

Machine Learning

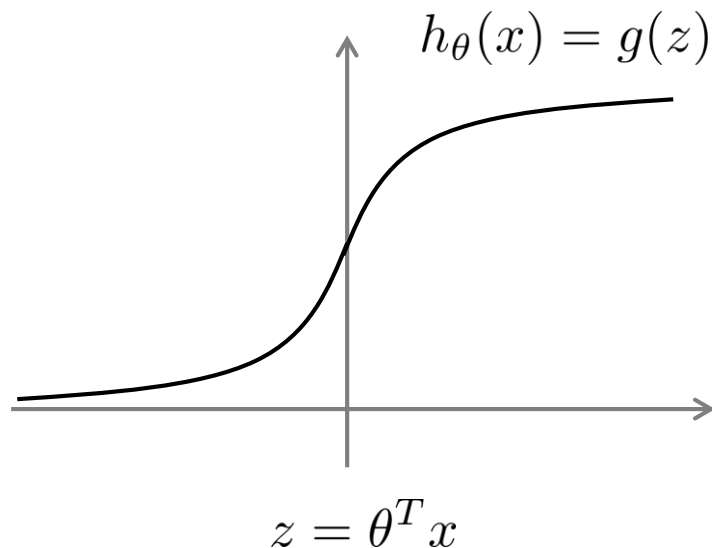
# Support Vector Machines

---

## Optimization objective

# Alternative view of logistic regression

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



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

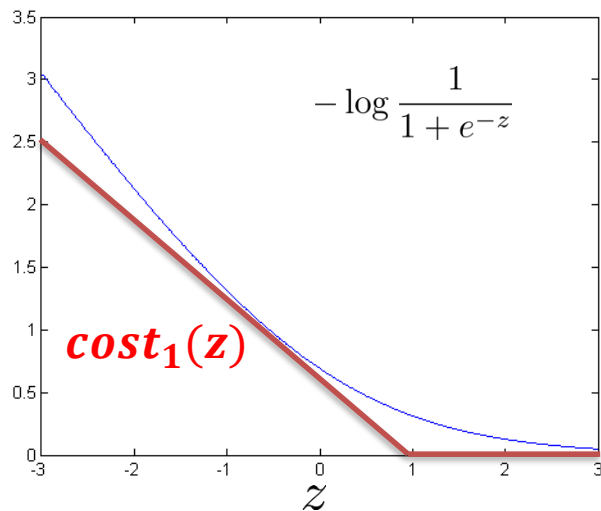
If  $y = 0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$

## Alternative view of logistic regression

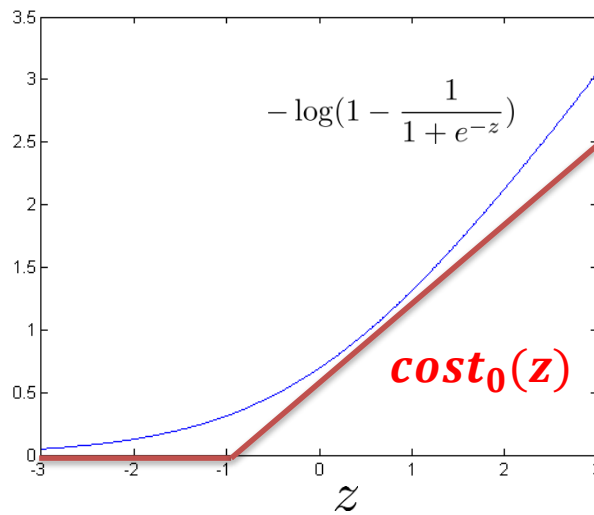
Cost of example:  $-(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$

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

If  $y = 1$  (want  $\theta^T x \gg 0$ ):



If  $y = 0$  (want  $\theta^T x \ll 0$ ):



# Support vector machine

Logistic regression:

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

Support vector machine:

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

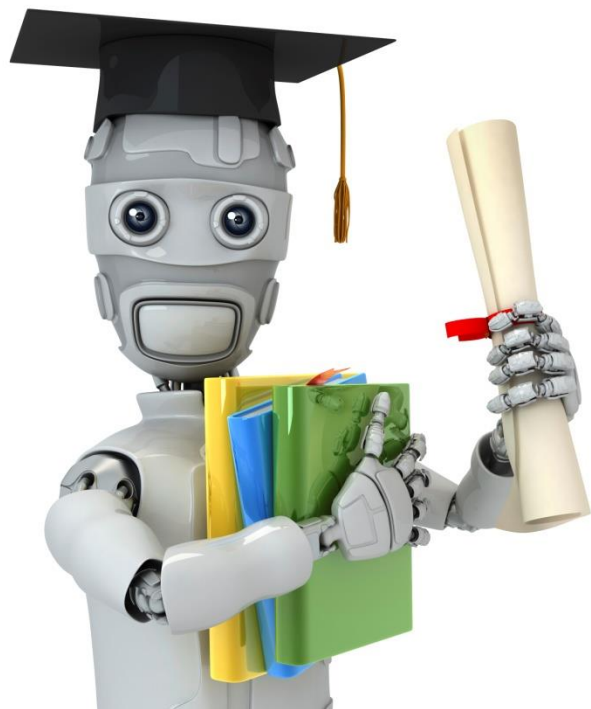
$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

## SVM hypothesis

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

Hypothesis:

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



Machine Learning

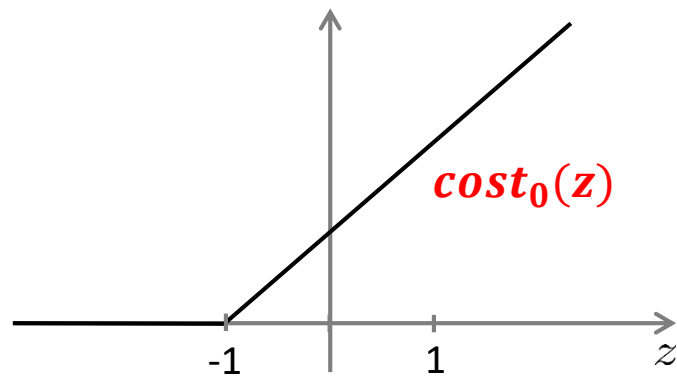
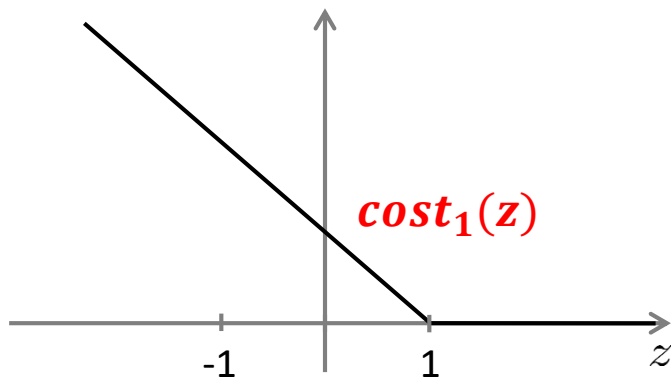
# Support Vector Machines

---

## Large Margin Intuition

# Support Vector Machine

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



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

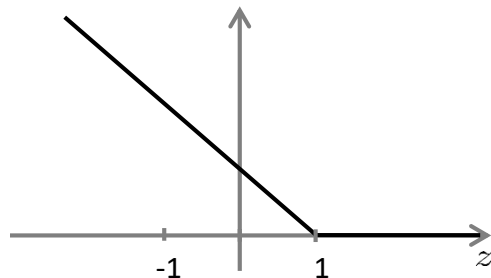
If  $y = 0$ , we want  $\theta^T x \leq -1$  (not just  $< 0$ )

# SVM Decision Boundary

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

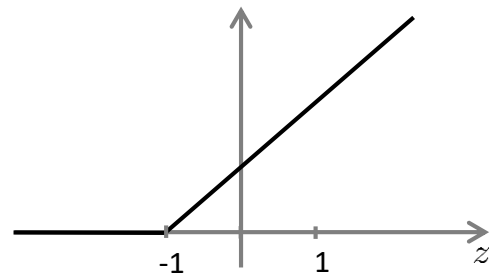
Whenever  $y^{(i)} = 1$ :

