

# Support vector machine

---

Ngô Minh Nhựt

2024

# Large margin classifier

---

## Part 1

# Logistic regression

---

## □ Hypothesis:

- $h_{\theta}(x) = g(\theta^T x)$
- $h_{\theta}(x) = P(y = 1|x; \theta)$

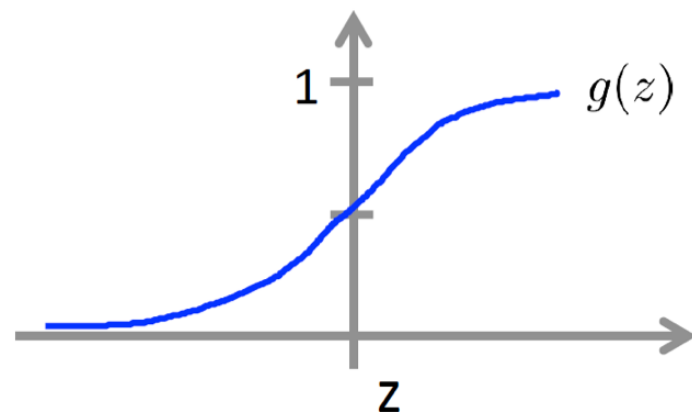
$$g(z) = \frac{1}{1 + e^{-z}}$$

## □ $\theta^T x \gg 0$

- Possibility  $y = 1$  is higher

## □ $\theta^T x \ll 0$

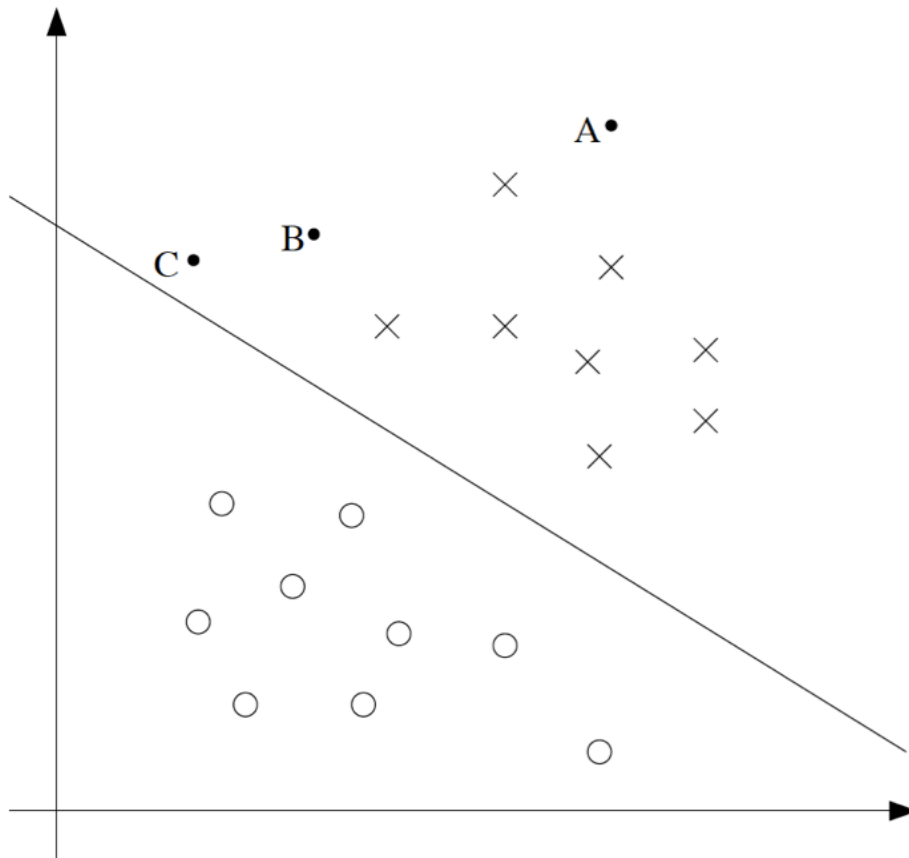
