# Pattern Recognition

## SVM 개념 잡기
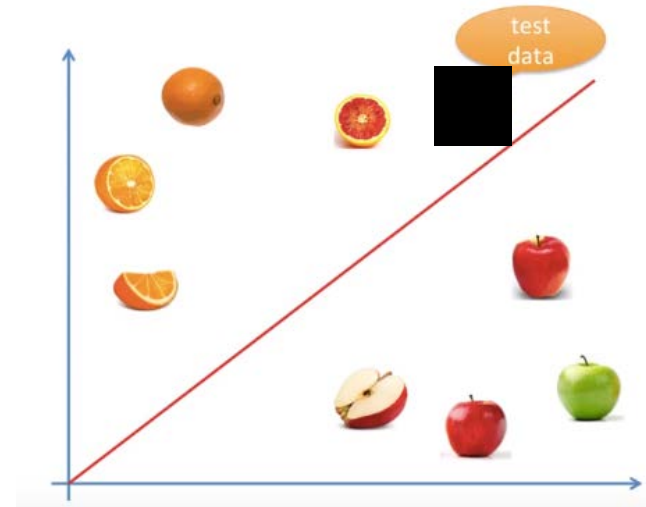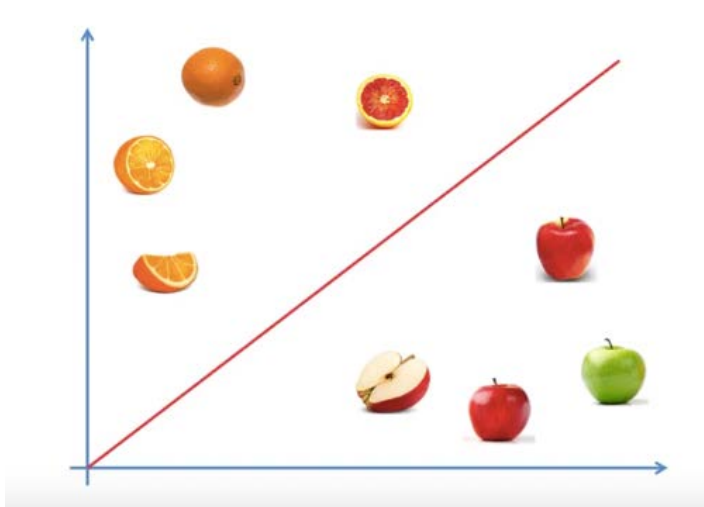
Yukyung Choi

yk.choi@rcv.sejong.ac.kr

# What is SVM?

- **S**upport **V**ector **M**achine ➔ SVM

- Traditional Classifier

- Until now, favorite classifier to everyone
  - Wondering why? **<u>Kernel Trick!!!</u>**

<span style="color:red">**"만약, 문제에 어떠한 알고리즘을 사용할지 모르겠다면, SVM은 좋은 출발선이 될 수 있음"**</span>
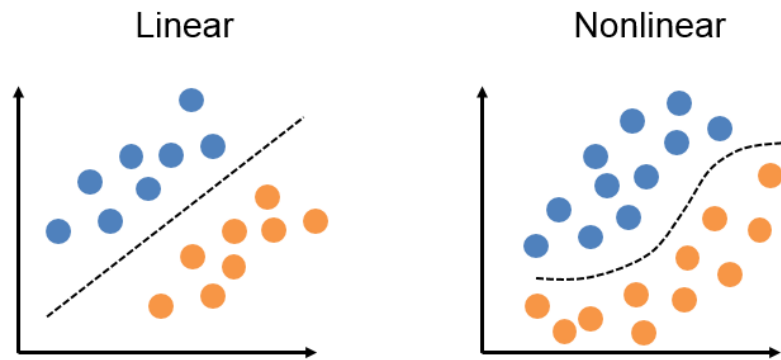
# Classifier

- **Classifier** is a <u>hypothesis</u> or <u>discrete-valued function</u> that is used to **assign (categorical) class labels to particular data points**.

- In the email classification example, this classifier could be a hypothesis for labeling emails as **spam** or **non-spam**.

# Classifier

- y = label, x = data, y = f(x), f: classifier

- If **decision function** is linear, this classifier (f) is linear classifier
- If not, this classifier (f) is non-linear classifier

Linear                    Nonlinear

y = f(x)

데이터를 구획해주는 이 **점선의 함수 (decision boundary)**를 우리는 **판별 함수 (decision function)**라 부른다.
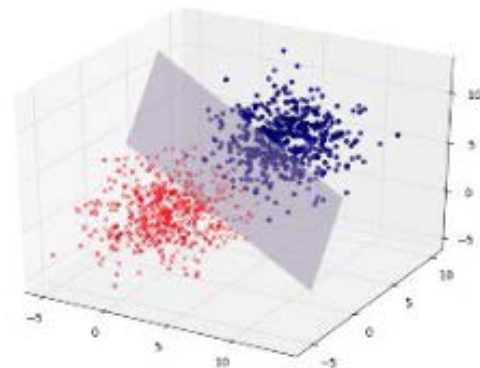
# Classifier

- **Hyperplane**
  - In geometry, a hyperplane is a subspace whose dimension is one less than that of its ambient space. If a space is **3-dimensional** then its hyperplanes are the **2-dimensional planes**, while if the space is **2-dimensional**, its hyperplanes are the **1-dimensional lines**.
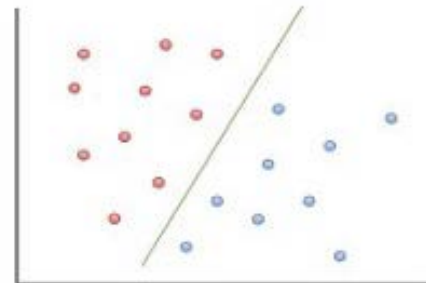
$$\mathbf{w}^T \mathbf{x} = 0 \qquad\qquad y = ax + b$$
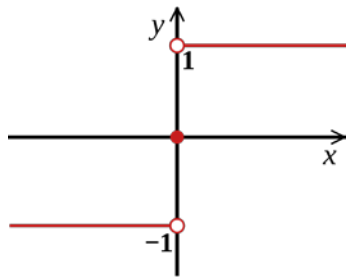
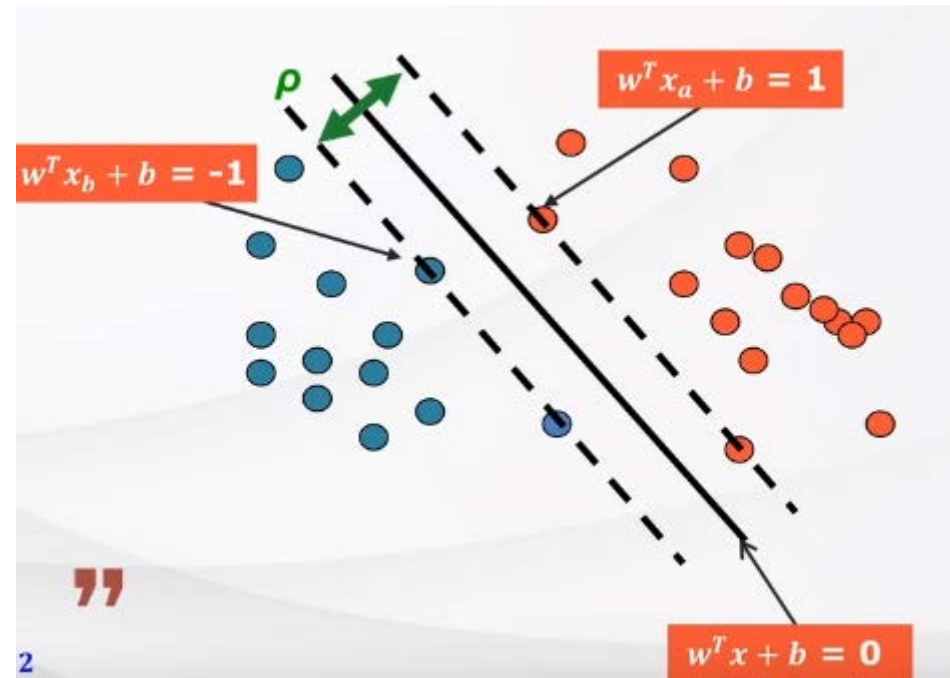Hyperplane                           Line