$$\theta^T x^{(i)} \geq 1$$



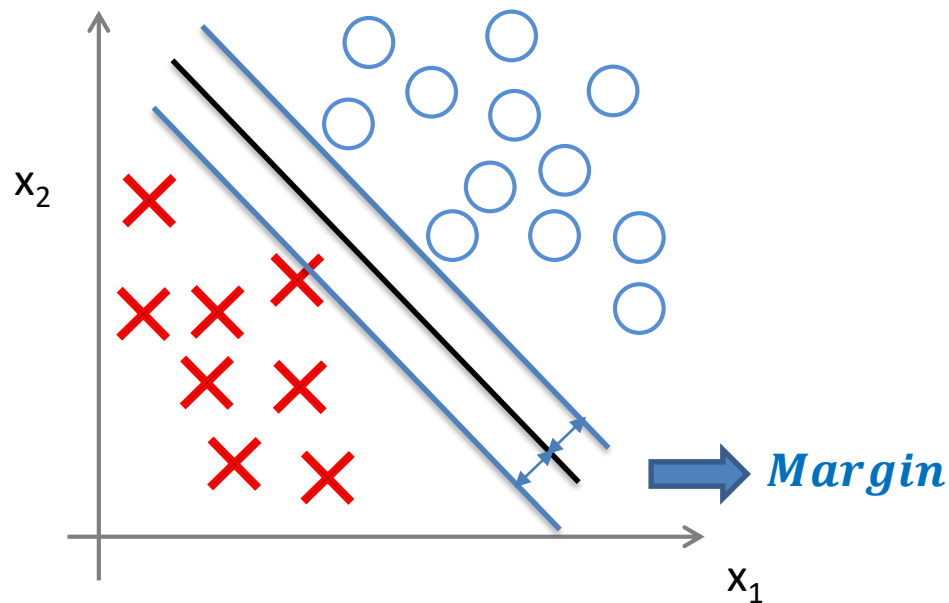
Whenever  $y^{(i)} = 0$ :

