

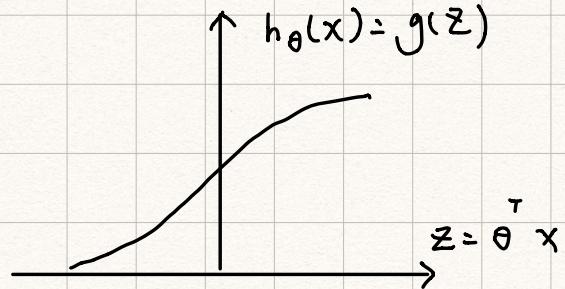
Model 表現狀況取決於：  
特徵量選擇 and 選擇正則化參數

## Support Vector Machines (SVM) 支持向量機

### Optimization Object

### Alternative view of logistic regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



If  $y=1$ , we want  $h_{\theta}(x) \approx 1$ ,  $\theta^T x \gg 0$

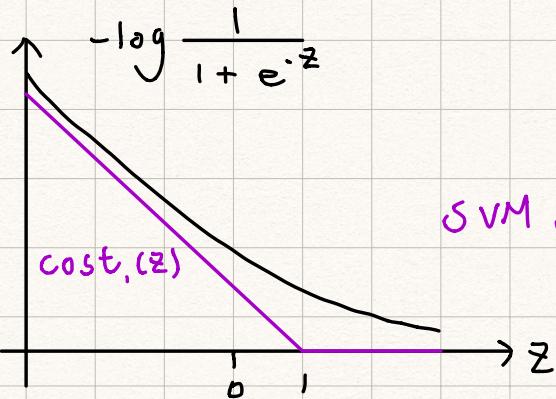
If  $y=0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$

### Alternative view of logistic regression

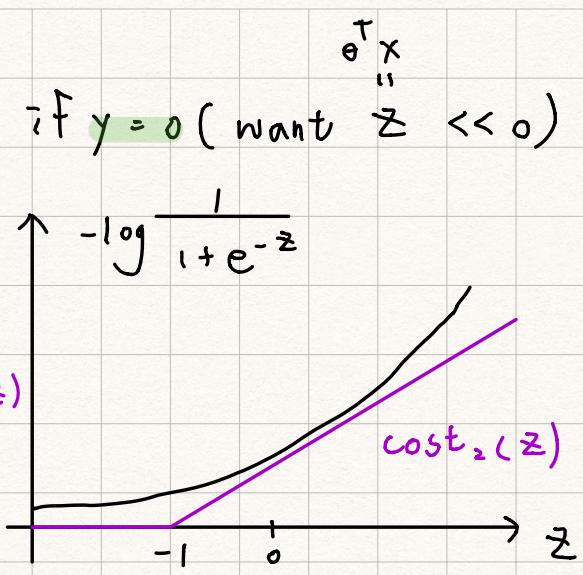
Cost of example:  $-(y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)))$

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

if  $y=1$  (want  $z \gg 0$ )      if  $y=0$  (want  $z \ll 0$ )



SVM  $J(z)$



$\text{cost}_2(z)$

## Support vector machine

Logistic regression:

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

## Support vector Machine:

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

if  $C = \frac{1}{\lambda}$ , 則以上兩條算式算出來  
的值會相等  
↓ 正則化因子

SVM Hypothesis:

SVM 在預測結果時所用的 hypothesis

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

SVM 用於找是否高工具次方程式有諸多幫助

### Large Margin Intuition

### SVM Decision Boundary

if  $y=1$ , we want  $\theta^T x \geq 1$  (not just  $\geq 0$ )  $\theta^T x \geq 0$  /

if  $y=0$ , we want  $\theta^T x \leq -1$  (not just  $< 0$ )  $\theta^T x \leq 0$  -1