# SVM Classifier

- W : vector for hyperplane
- $x_i$ : $i_{th}$ data, $y_i$ : label (class) of $i_{th}$ data

- **Y** = **sign**$(W^T X + b)$ = $f(X)$
  - $Y_i$ = +1 when $W^T X_i + b > 1$
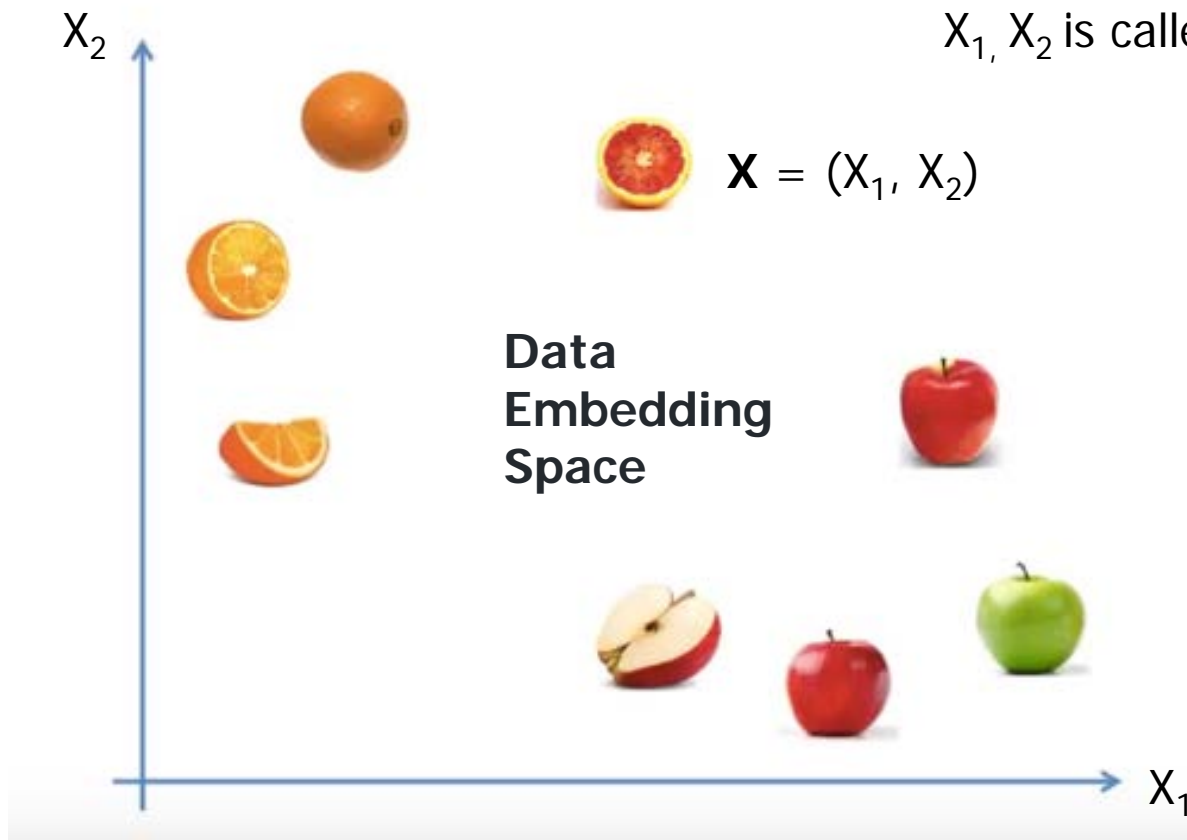  - $Y_i$ = -1 when $W^T X_i + b < -1$



Sign function

# Apple, orange classifier

- **Data Embedding Space**

$X_1, X_2$ is called feature or attribute.

