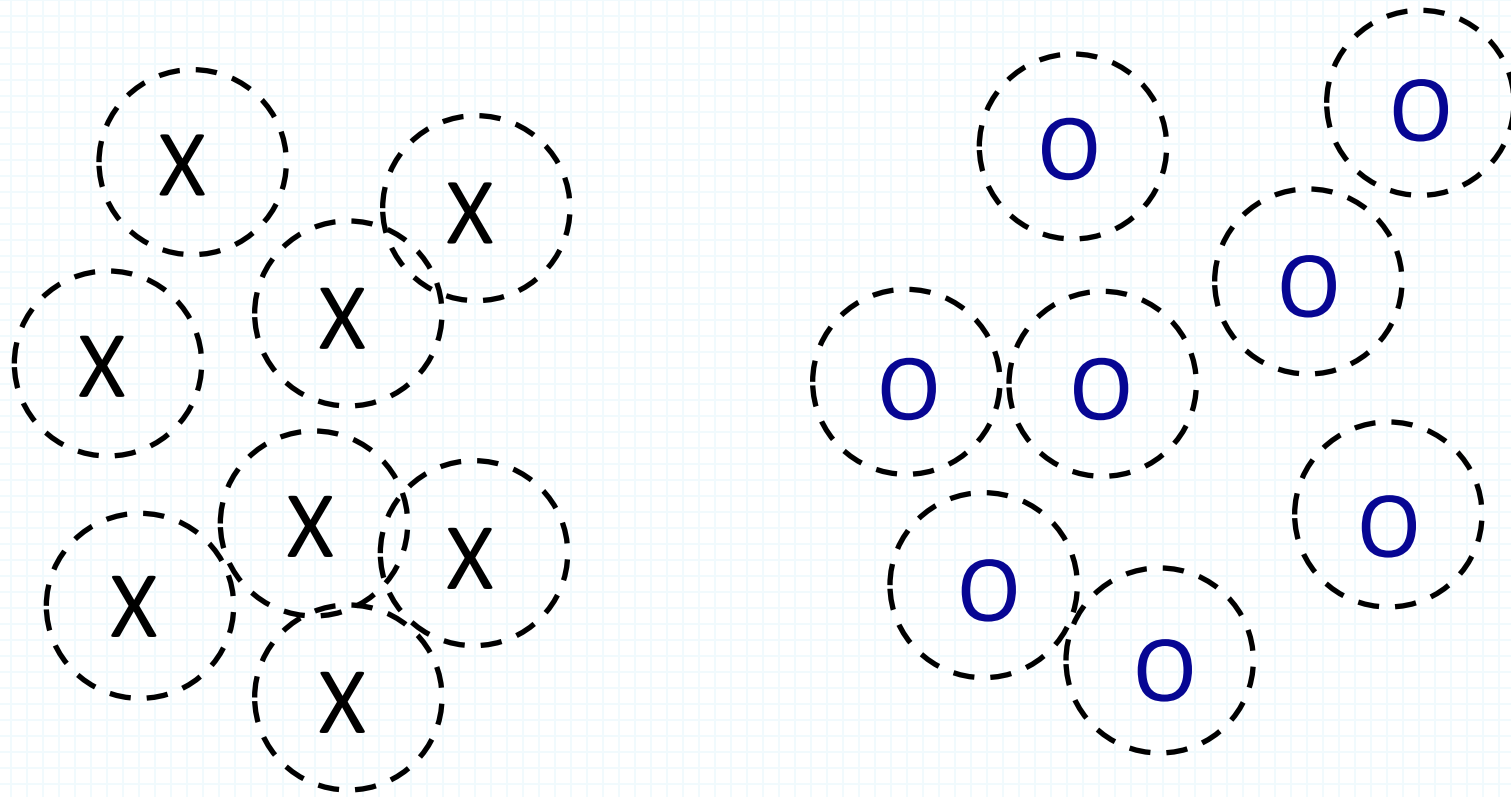
# Why SVM? ... Why not decision trees?

- Can Decision Trees detect non-linear models?
  - Yes, Decision trees can detect non-linear relationships

- What type of boundaries can be detected using decision trees in each step?
  - The decision boundary in a Decision Tree is linear and perpendicular to one of the input dimensions, which means that it is limited to finding only axis-parallel splits.