- Possibility  $y = 0$  is higher



# Logistic regression

---

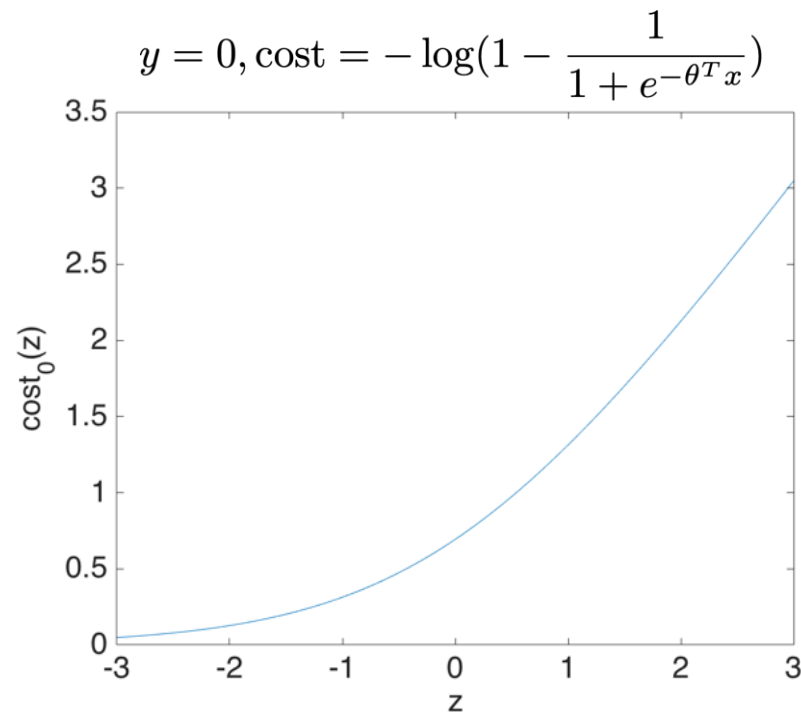
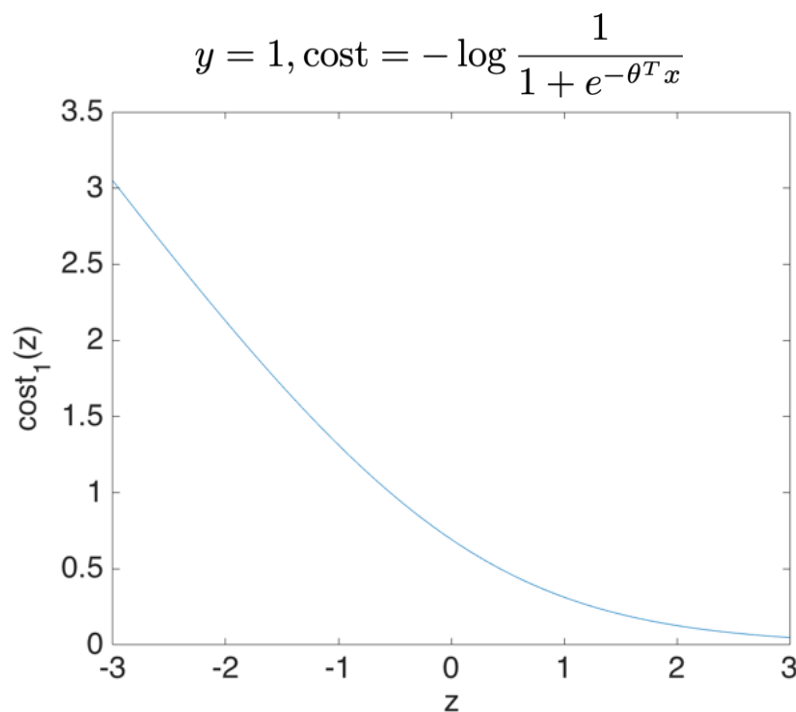
## □ Margin



# Logistic regression

## □ Cost function for a sample

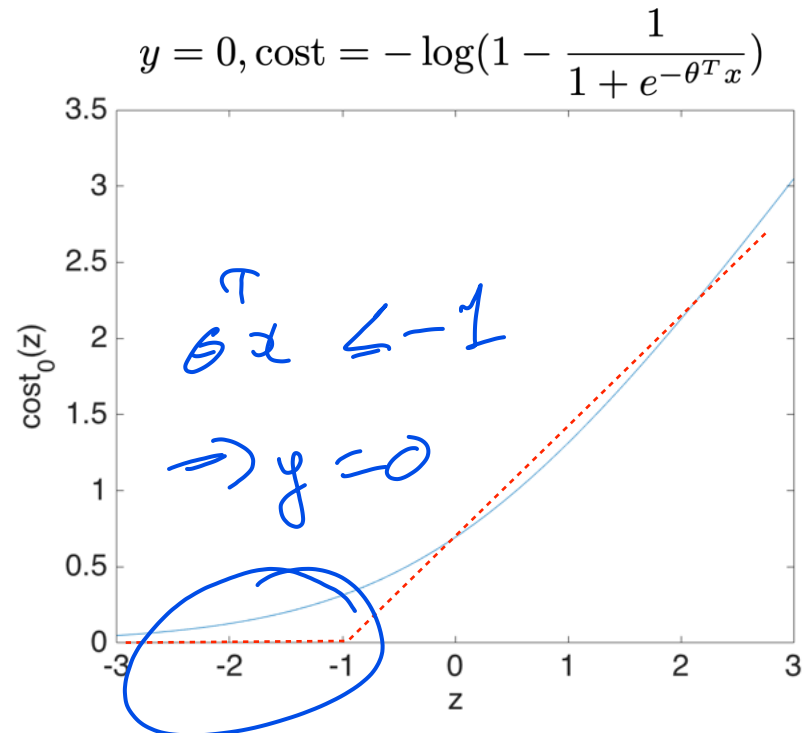
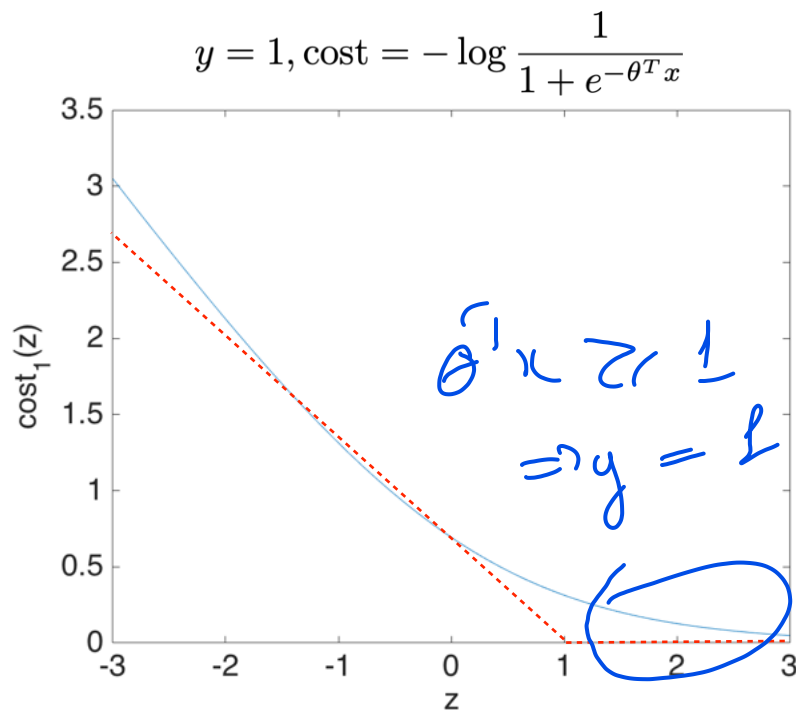
$$-y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left( 1 - \frac{1}{1 + e^{-\theta^T x}} \right)$$