$X_2$

$\mathbf{X} = (X_1, X_2)$

Data
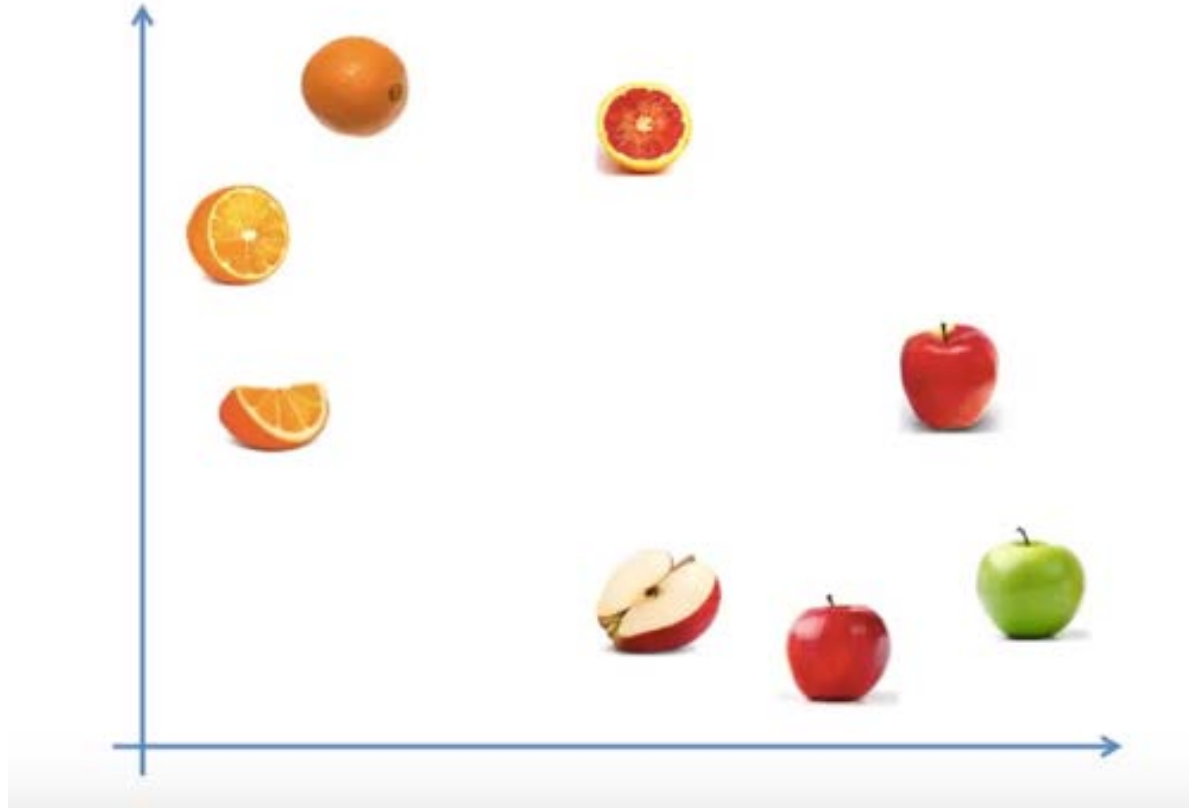Embedding
Space

$X_1$

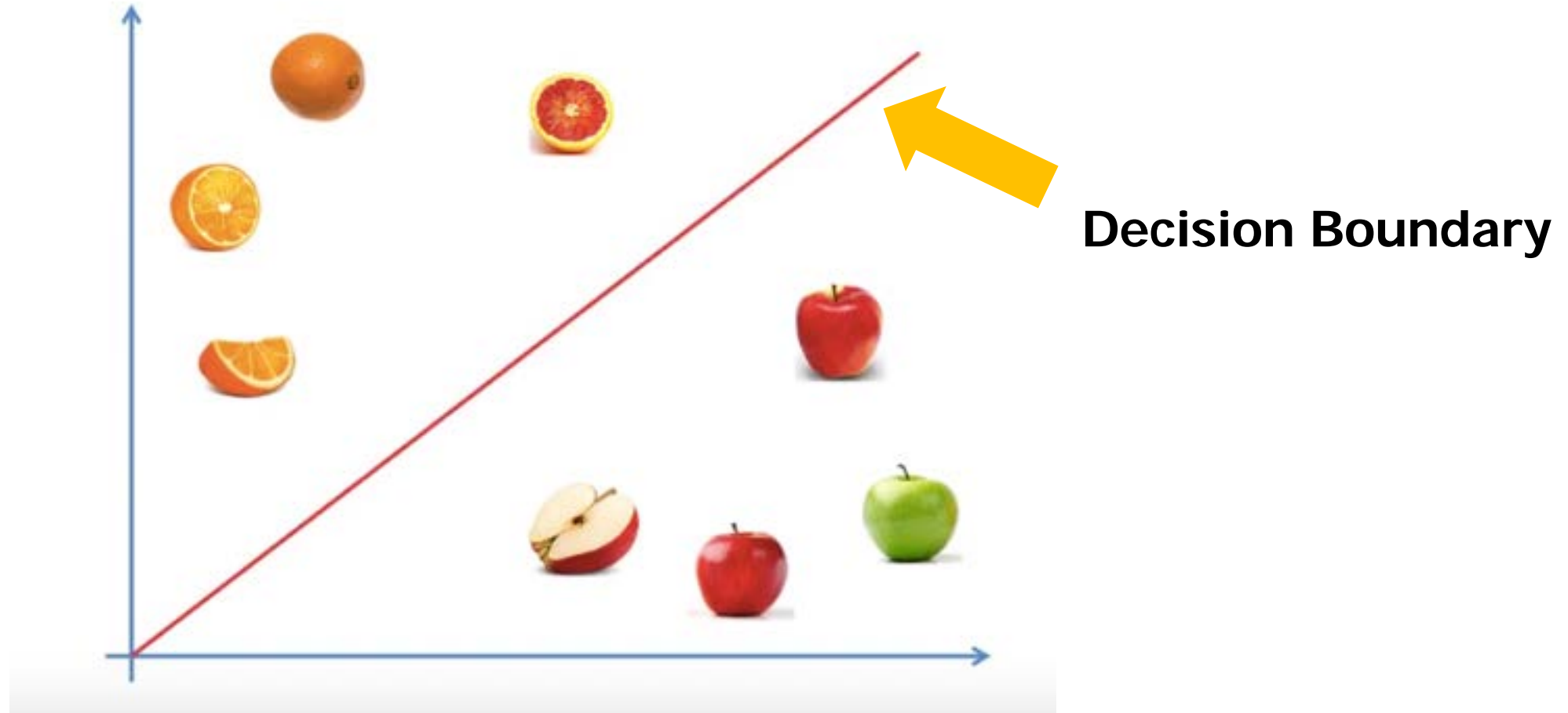**Data Embedding**
범주형 자료를 **벡터 형태**로 바꾸는 것