↳ 如此一來，誤才會等於 0

Whenever  $y^{(i)} = 1$ :  
 $\theta^T x^{(i)} \geq 1$

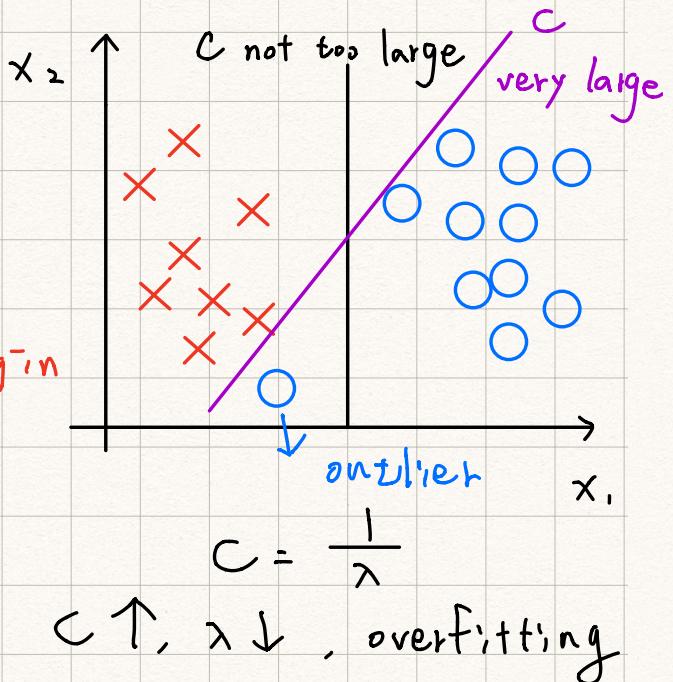
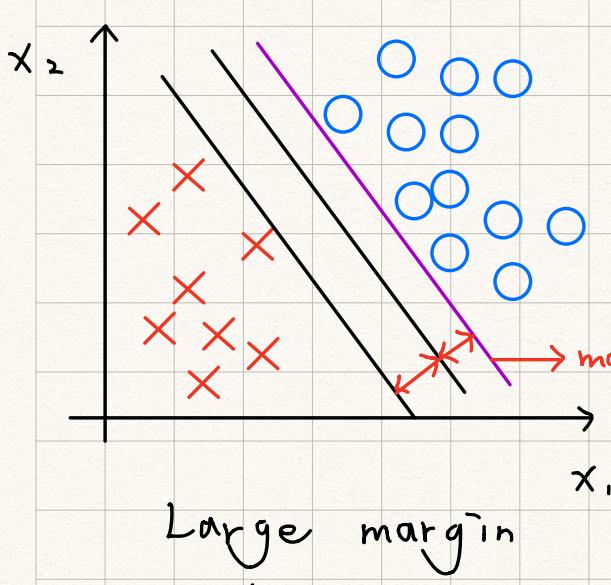
Whenever  $y^{(i)} = 0$ :  
 $\theta^T x^{(i)} \leq -1$

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

$$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

產生大間距分類器

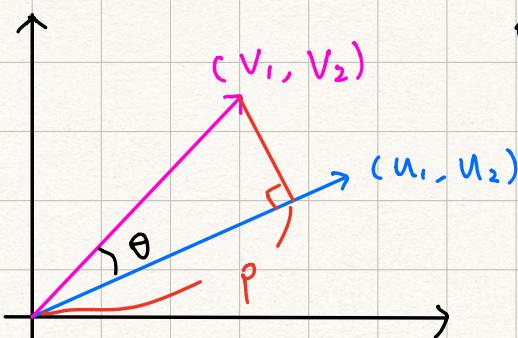
## SVM Decision Boundary : Linearly separable case



## Mathematics Behind Large Margin Classification

本部分選學課程，利用內積理解 SVM

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad |u| = \sqrt{u_1^2 + u_2^2}$$



$p = \text{length of projection of } v \text{ onto } u$

$$\begin{aligned} u^T v &= |u| |v| \cos \theta = p \cdot |u| \\ &= u_1 v_1 + u_2 v_2 \end{aligned}$$