- What if we have higher-dimensional feature space, more complex relationships between input features and target class?
  - In the higher-dimensional feature space, the decision boundary can take on a more complex shape, such as a curved or nonlinear boundary.
  - More problems when the relationship between the input features and the target variable is complex (ex: image classification, sentiment analysis, etc)
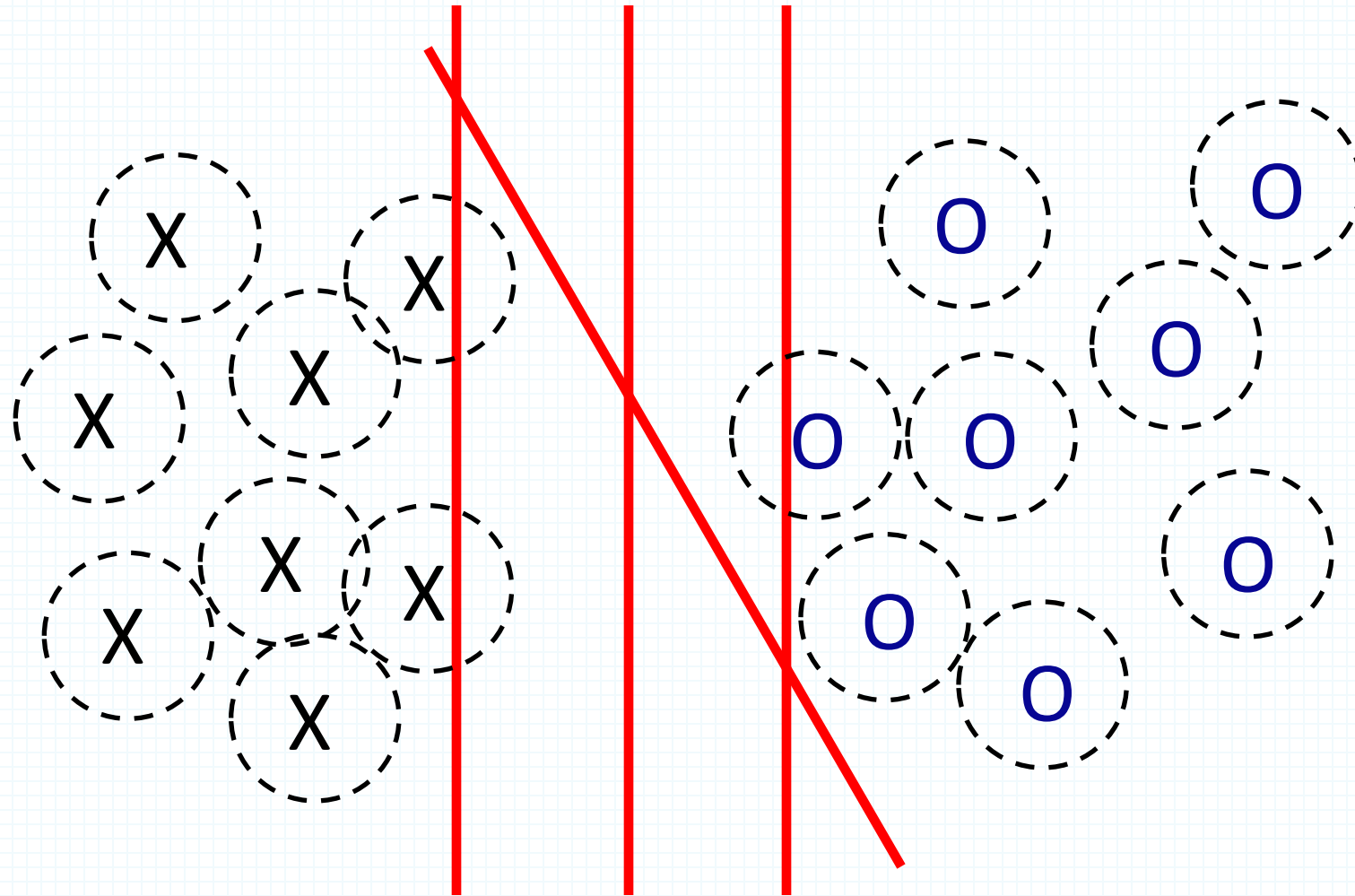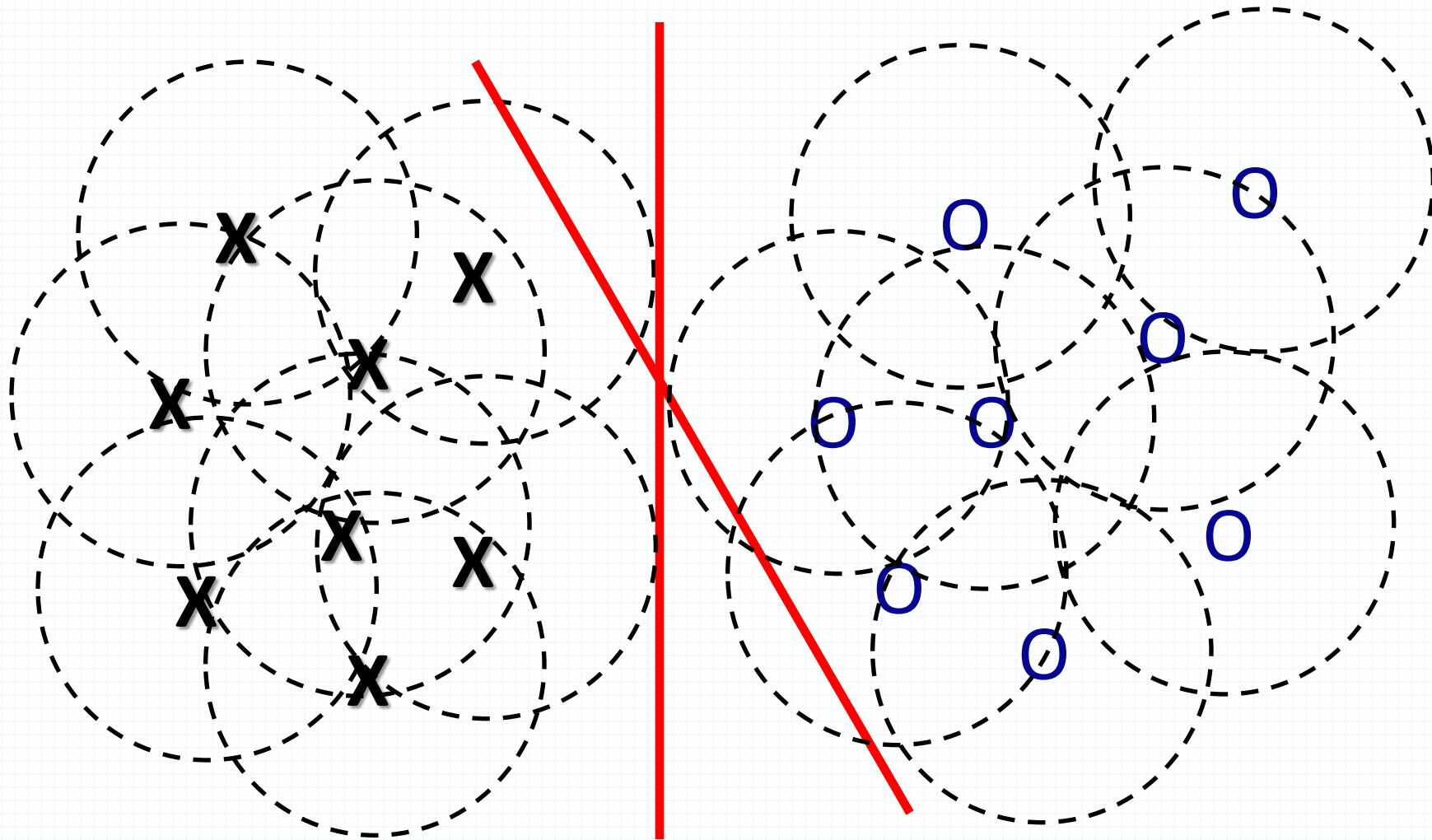
# A "Good" Separator