$$\theta^T x^{(i)} \leq -1$$



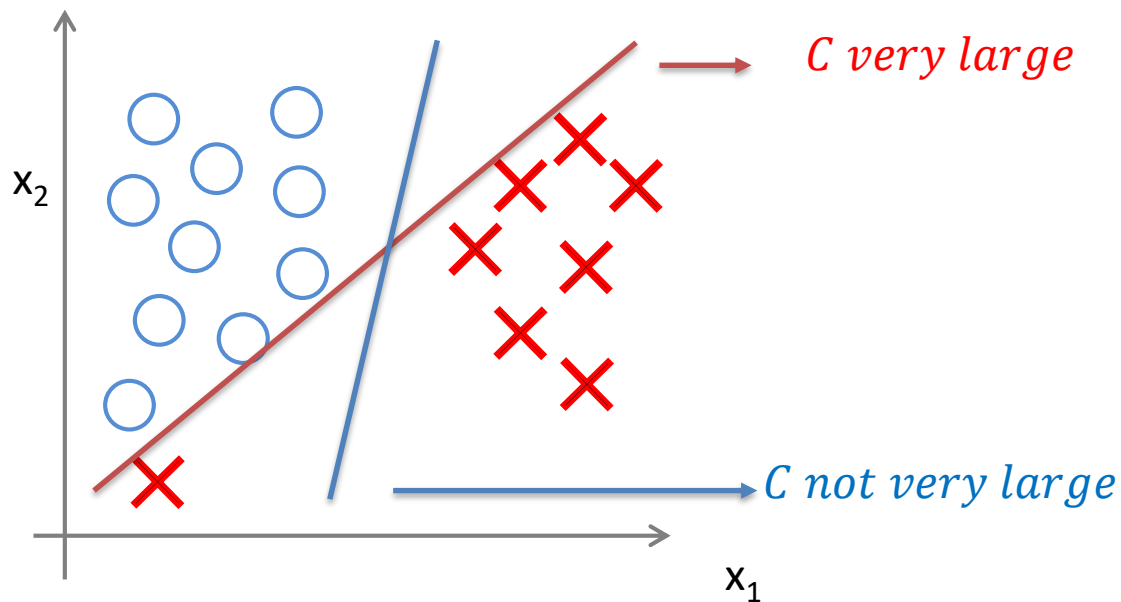


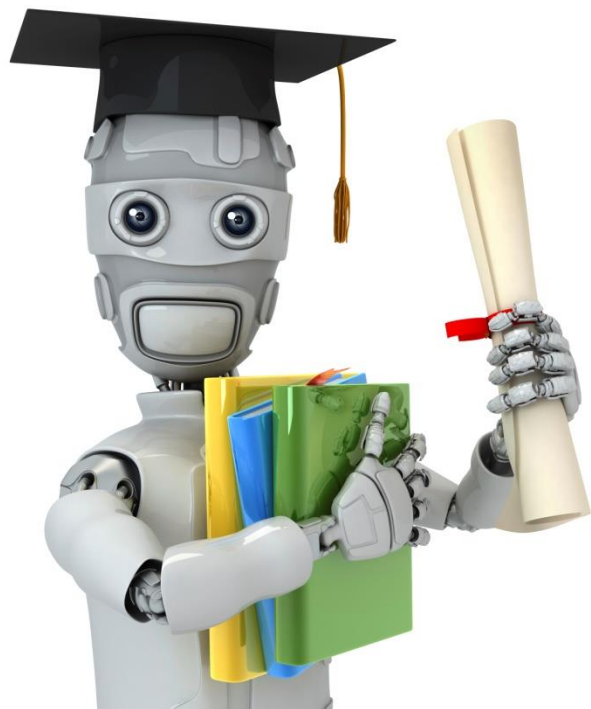
# SVM Decision Boundary: Linearly separable case



Large margin classifier

# Large margin classifier in presence of outliers





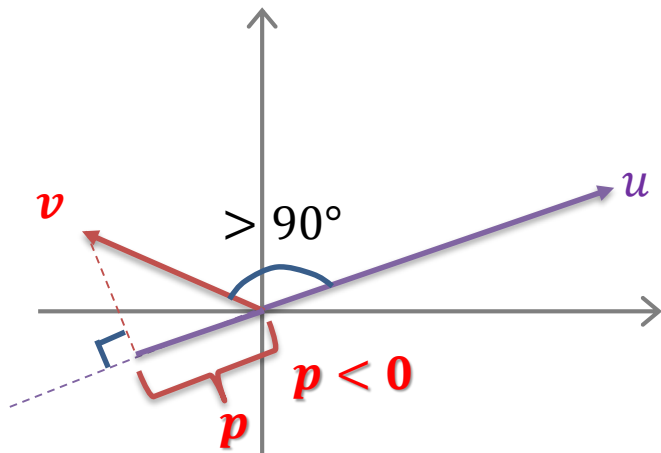
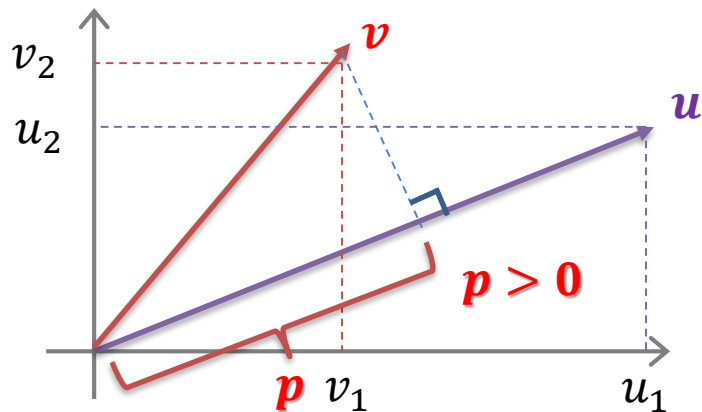
Machine Learning

# Support Vector Machines

---

The mathematics  
behind large margin  
classification (optional)

# Vector Inner Product



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

$u, v$ : two vectors

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

$\|u\|$ : length of vector  $u$

$p$ : length of projection of vector  $v$  to  $u$

$$\begin{aligned} u^T v &= [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \\ &= u_1 \times v_1 + u_2 \times v_2 \\ &= p \times \|u\| \\ &= v^T u \end{aligned}$$