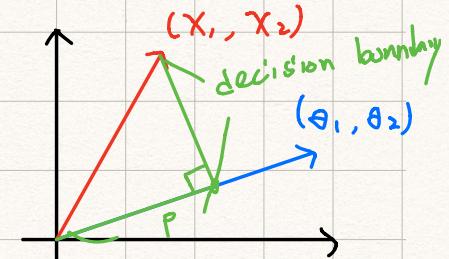
# SVM Decision Boundary

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

視為零

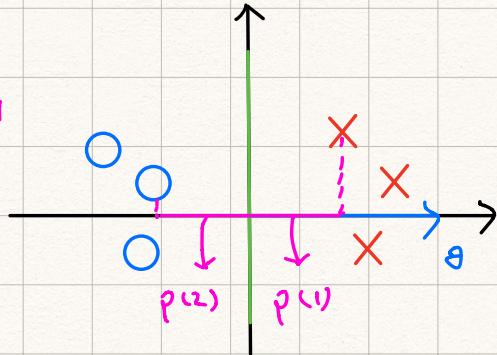
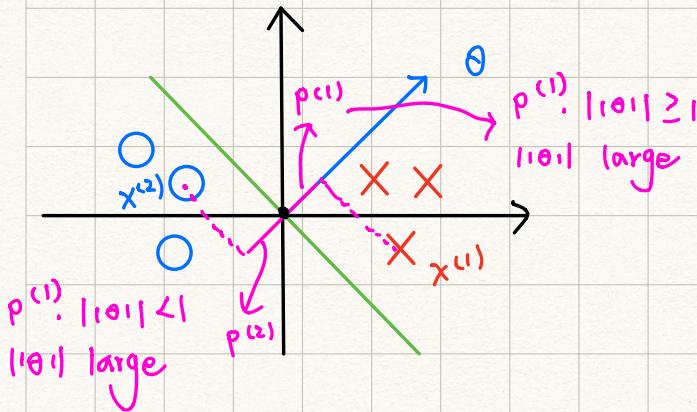
$$\boxed{\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2} = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

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



s.t.  $p^{(i)} \cdot \|\theta\| \geq 1$  if  $y^{(i)} = 1$   
 $p^{(i)} \cdot \|\theta\| \leq -1$  if  $y^{(i)} = 0$

where  $p^{(i)}$  is the projection of  $x^{(i)}$  onto the  
Simplification:  $\theta_0 = 0 \rightarrow$  因此 decision boundary 不通過原點



目標是



最小化  $\min \frac{1}{2} \sum_{j=1}^n \theta_j^2$ , 但比左圖的大得多。

$p^{(1)}$  和  $p^{(2)}$



由此圖分析得之,  $\|\theta\|$  is 因此大, 故此條 decision boundary 非最佳化

$\|\theta\|$  can be smaller

## Kernel | 核密度函數

透過此函數, 可以訓練出複雜的  
非線性 Decision Boundary

通過標記點和相似性函數定義新的  
特徵變量

## Non-linear Decision Boundary

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, \dots$$

Is there a different/better choice of  
the features  $f_1, f_2, f_3 \dots$ ?

## Kernel



Given  $x$ , compute new feature depending on proximity to landmarks.  
 $\mathcal{L}^{(1)}, \mathcal{L}^{(2)}, \mathcal{L}^{(3)}$

Given  $x$ :

$$f_1 = \text{similarity}(x, \mathcal{L}^{(1)}) = \exp\left(-\frac{\|x - \mathcal{L}^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, \mathcal{L}^{(2)}) = \exp\left(-\frac{\|x - \mathcal{L}^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, \mathcal{L}^{(3)}) = \exp(\dots)$$

kernel (Gaussian kernels) =  $k(x, \mathcal{L}^{(1)})$

## Kernels and Similarity

$$f_1 = \text{similarity}(x, \mathcal{L}^{(1)}) = \exp\left(-\frac{\|x - \mathcal{L}^{(1)}\|^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{\sum_{j=1}^n (x_j - \mathcal{L}_j^{(1)})^2}{2\sigma^2}\right)$$

$$x \approx \mathcal{L}^{(1)}, f_1 \approx \exp\left(-\frac{0}{2\sigma^2}\right) \approx 1$$