# Logistic regression

## □ Cost function for a sample

$$-y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$



# Support vector machine

---

## □ Cost function

$$J(\theta) = C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

## □ In comparison with logistic regression

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

# SVM model

---

## □ Hypothesis

- $y = 1$  if  $\theta^T x \geq 1$
- $y = 0$  if  $\theta^T x \leq -1$

## □ Cost function

$$J(\theta) = C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

## □ Learning:

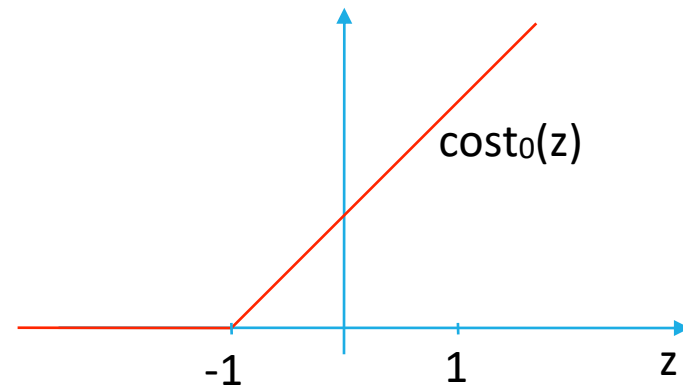
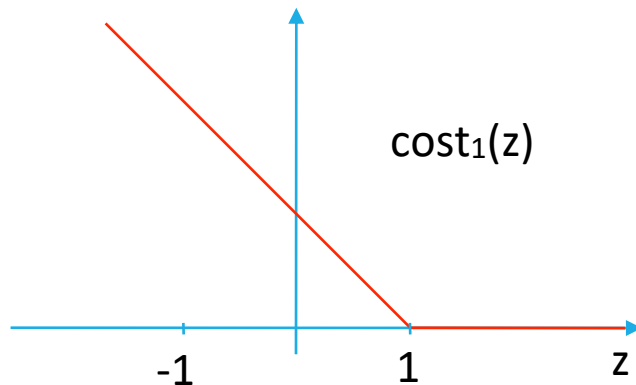
$$\min_{\theta} J(\theta)$$



# Learning

## □ Training:

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



- If  $y = 1$ , we need  $\theta^T x \geq 1$
- If  $y = 0$ , we need  $\theta^T x \leq -1$

# Learing

## □ Goal

$$\tau_1 \min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \overset{\tau_1}{\text{cost}_1}(\theta^T x^{(i)}) + (1 - y^{(i)}) \overset{\tau_0}{\text{cost}_0}(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

■ If  $y = 1$ , we need

$$\theta^T x \geq 1 \Rightarrow \tau_0 = 0$$

■ If  $y = 0$ , we need

$$\theta^T x \leq -1 \Rightarrow \tau_1 = 0$$

$$\rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

# Why large margin

---

- Dot product

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = u_1 v_1 + u_2 v_2$$

- Magnitude of vector  $u$ :  $\|u\| = \sqrt{u_1^2 + u_2^2}$

- Projection of  $v$  on  $u$ :  $p$

- Then:  $u^T v = p \|u\|$

# Why large margin

---



$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$



Dot product:

$$\begin{aligned} \theta^T x^{(i)} &= p^{(i)} \|\theta\| \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