**Categorical Data**
**범주형 데이터**란 몇 개의 범주로 나누어진
데이터 예) 남/여, A/B/O/AB

# Apple, orange classifier

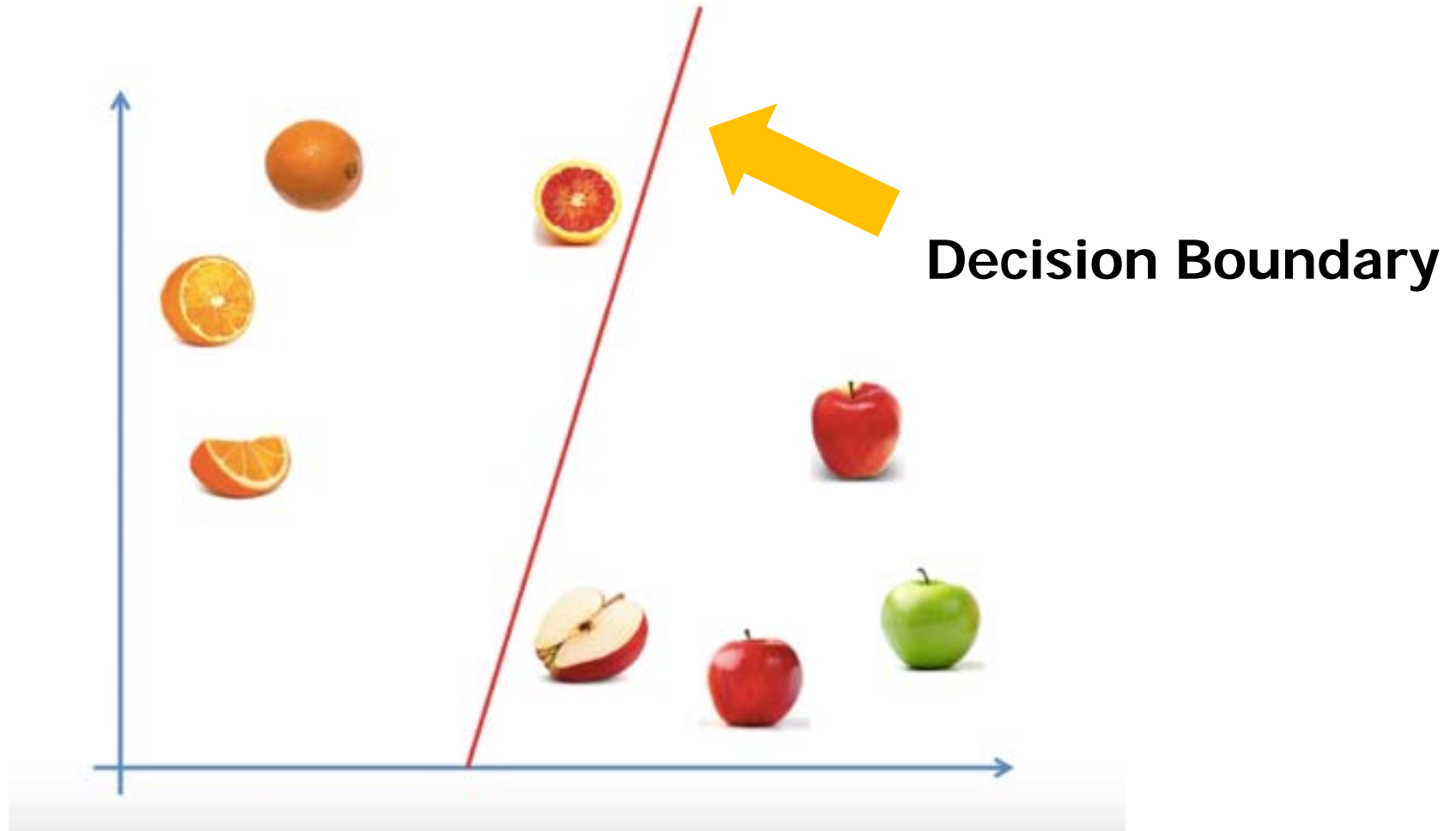- Which hyperplane can we choose?