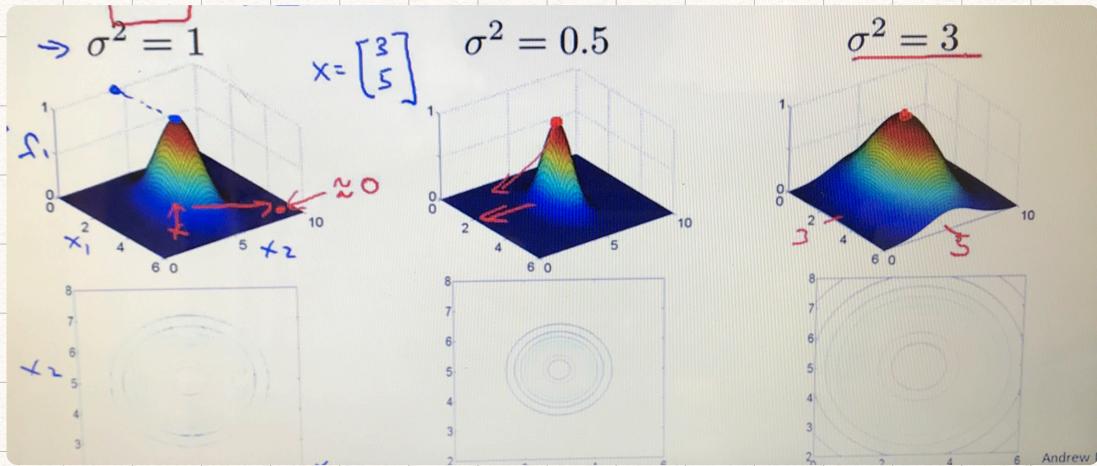
$$x \text{ far from } \mathcal{L}^{(1)}, f_1 \approx \exp\left(-\frac{\text{large Num}}{2\sigma^2}\right) \approx 0$$

Example

$$f_i = \exp\left(-\frac{\|x - \boldsymbol{\lambda}^{(i)}\|^2}{2\sigma^2}\right)$$

$$\boldsymbol{\lambda}^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

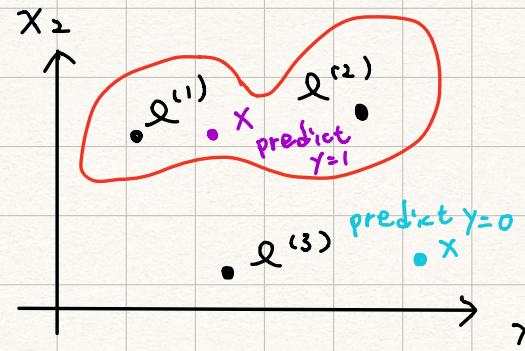
$x_1, x_2, f_i$  的函數隨  $\sigma^2$  而變



依上圖所示， $\sigma^2$  越大，核函數坡度越緩  
 $\sigma^2$  越小，核函數坡度越陡

Predict '1' when  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$



$$f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$$

$$\theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0 \geq 0$$

$$f_1, f_2, f_3 \approx 0 \rightarrow \theta_0 + 0 + 0 \approx -0.5 < 0$$

## Choosing the landmarks

Where to get landmark,  $\mathcal{L}^{(1)}, \mathcal{L}^{(2)} \dots$ ?

直接將樣本點轉為 landmark

## SVM with kernels

Given  $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ , ...,  $(x^{(m)}, y^{(m)})$   
choose  $\mathcal{L}^{(1)} = x^{(1)}$ ,  $\mathcal{L}^{(2)} = x^{(2)}$ , ...,  $\mathcal{L}^{(m)} = x^{(m)}$

- example  $x$ :

$$f_1 = \text{similarity}(x, \mathcal{L}^{(1)})$$

$$f_2 = \text{similarity}(x, \mathcal{L}^{(2)})$$

...

$$f = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{pmatrix}$$

$$f_0 = 1$$

- For training example  $(x^{(i)}, y^{(i)})$ :

$$f_{1(i)} = \text{sim}(x^{(i)}, \mathcal{L}^{(1)})$$

$$x^{(i)} \rightarrow f_{2(i)} = \text{sim}(x^{(i)}, \mathcal{L}^{(2)})$$

⋮

$$f_{1(i)} = \text{sim}(x^{(i)}, \mathcal{L}^{(1)}) = 1$$

$$f_{m(i)} = \text{sim}(x^{(i)}, \mathcal{L}^{(m)})$$

$$f = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{pmatrix}$$

$$f_{0(i)} = 1$$

Hypothesis : Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$   
 Predict " $y=1$ " if  $\theta^T f \geq 0$

Training :

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

↓                                  ↓  
 取正則化类似，不對下標為 0 的數  
 進行計算 ignore  $\theta_0$ .