# Why large margin

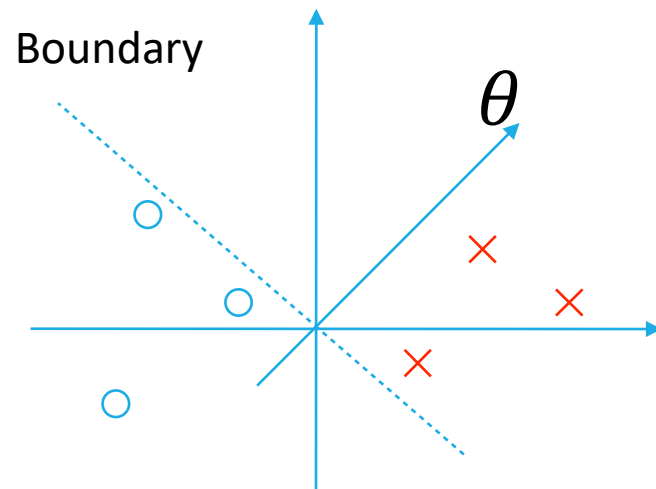
---

- $$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

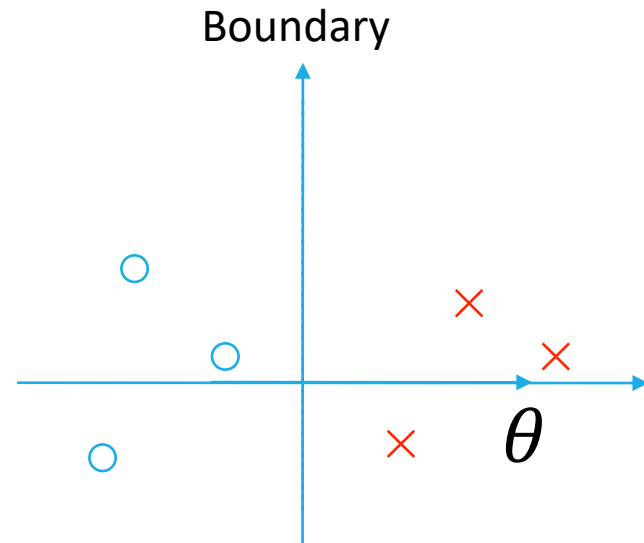
s.t.  $p^{(i)} \|\theta\| \geq 1$  if  $y^{(i)} = 1$   
 $p^{(i)} \|\theta\| \leq -1$  if  $y^{(i)} = 0$
- Given  $p^{(i)}$  to be length of projection of  $x^{(i)}$  on  $\theta$
- $p^{(i)}$  is distance between sample and boundary
- Do  $p^{(i)} \|\theta\| \geq 1$ 
  - $p^{(i)}$  is smaller, length of  $\theta$  is larger
  - $p^{(i)}$  is large, length of  $\theta$  can be small

# Why large margin

---



Narrow margin



Large margin

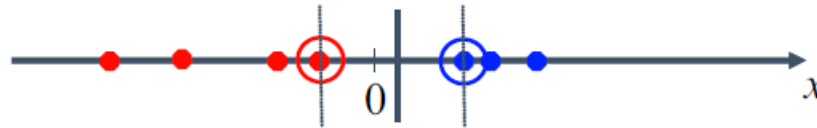
# Feature mapping with Kernel

---

Part 2

# Non-linear classification

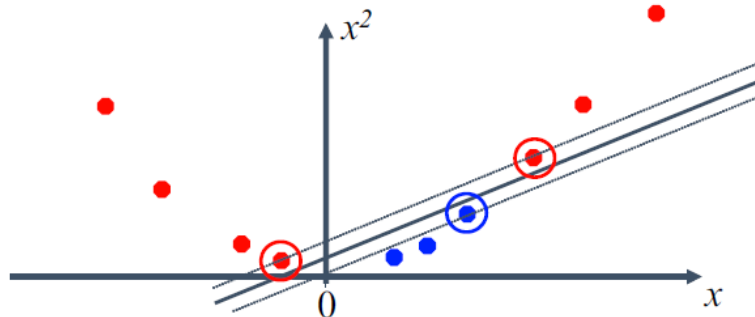
- Linear classifier works well on linearly separated data