# Apple, orange classifier
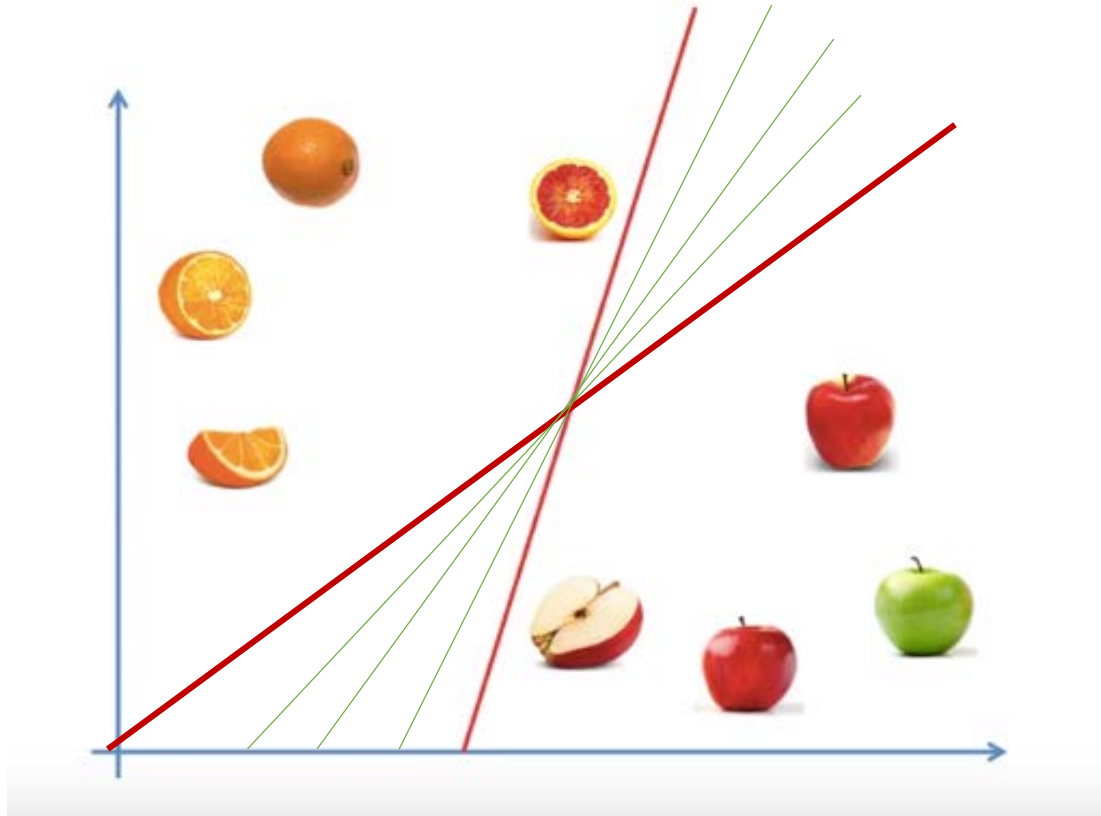


**Decision Boundary**

# Apple, orange classifier
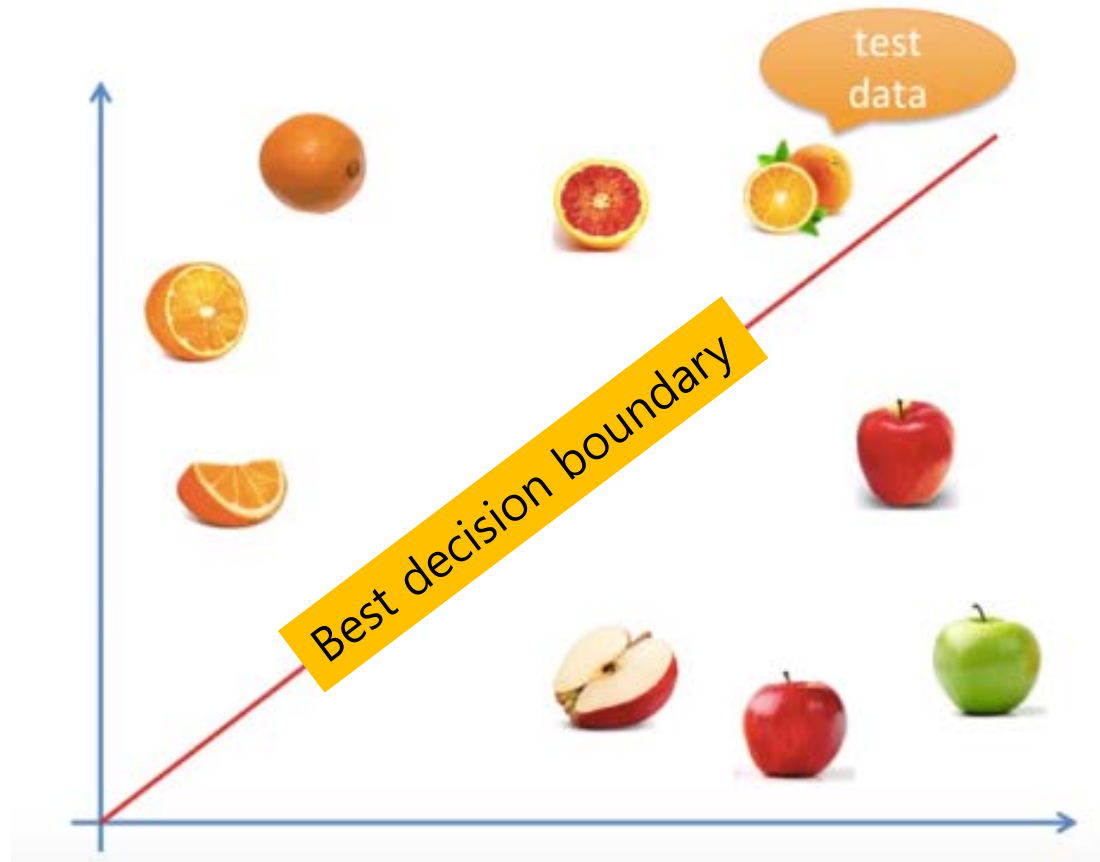


Decision Boundary

# Apple, orange classifier

- Which one is better?
  - Classifier should have dealt with **unseen data**

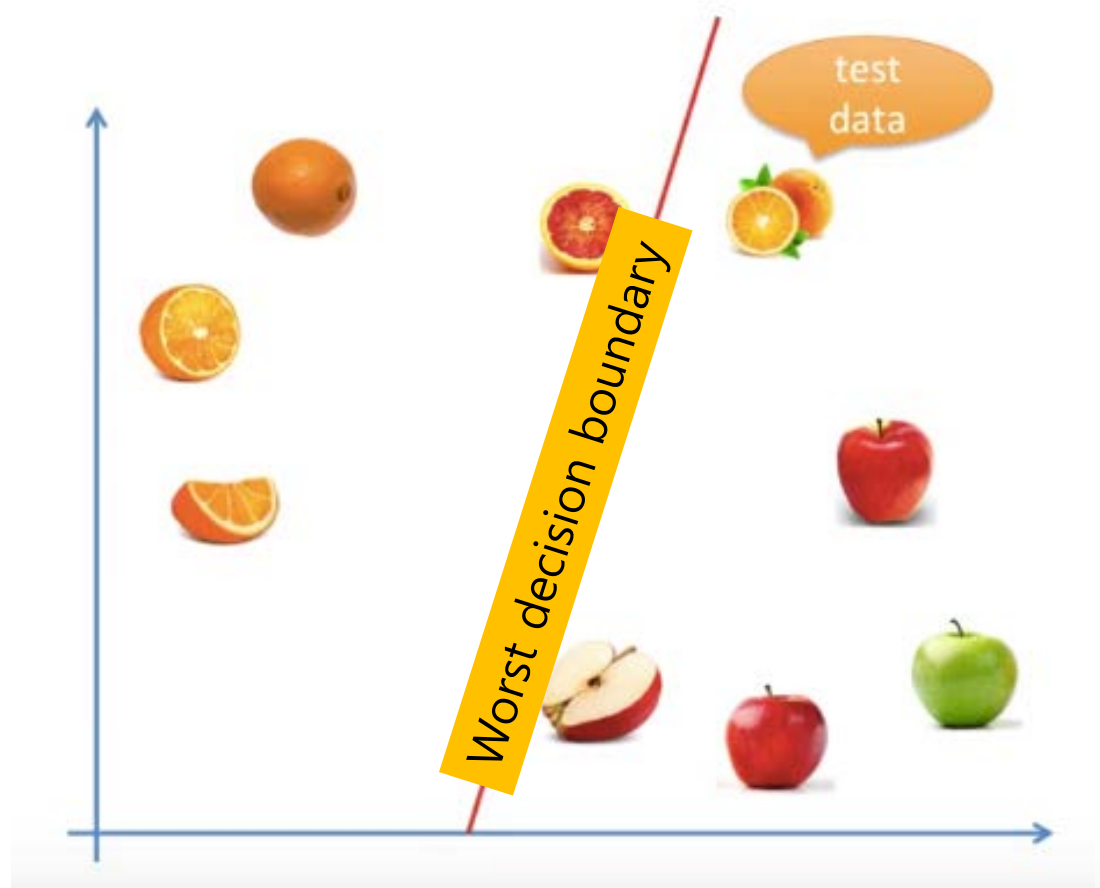Train sample data ➔ seen data
Test sample data ➔ unseen data

# How can we decide decision boundary?

- Test data predicted well (O)

# How can we decide decision boundary?

- Test data predicted well (X)
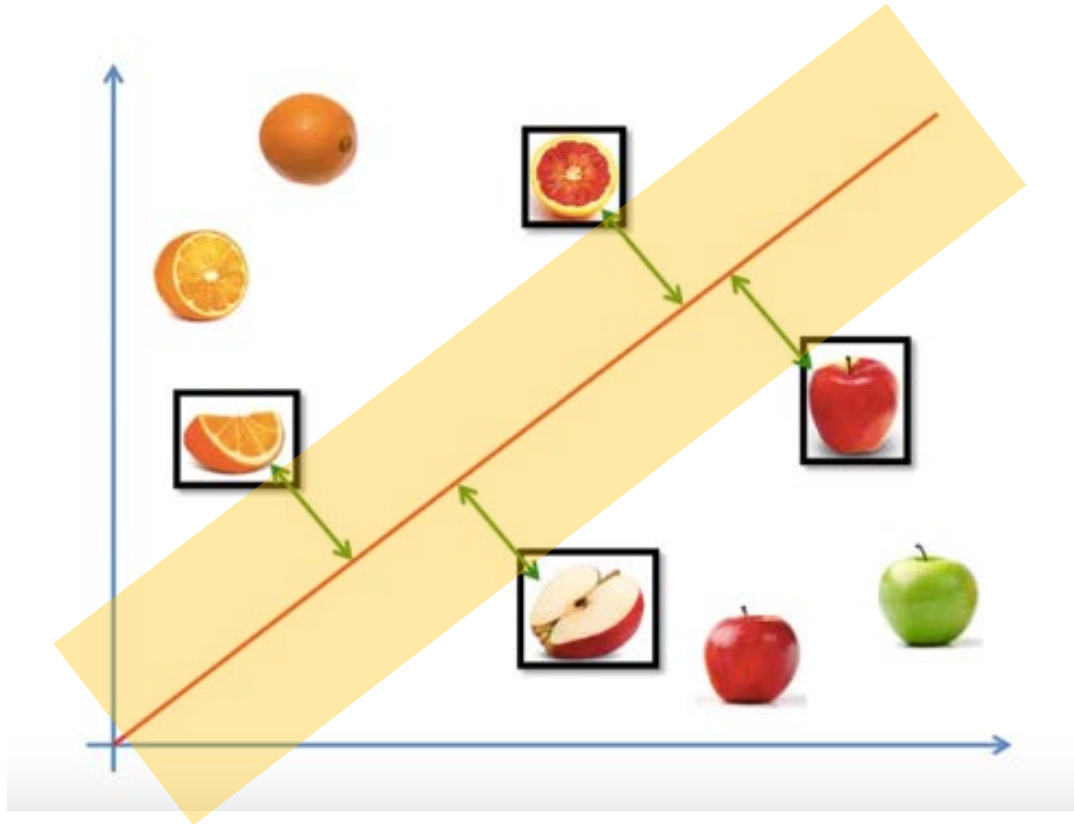
# How can we decide decision boundary?