$$\sum_j \theta_j^2 = \theta^T \theta$$

SVM 是這樣實現的，多  
 乘一個  $M$  使 SVM 更有效  
 進行運算

$$= \theta^T M \theta$$

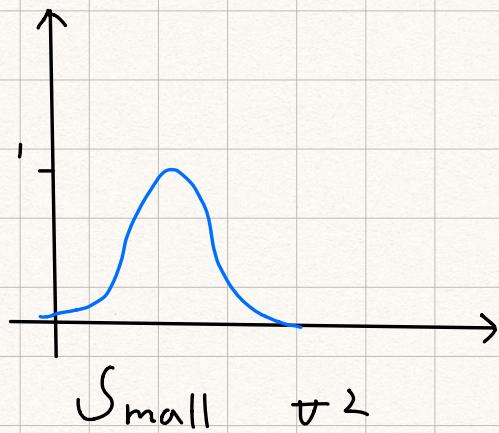
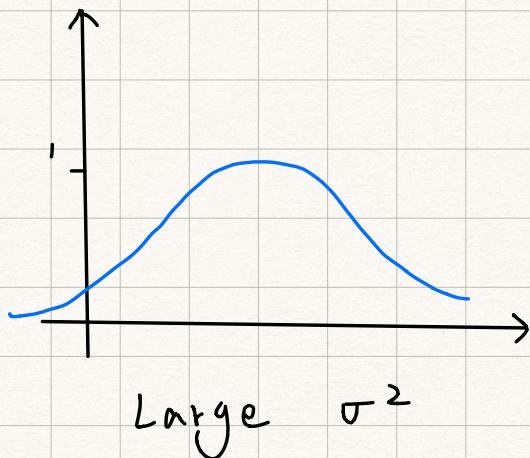
另，核函數 (kernel) 是特地為 SVM 設計的。  
 用在其他算法如 logistic regression 會使之運  
 算緩慢。

## SVM Parameters

- $C (= \frac{1}{\lambda})$ . large  $C$ : Lower bias, high variance  
 Small  $C$ : Higher bias, low variance  
 $C$  和  $\lambda$  成反比

- $\sigma^2$ . Large  $\sigma^2$ : Features  $f_i$  more smoothly.  
Higher bias, lower variance

Small  $\sigma^2$ : Features  $f_i$  less smoothly.  
Lower bias, higher variance.



### Using an SVM

Use SVM software package (e.g. liblinear, libsvm, ...)

- 不建議自己寫算法，調用現有函式庫即可

Need to specify:

Choice of parameter C

Choice of kernel (Similarity function)



①

No kernel ("linear kernel")

Predict "y=1" if  $\theta^T x \geq 0$

使用時機：  
特徵數量多但資料量少的時候  
即以線性子集之，而不造成 overfitting

②

Gaussian Kernel:

Need to choose  $\sigma^2$ .

Kernel (similarity) functions

function  $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

return

$x \rightarrow$   
 $f_1$   
 $f_2$   
 $\vdots$   
 $f_m$

Note: Do perform **feature scaling** before using the Gaussian kernel.

因為各特徵的資料大小有極大的差別

## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, \mathbf{x})$  make valid kernels.

Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM Packages' optimizations run correctly, and do not diverge.

Many off-the-shelf kernel available:

- Polynomial Kernel
- More esoteric: String Kernel, chi-square Kernel histogram intersection kernel.

## Multi-class classification

Many SVM packages already have build-in multi-class classification functionality.

Pick class  $i$  with largest  $(\boldsymbol{\theta}^{(i)})^T \mathbf{x}$

## Logistic regression vs. SVMs

$n$  是特徵數量的大小， $m$  是訓練集的數量

- $n$  is large relative to  $m$ :

Use logistic regression or SVM without kernel

$$n \geq m, n = 10000, m = 10 \dots 1000$$

- $n$  is small,  $m$  is intermediate:

Use SVM with Gaussian Kernel

$$n = 1 - 1000, m = 10 - 10000$$

- if  $n$  is small,  $m$  is large:

Create / Add more features, then use logistic regression or SVM without a kernel

$$n = 1 - 1000, m = 50000$$

NN likely to work well for most of these settings, but slower to train.

對於 SVM, 是凸優化問題, 不必擔心找到局部最小值, SVM 之解通常為全局最小值