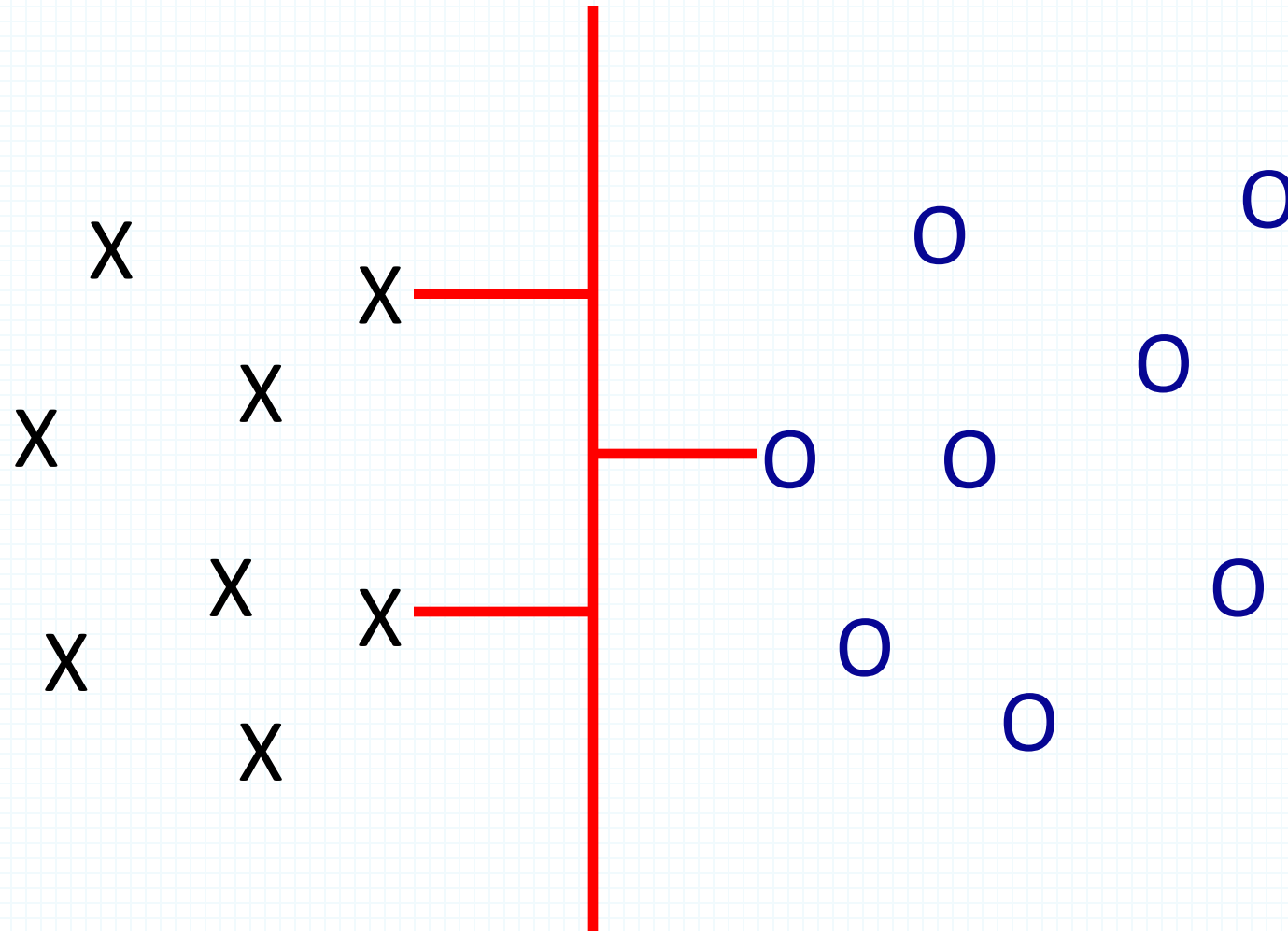# Noise in the Observations

# Ruling Out Some Separators

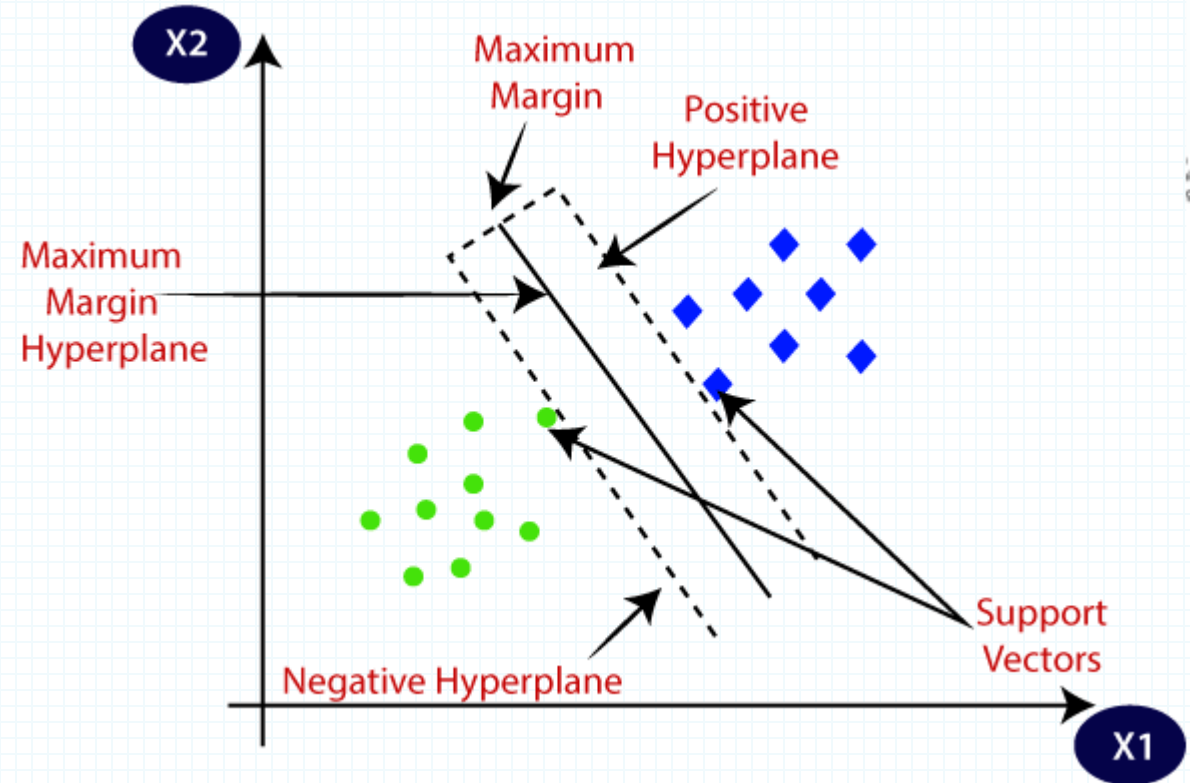# Lots of Noise

# Maximizing the Margin

# Terms

- **Support Vectors:**
  - These are the points that are closest to the hyperplane.
  - A separating line will be defined with the help of these data points.