- Some data are not linearly separated



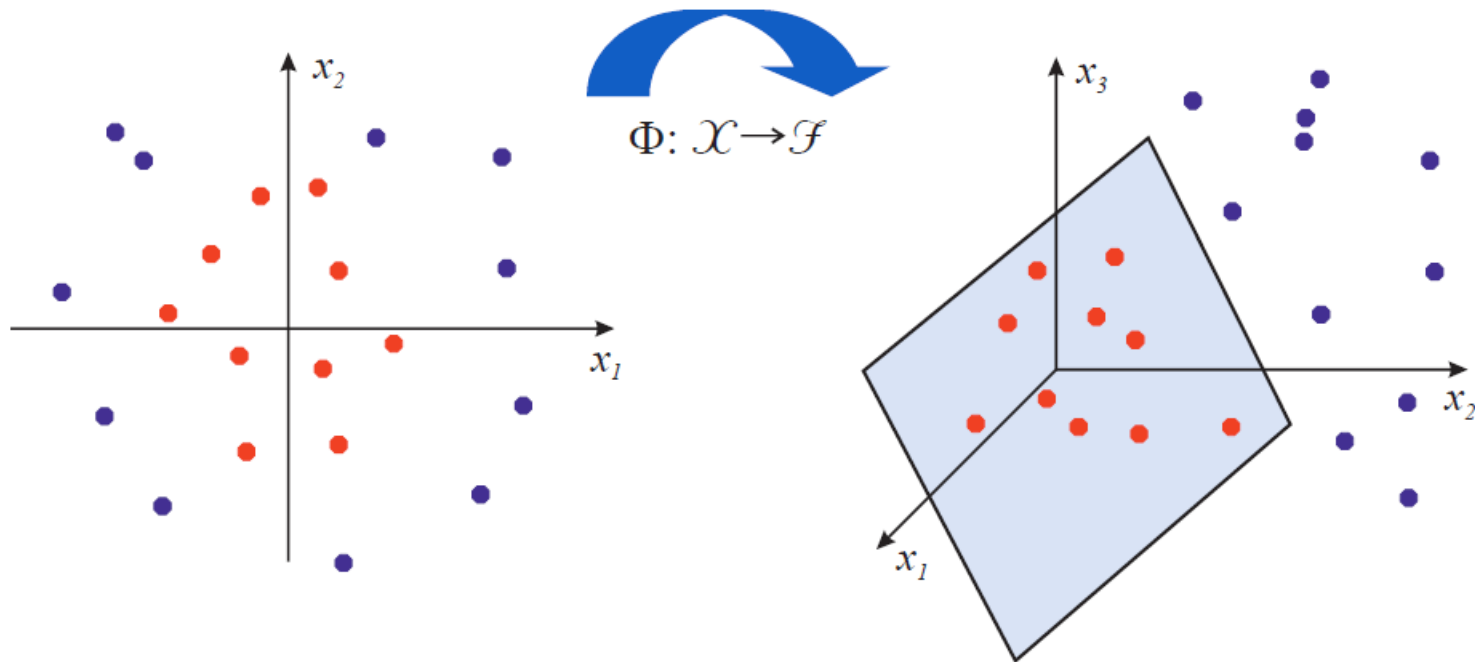
- Solution: map data into higher dimension spaces





# Non-linear classification

- Solution: map data into higher dimension spaces



How is data mapping is done with SVM?

# Learning algorithm

---

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^n \theta_j^2 \\ \text{s.t.} \quad & \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ & \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0 \end{aligned}$$

Given  $y = -1$  for negative class

$$\begin{aligned} \Rightarrow \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^n \theta_j^2 \\ \text{s.t.} \quad & y^{(i)} (\theta^T x^{(i)}) \geq 1, i = 1, 2, \dots, m \end{aligned}$$

# Learning algorithm

---

□ Given  $g_i(\theta) = -y^{(i)}(\theta^T x^{(i)}) + 1$

$$\Rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t. } g_i(\theta) = -y^{(i)}(\theta^T x^{(i)}) + 1 \leq 0$$

Switch to Lagrange problem:

$$\mathcal{L}(\theta, \alpha) = \frac{1}{2} \sum_{j=1}^n \theta_j^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(\theta^T x^{(i)}) - 1]$$

$$\text{Goal: } \min_{\theta, \alpha} \mathcal{L}(\theta, \alpha)$$

# Learning algorithm

---

$$\mathcal{L}(\theta, \alpha) = \frac{1}{2} \sum_{j=1}^n \theta_j^2 - \sum_{i=1}^m \alpha_i [y^{(i)} (\theta^T x^{(i)}) - 1]$$

Goal:  $\min_{\theta, \alpha} \mathcal{L}(\theta, \alpha)$

Conditions Karush-Kuhn-Tucker (KKT)

Suppose  $\theta^*, \alpha^*$  to be solution

- (1)  $\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta^*, \alpha^*) = 0, j = 1, \dots, n$
- (2)  $\alpha_i^* g_i(\theta^*) = 0, i = 1, \dots, m$
- (3)  $g_i(\theta^*) \leq 0, i = 1, \dots, m$
- (4)  $\alpha_i^* \geq 0, i = 1, \dots, m$

# Learning algorithm

---

$$\mathcal{L}(\theta, \alpha) = \frac{1}{2} \sum_{j=1}^n \theta_j^2 - \sum_{i=1}^m \alpha_i [y^{(i)} (\theta^T x^{(i)}) - 1]$$

Goal:  $\min_{\theta, \alpha} \mathcal{L}(\theta, \alpha)$

Solution for theta:

$$(1) \quad \nabla_{\theta} \mathcal{L}(\theta, \alpha) = \theta - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0$$

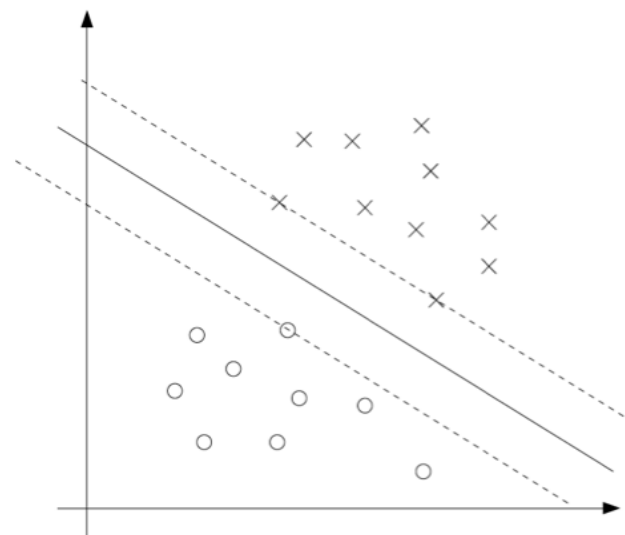
$$\Rightarrow \theta = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