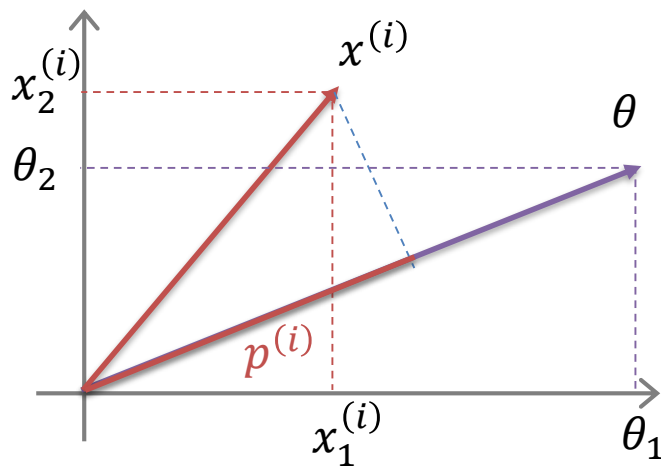
## SVM Decision Boundary

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

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

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

*simplification:  $n = 2, \theta_0 = 0$*



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

# SVM Decision Boundary

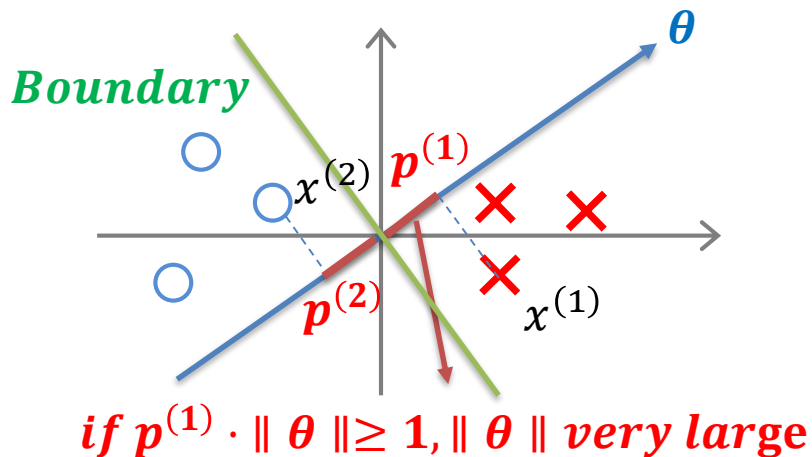
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

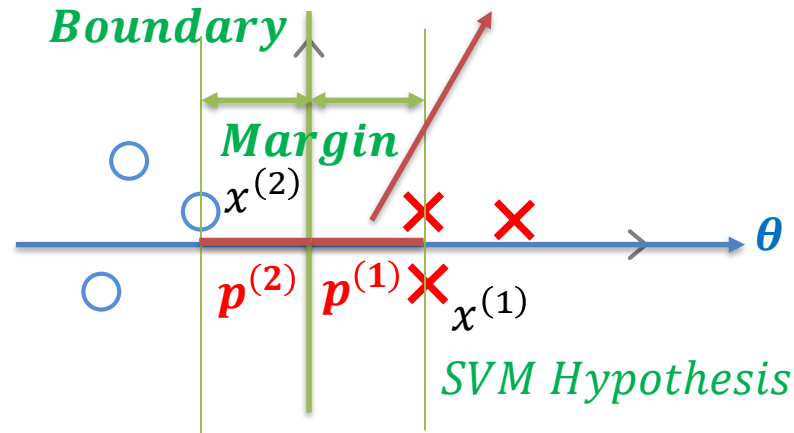
$$p^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = -1$$

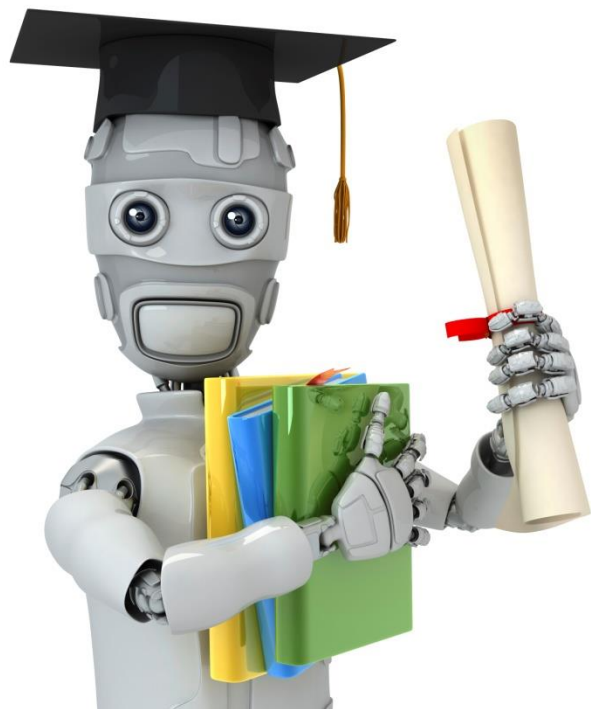
where  $p^{(i)}$  is the projection of  $x^{(i)}$  onto the vector  $\theta$ .