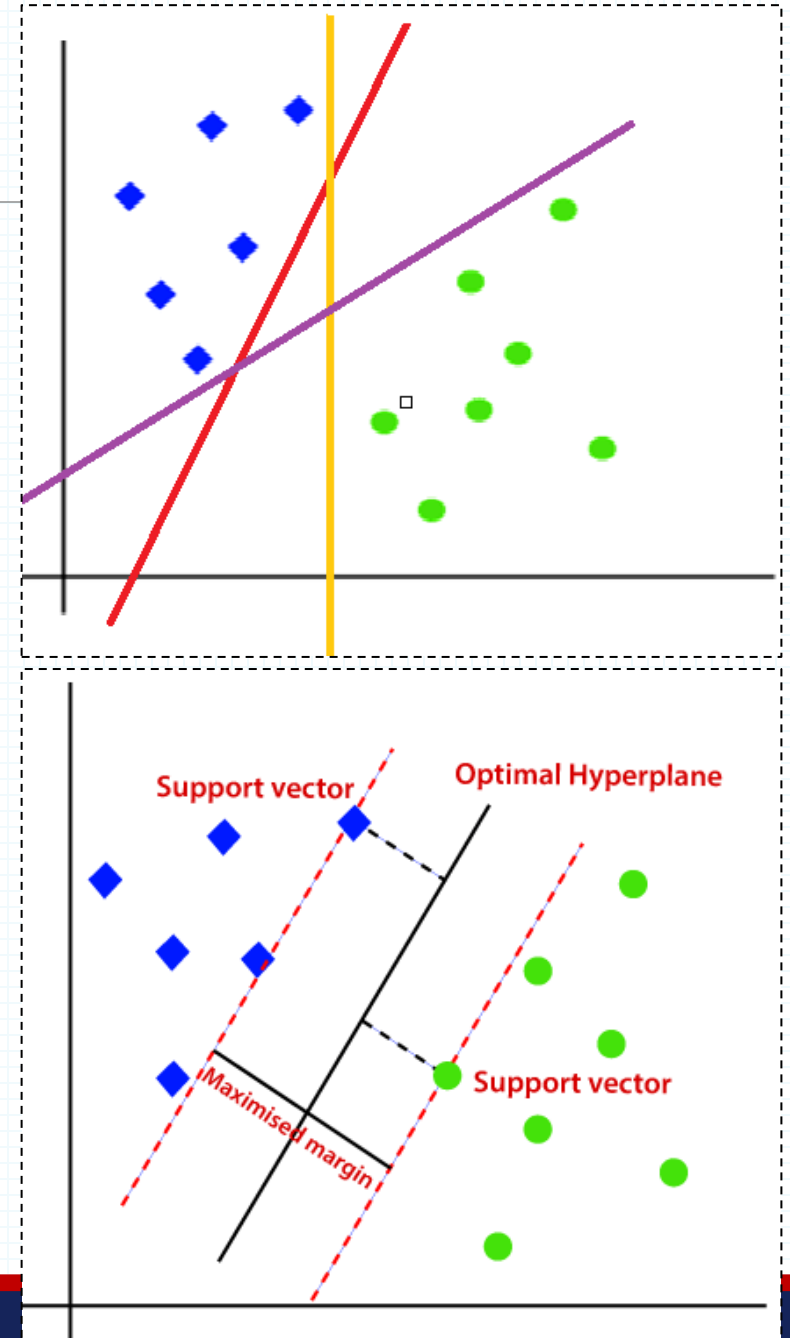
- **Margin:**
  - It is the distance between the hyperplane and the observations closest to the hyperplane (support vectors).
  - In SVM large margin is considered a good margin.
  - There are two types of margins **hard margin** and **soft margin.**