# Support vector

Conditions Karush-Kuhn-Tucker (KKT)

Suppose  $\theta^*, \alpha^*$  to be solution

- (1)  $\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta^*, \alpha^*) = 0, j = 0, \dots, n$
- (2)  $\alpha_i^* g_i(\theta^*) = 0, i = 1, \dots, m$
- (3)  $g_i(\theta^*) \leq 0, i = 1, \dots, m$
- (4)  $\alpha_i^* \geq 0, i = 1, \dots, m$



From (2), (3) and (4):

if  $g_i(\theta^*) < 0$ , then  $\alpha_i^* = 0, i = 1, \dots, m$

Samples outside margins:  $\alpha_i^* = 0$

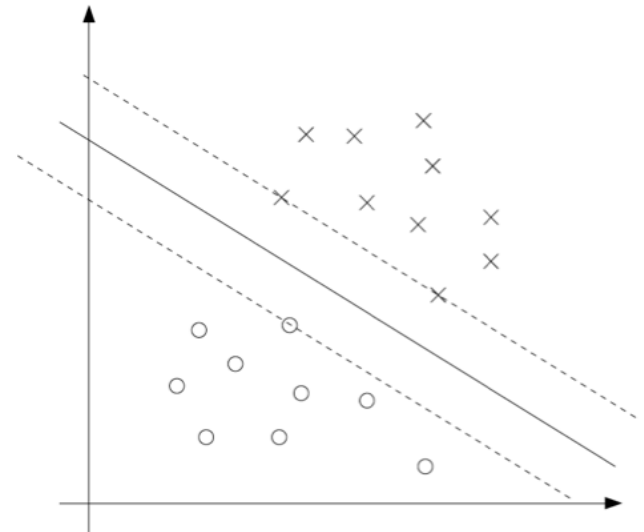
# Support vector

Samples outside margins:  $\alpha_i^* = 0$

$$\Rightarrow \theta = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

Predict:

$$\begin{aligned} \theta^T x &= \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x \\ &= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle \end{aligned}$$



Only calculated with data points on margins  $\Rightarrow$  support vector

# Feature mapping

---

- Given a function

$$\phi(x) = [x \quad x^2 \quad x^3]^T$$

Predict:

$$\begin{aligned}\theta^T x &= \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x \\ &= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle\end{aligned}$$

Replace  $\langle x^{(i)}, x \rangle$  by  $\langle \phi(x^{(i)}), \phi(x) \rangle$



# Kernel

---

□ Given

$$K(x^{(i)}, x) = \langle \phi(x^{(i)}), \phi(x) \rangle$$

Predict:

$$\begin{aligned}\theta^T x &= \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x \\ &= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle\end{aligned}$$

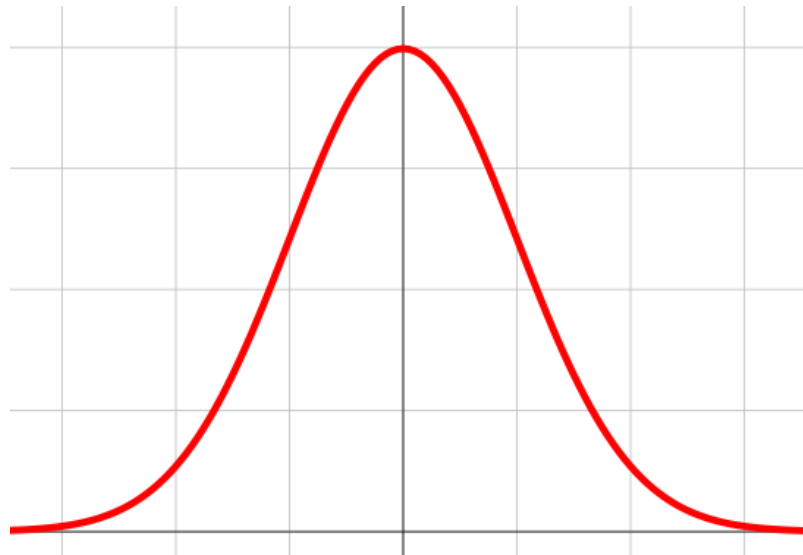
Replace  $\langle x^{(i)}, x \rangle$  by  $K(x^{(i)}, x)$