- The answer is "Large Margin"!!

# Support Vector

- **Support Vector**
  - Samples on the margin are called the support vectors.

# Support Vector

- SVM only uses support vector for prediction
  - Less computation!!!

# Linearly Separable or not

# What if data is not linearly separable?



**Data Embedding Space ➔ 1D**

**Hyperplane ➔ 0D ➔ y = b**

# What if data is not linearly separable?

- Mapping <u>lower dimension</u> to <u>high dimension</u>



$$y = x^2$$

**Data Embedding Space ➜ 2D**

**Hyperplane ➜ 1D ➜ y = ax+b**

# What if data is not linearly separable?

- Now it is linearly separable in higher dimension
  - Mapping to high dimension requires **much computation!**

# What if data is not linearly separable?

- **Kernel trick** in SVM do this without explicitly
  - Move data point to higher dimension with **low computation!**

# Kernel Trick

- The **kernel trick** avoids the explicit mapping that is needed to get linear learning algorithms.

- **Kernel methods** owe their name to the use of kernel functions, which enable them to operate in a high-dimensional, implicit feature space without ever computing the coordinates of the data in that space, but rather by **simply computing the inner products** between the images of all pairs of data in the feature space

# Kernel Trick

- Kernel Function ➜ **simply computing the inner products**

The *kernel function* can be any of the following:

- linear: $\langle x, x' \rangle$.
- polynomial: $(\gamma \langle x, x' \rangle + r)^d$. $d$ is specified by keyword `degree`, $r$ by `coef0`.
- rbf: $\exp(-\gamma \|x - x'\|^2)$. $\gamma$ is specified by keyword `gamma`, must be greater than 0.
- sigmoid $(\tanh(\gamma \langle x, x' \rangle + r))$, where $r$ is specified by `coef0`.

Mapping 함수의 inner-product.. Mapping (m➜n)

$$K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j) = x_i^T A^T A x_j$$

# SVM Parameter - Cost

- Cost is small == Margin is large



train error

**C** is small

Training error is allowed

Overfitting is not allowed

Margin is large

Testing error is small

$$\min_{w,b,\zeta} \frac{1}{2} w^T w + C \sum_{i=1}^{n} \zeta_i \qquad \text{subject to } y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i,$$
$$\zeta_i \geq 0, i = 1, \ldots, n$$

Margin width          misclassification

# SVM Parameter - Cost

- Cost is large == Margin is small



$$\min_{w,b,\zeta} \frac{1}{2} w^T w + C \sum_{i=1}^{n} \zeta_i \qquad \text{subject to } y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i,$$
$$\zeta_i \geq 0, i = 1, \ldots, n$$

Margin width       misclassification

**C** is large
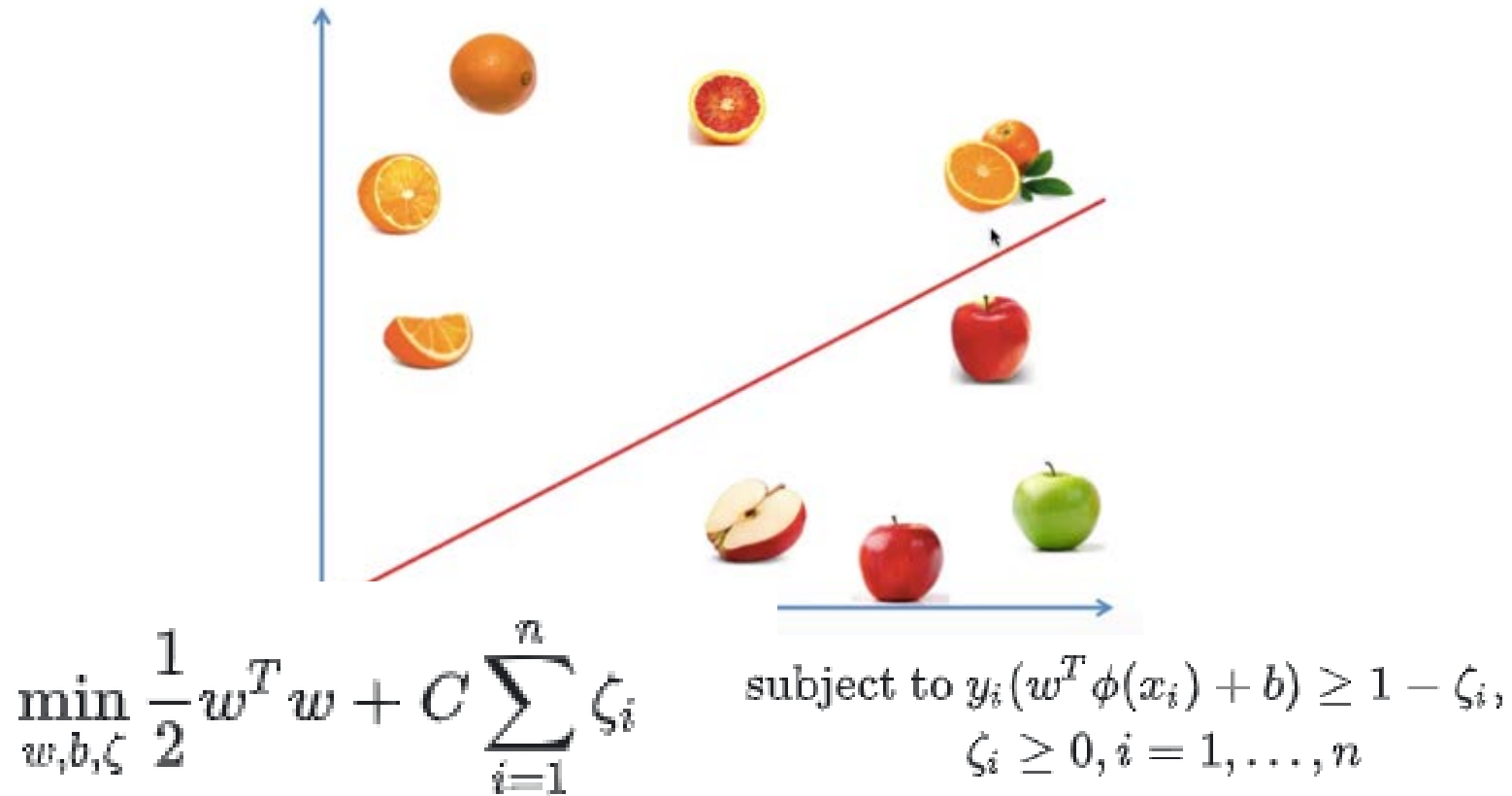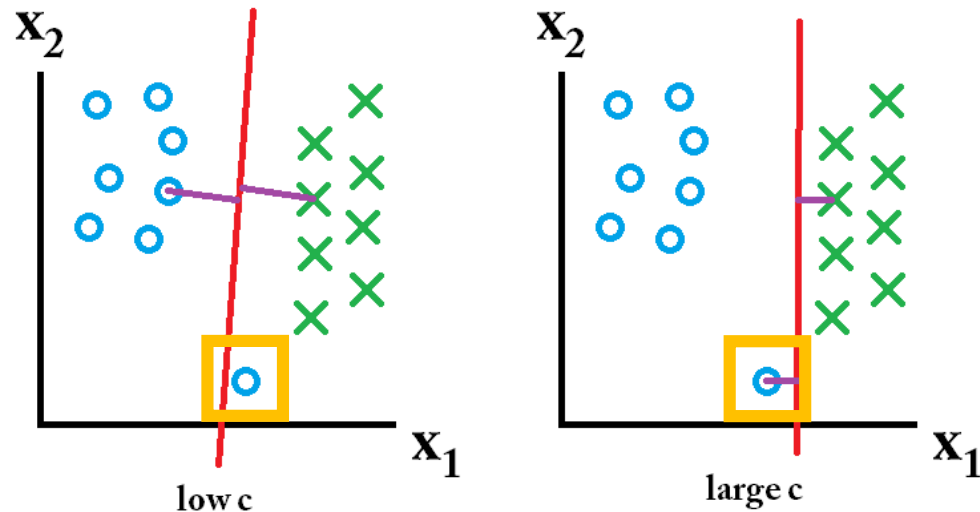
Training error is not allowed