Simplification:  $\theta_0 = 0$



*if  $p^{(1)} \cdot \|\theta\| \geq 1$ ,  $\|\theta\|$  can be smaller*





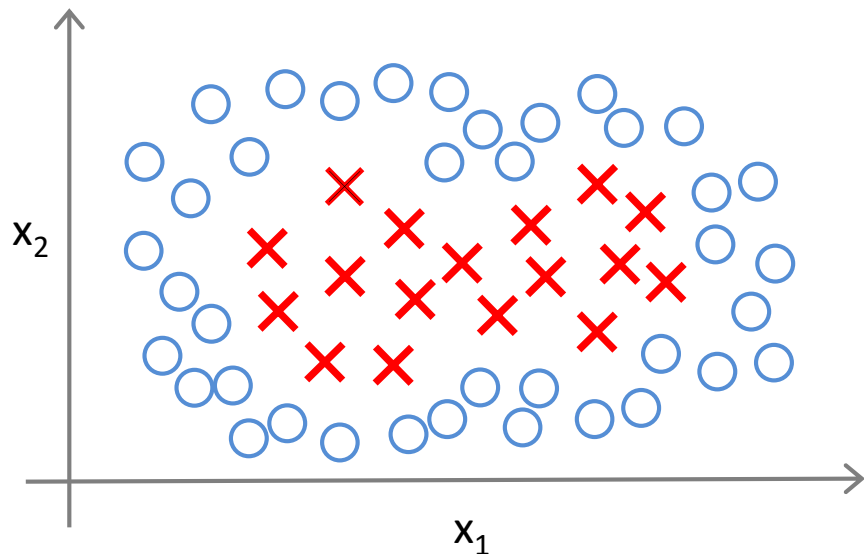
Machine Learning

# Support Vector Machines

---

# Kernels I

## Non-linear Decision Boundary



Predict  $y = 1$  if

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 \\ + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

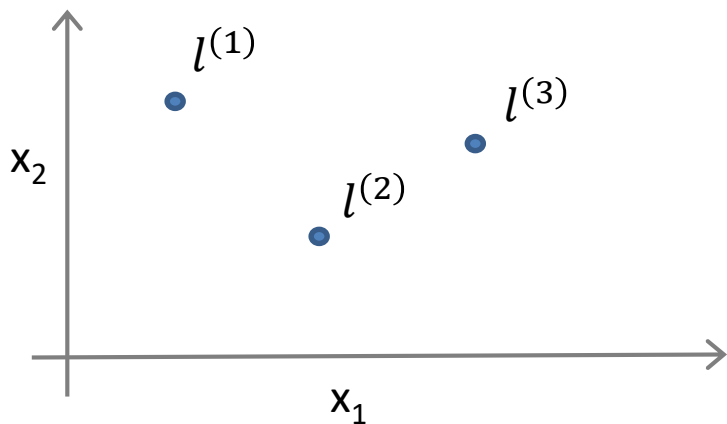
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \dots$$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2$$

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  $l^{(1)}, l^{(2)}, l^{(3)}$

$$\|x - l^{(i)}\|^2 = \sum_{j=1}^n (x_j - l_j^{(i)})^2$$

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

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

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|^2}{2\sigma^2}\right)$$

Kernel

## Kernels and Similarity

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

If  $x \approx l^{(1)}$  :