# Kernel

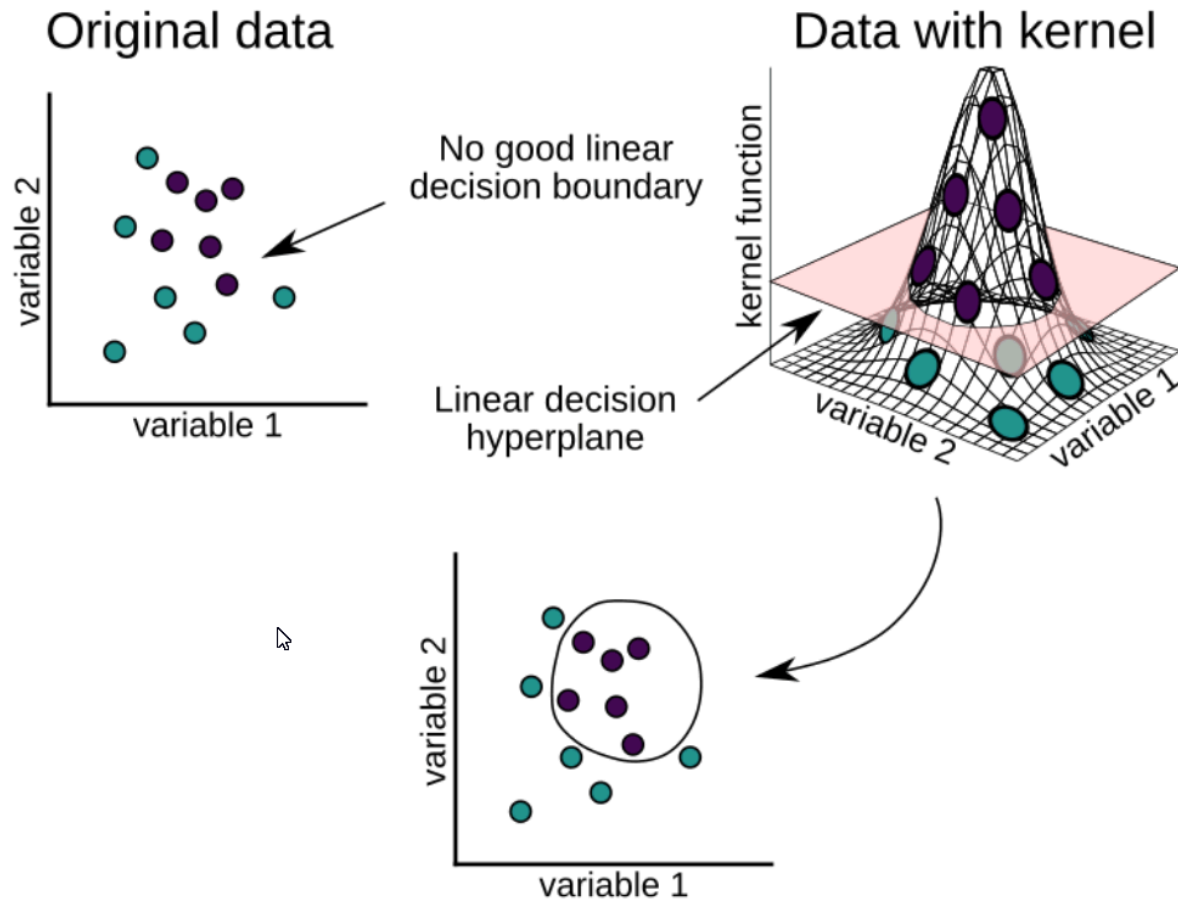
---

## □ Gaussian kernel

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



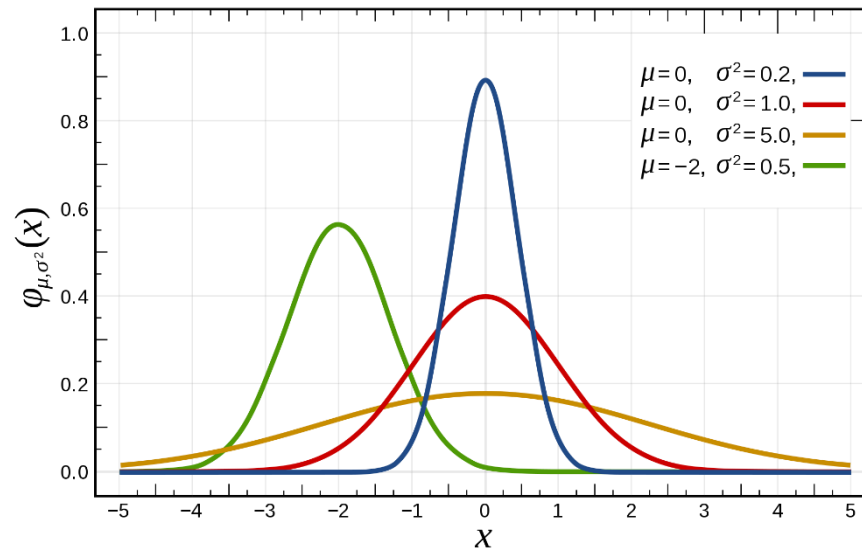
# Kernel



# Kernel

## □ Gaussian (RBF) kernel

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



- Larger  $\sigma$ , smaller magnitude

# Kernel

---

- Linear kernel