# How does SVM work?

- SVM is defined such that it is defined in terms of the support vectors only.
  - The margin is made using the <u>points</u> which are <u>closest</u> to the <u>hyperplane</u> (support vectors).
  - We don't have to worry about other observations
  - Hence SVM enjoys some natural <u>speed-ups</u>!
- The <u>best hyperplane </u>is that plane that has the maximum distance from both the classes.
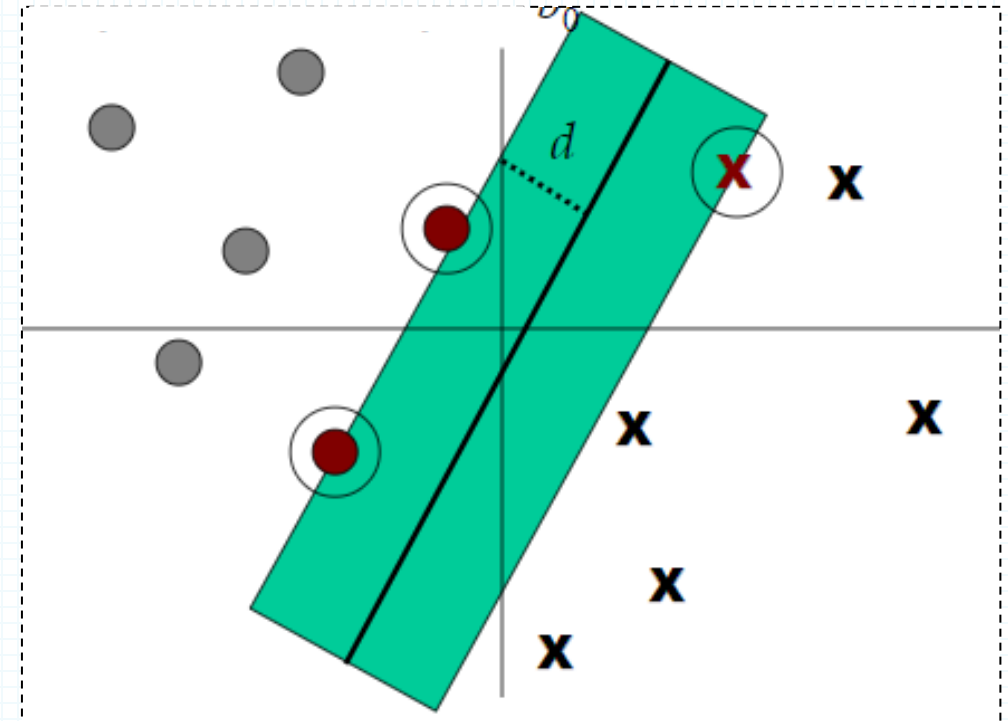
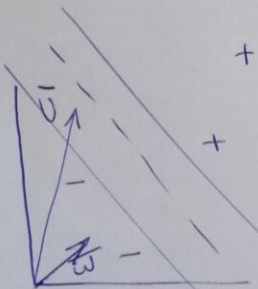# So.. What is our optimization problem?