Overfitting is allowed
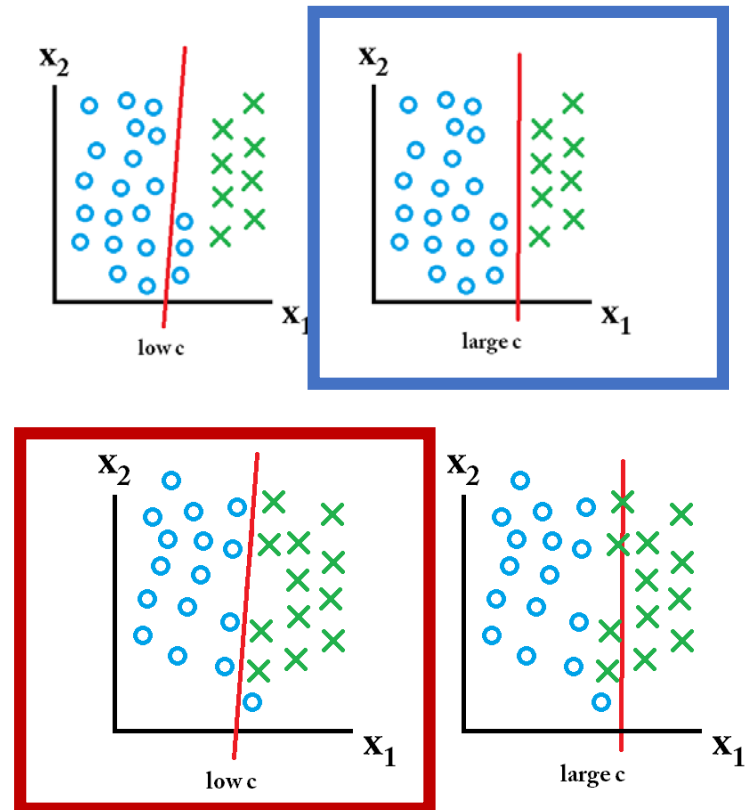
Margin is small

Testing error is large

# SVM Parameter - Cost

- We **assume** that some samples caused by train error are the **outlier**.
- Therefore, we generally select a **large margin** for decision boundary.
- But, if not?

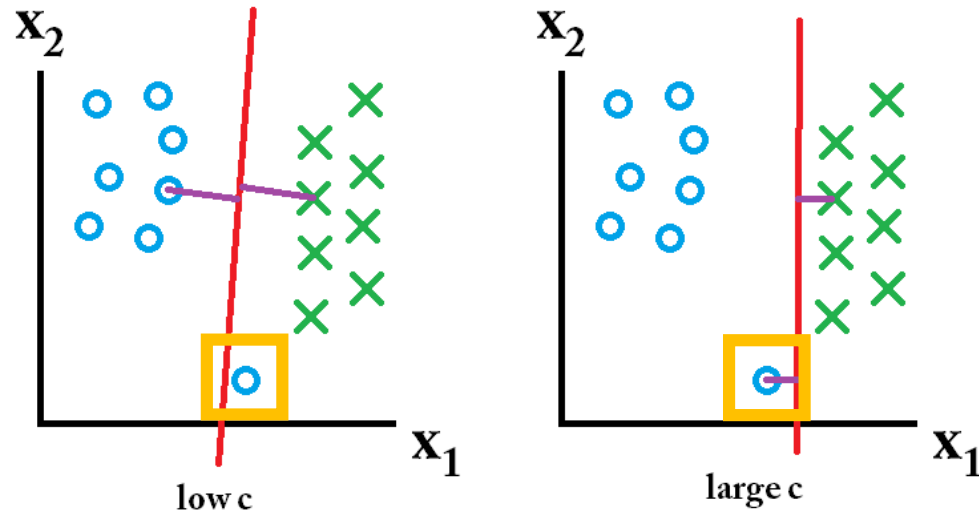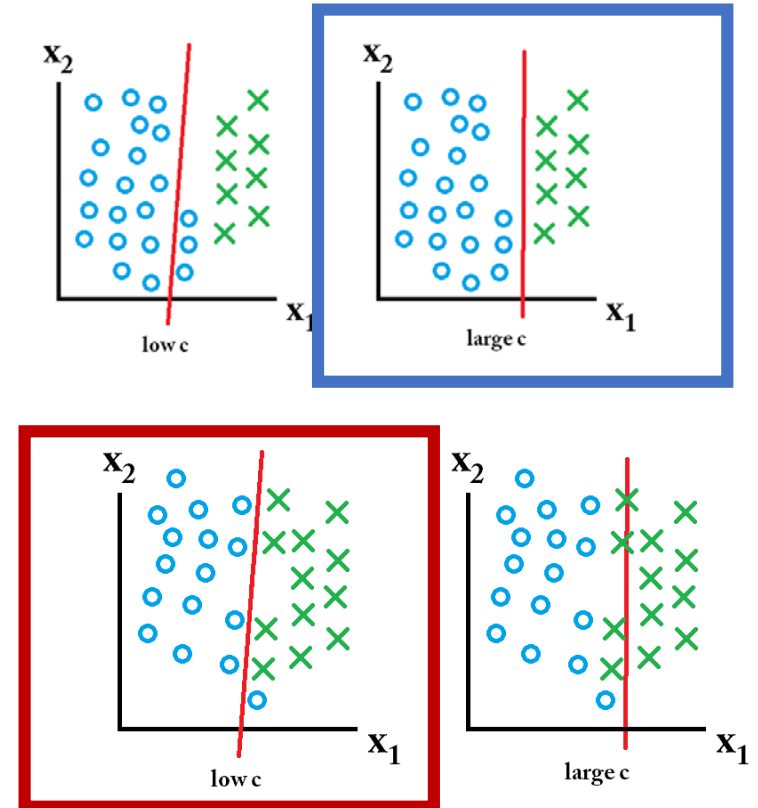# SVM Parameter - Cost