$$K(x, z) = x^T z$$

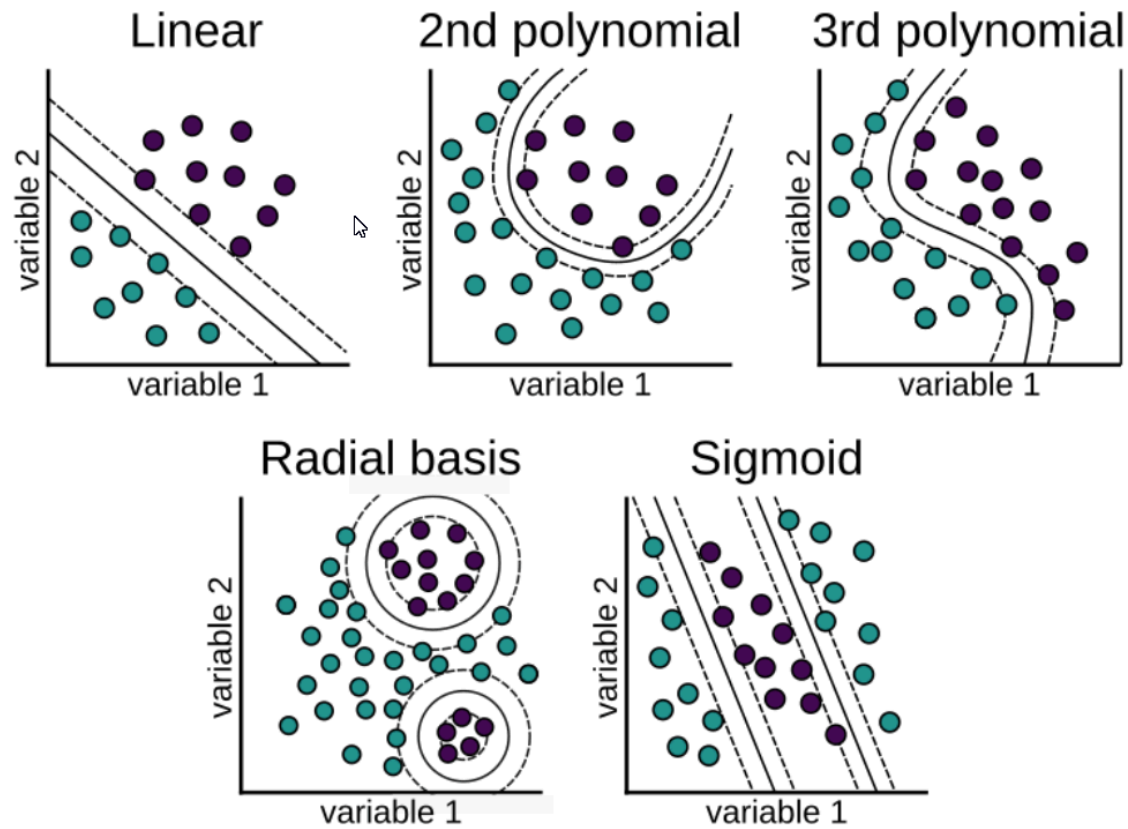
- Polynomial kernel

$$K(x, z) = (x^T z + c)^d$$

- Sigmoid kernel

$$K(x, z) = \tanh(ax^T z + b)$$

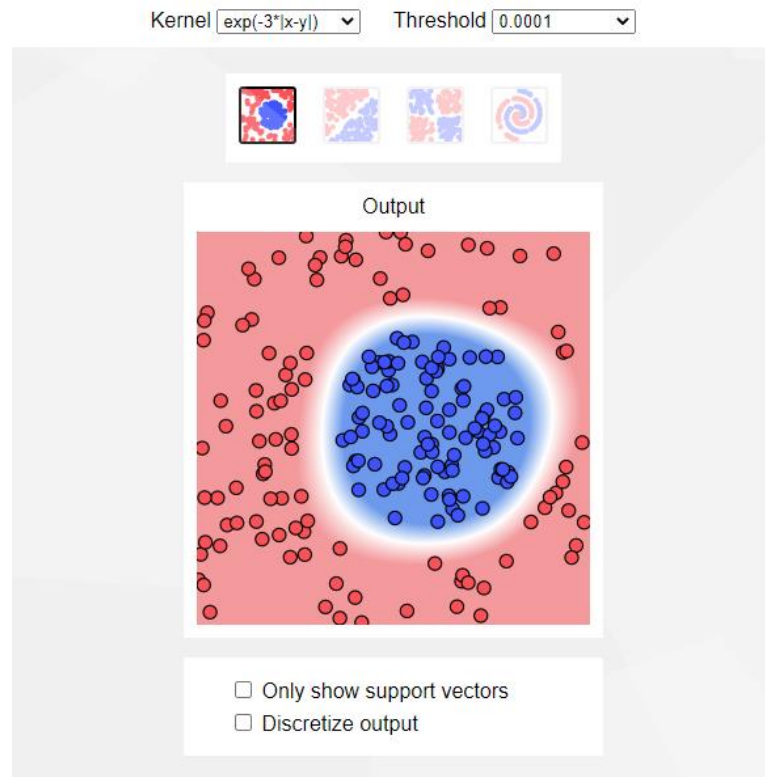
# Kernel



# SVM Playground

---

❑ <http://macheads101.com/demos/svm-playground>



# SVM Playground

---

❑ <http://macheads101.com/demos/svm-playground>