- **Our problem:**

*Maximizing the shortest distance to the closest positive or negative point*

$$w^* = arg_w max \left[ min_n \, d_H(\phi(x_n)) \right]$$



Note that W represents all parameters i.e. w and b

$\overline{w} \cdot \overline{u} \geqslant C \qquad$ the $N^+$

$\overline{w} \cdot \overline{u} + b \geqslant \phi$

---

$\overline{w} \, x_+^* + b \geqslant 1$

$w \, x_- + b \leqslant -1$

---

$y_i = +1$ for $+$ samples

$\qquad -1$ for $-$ samples

constraint

$\boxed{y_i (\overline{w} \cdot \overline{x}_i + b) \geqslant 1}$ for all points

$\circledast \quad y_i (\overline{w}\,\overline{x}_i + b) = 1$ for support vector

---

$+$

$\text{width} = (\overline{x}_+ - \overline{x}_-) \cdot \dfrac{\overline{w}}{\|w\|}$

support vector $\circledast$

$= (\overline{x}_+ \overline{w} - \overline{x}_- \overline{w}) \dfrac{1}{\|w\|}$

---

$N^+$

$x_+ \rightarrow \quad \overline{w} x_+ + b = 1 \longrightarrow \overline{w} x_+ = 1 - b$

$x_- \rightarrow \quad \overline{w} x_- + b = -1 \Rightarrow -\overline{w} x_- = 1 + b$

$\Downarrow$

$\text{width} = \dfrac{2}{\|w\|}$

---

Goal $\quad$ Maximize width

$\dfrac{2}{\|w\|}$

$\Downarrow$

Minimize $\|w\|$

$\Downarrow$

Goal $\quad \boxed{\text{Minimize} \ \dfrac{1}{2} \|w\|^2}$

$\uparrow$

true only if the constraint is satisfied

$\Downarrow$

use lagrange Multiplier

$L = \dfrac{1}{2} \|\overline{w}\|^2 - \sum \alpha_i \left[ y_i (\overline{w}\,\overline{x}_i + b) - 1 \right]$

w.r.t

# SVM Optimization

$$w^*, b^* = \arg \underset{w,b}{Min} \frac{1}{2} \|w\|^2 \ , \quad s.t. \quad y_n \left(w^T\big(\emptyset(x_n)\big) + b\right) \geq 1 \quad \forall n$$

Solved by Lagrange multiplier method:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 \ - \sum_n \alpha_n [y_n\big(w^T\big(\emptyset(x_n)\big) + b\big) - 1]$$

where $\boldsymbol{\alpha}$ is the Lagrange multiplier

The optimization problem can be solved by setting **derivatives** of *Lagrangian* to 0

# SVM Optimization

$$L(w, b, \alpha) = \frac{1}{2}\|w\|^2 - \sum_n \alpha_n[y_n(w^T(\emptyset(x_n)) + b) - 1]$$

$$\frac{\partial L}{\partial w} = w - \sum_n \alpha_n y_n \emptyset(x_n) = 0 \Rightarrow w = \sum_n \alpha_n y_n \emptyset(x_n)$$