- Therefore, we cannot argue that we should choose large C, but we must make a decision through **data analysis**.



Inlier or outlier

# SVM Parameter - Cost



cost= 0.01

Cost is small
Training error is allowed
Overfitting is not allowed
**Decision boundary is simple**

cost= 1

Cost is large
Training error is allowed
Overfitting is allowed
**Decision boundary is complex**

# SVM parameter – Gamma in RBF kernel

- Intuitively, the **gamma parameter** defines how far the influence of a single training example reaches, with low values meaning 'far' and high values meaning 'close'.

Radial Base Function (also called Gaussian Kernel)

```
K(x,x') = exp(-gamma * ||x-x'||^2)
```



Far

Close

gamma = 0.01

gamma = 1

gamma = 100

# SVM parameter – Gamma in RBF kernel



Far

Close

gamma = 0.01

gamma = 1

gamma = 100

Gamma is small
Influence is large
Margin is large
**Similarly to a linear model**

Gamma is large
Influence is small
Margin is small
Overfitting is allowed

# Find optical parameter – **data analysis**

- **Grid Search**
  - **Grid search** builds a model for **every combination** of hyper-parameters specified and evaluates each model.

|  | 1 | 10 | 100 |
|---|---|---|---|
| 1 | 0.7 | 0.8 | 0.7 |
| 10 | 0.8 | 0.8 | 0.9 |
| 100 | 0.6 | 0.8 | 0.8 |

Gamma

cost

# Binary Classification

- One vs One

- One vs Rest

# Multiple Classification

- How to extend binary to multiple classifier

# Unbalanced problems

- Sklearn: class_weight

# How to design custom kernel

- [https://scikit-learn.org/stable/auto_examples/svm/plot_custom_kernel.html#sphx-glr-auto-examples-svm-plot-custom-kernel-py](https://scikit-learn.org/stable/auto_examples/svm/plot_custom_kernel.html#sphx-glr-auto-examples-svm-plot-custom-kernel-py)

# SVM Optimization



최적화 문제를 사용한 파라미터 계산 (1/3)

- 서포트 벡터 머신의 파라미터를 찾기 위해서 최적화 문제로 변형시킬 수 있

  **Find w and b such that**

  $1 / \|w\|$ **is maximized; and for all** $\{(x_i, y_i)\}$

  $w^T x_i + b \geq 1$ **if** $y_i = 1$ **;** $w^T x_i + b \leq -1$ **if** $y_i = -1$

- 보다 나은 형식으로 변형 (min $\|w\|$ = max $1/\|w\|$ )

  **Find w and b such that**

  $\Phi(w) = \frac{1}{2} w^T w$ **is minimized;**

  **and for all** $\{(x_i, y_i)\} : y_i (w^T x_i + b) \geq 1$

https://www.youtube.com/watch?v=POqtUhBhiP8

# 최적화 문제를 사용한 파라미터 계산 (2/3)

> **Find w and b such that**
> $\Phi(w) = \frac{1}{2} w^T w$ **is minimized** ;
> and for all $\{(x_i, y_i)\}$ ; $y_i (w^T x_i + b) \geq 1$

- 선형 조건에 부합하도록 이차함수를 최적화 시키는 문제
- 이차함수의 최적화 문제는 수학적 프로그래밍 문제에서 잘 알려진 분야로, 해결할 수 있는 많은 알고리즘이 존재함
- Lagrangian multiplier $\alpha_i$ 을 사용하여 다음의 primal과 dual problem으로 변형 가능

> **Maximize**
> $L(w,b) = 1/2 w^T w - \Sigma \alpha_i \{ y_i(w^T x_i - b) - 1 \}$
> (1) $\alpha_i \geq 0$ for all $\alpha_i$

> **Find** $\alpha_1 \ldots \alpha_N$ **such that**
> $Q(\alpha) = \Sigma \alpha_i - \frac{1}{2} \Sigma\Sigma \alpha_i \alpha_j y_i y_j x_i^T x_j$ **is maximized and**
> (1) $\Sigma \alpha_i y_i = 0$
> (2) $\alpha_i \geq 0$ for all $\alpha_i$

# 최적화 문제를 사용한 파라미터 계산 (3/3)

**솔루션은 다음과 같은 형식을 가짐**

$$W = \sum \alpha_i y_i X_i \quad b = y_k - w^T X_k \,, \, k \text{는 } \alpha_k \neq 0 \text{ 을 만족}$$

- 0이 아닌 $\alpha_i$는 해당하는 $x_i$가 서포트 벡터임을 의미

**그러므로 분류함수는 다음과 같은 형식임**

$$f(\mathbf{X}) = \sum \alpha_i y_i X_i^T X + b$$

- 분류는 새로운 테스트 데이터 x와 서포트 벡터 $x_i$의 내적에 의해 계산됨

" 하지만, 모델의 훈련 과정 때는
모든 훈련 데이터 쌍 $(x_i, x_j)$에 대해 내적 $x_i^T x_j$을 계산 "

$$W = \sum \alpha y x$$
$$b = y - w$$

# 소프트 마진 분류 (soft margin classification)

◉ 만약 훈련 데이터가 선형으로 분리되지 않을 경우,
슬랙 변수 $\xi_i$가 잘못 분류되거나 노이즈가 포함된 데이터에 추가됨

◉ 잘못 분류된 데이터 포인트를 본래 속하는 클래스로
비용을 들여 이동시켜줌

$$y_i(w^T x_i + b) \geq 1$$

$$min\|w\|$$

➡

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$min\|w\|+C\|\xi\|$$

◉ 모델의 학습 방법은 여전히 결정 영역을
각 클래스로부터 가장 멀리 위치하는 것임
(large margin)