$$\frac{\partial L}{\partial b} = \sum_n \alpha_n y_n = 0 \Rightarrow \sum_n \alpha_n y_n = 0$$
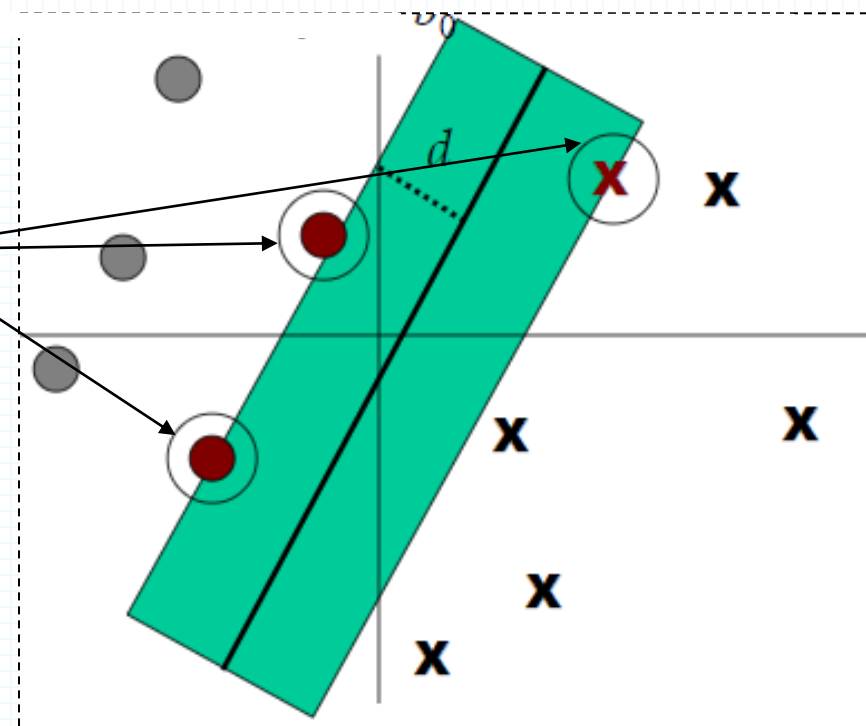
# SVM Optimization

$$w^*, b^* = \arg \underset{w,b}{Min} \frac{1}{2} \|w\|^2 , \quad s.t. \quad y_n \left(w^T\left(\emptyset(x_n)\right) + b\right) \geq 1 \quad \forall n$$

Y = $w^T\left(\emptyset(x)\right) + b$ = $\sum_n \alpha_n y_n \emptyset^T(x_n)\, \emptyset(x)$

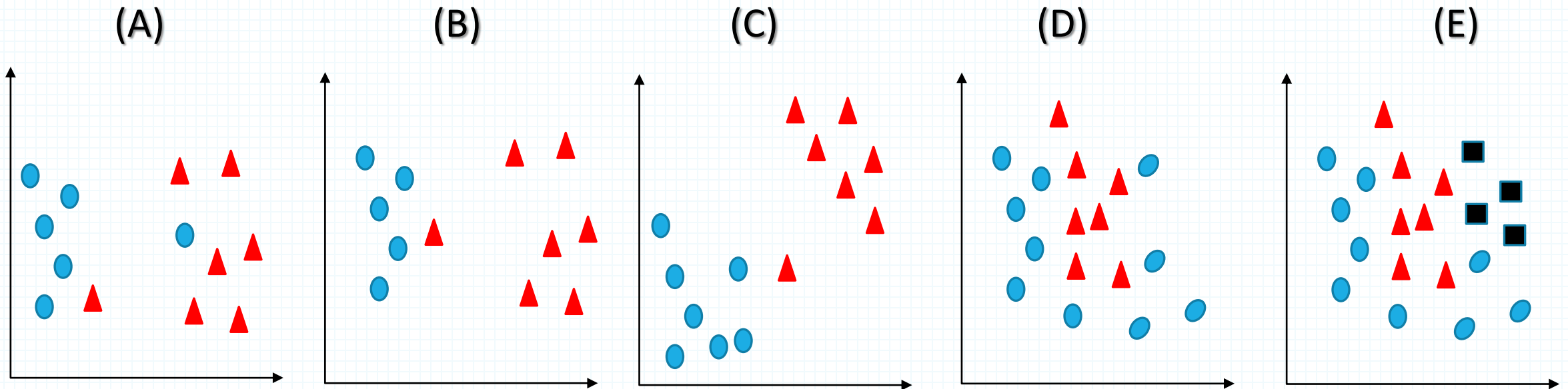The decision rule in SVMs only depends on the dot product with support vectors

several important implications

Computational efficiency
Memory efficiency
Robustness to noise and outliers

# What if?

What are the problems of the current version for SVM?

(A)    (B)    (C)    (D)    (E)

# 1ˢᵗ Improvement
## Soft Margin SVM
## (allows few misclassifications)

# C Hyper-parameter

- When **C** is <u>high</u> it will <u>classify all the data points correctly</u>, also there is a chance to overfit.

$$\mathrm{argmin}\left(\mathrm{w}^*, \mathrm{b}^*\right) \frac{\|\mathrm{w}\|}{2} + c \sum_{i=1}^{n} \zeta_i$$

- ***SVM Error = Margin Error + Classification Error***