$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) \approx 1$$

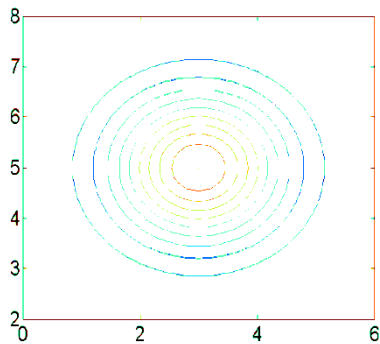
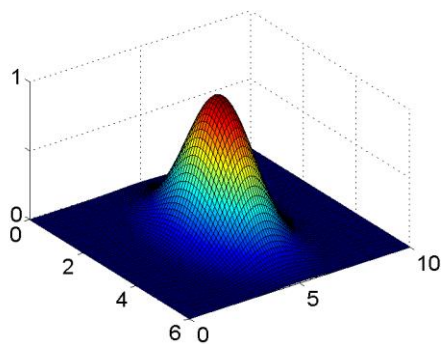
If  $x$  is far from  $l^{(1)}$  :

$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) \approx 0$$

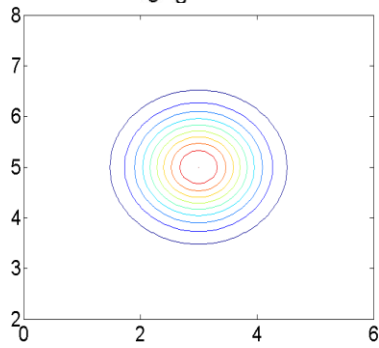
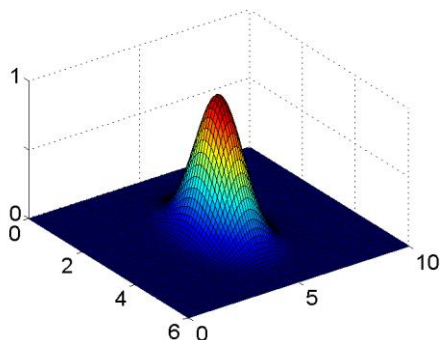
**Example:**

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

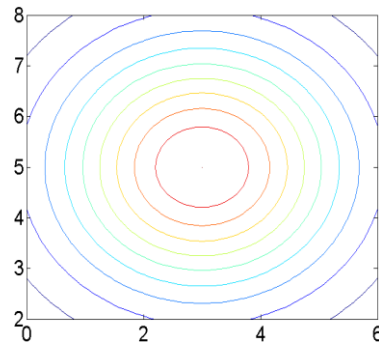
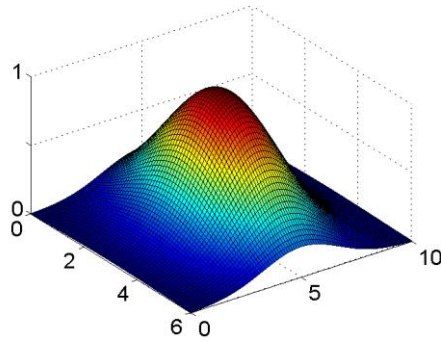
$$\sigma^2 = 1$$



$$\sigma^2 = 0.5$$



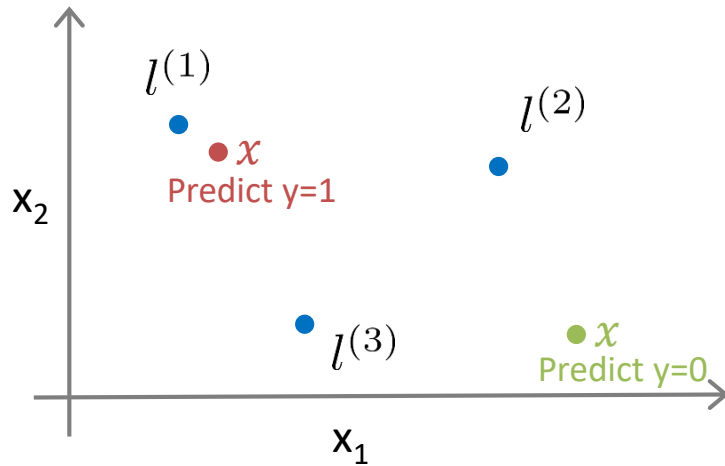
$$\sigma^2 = 3$$



Predict “1” when

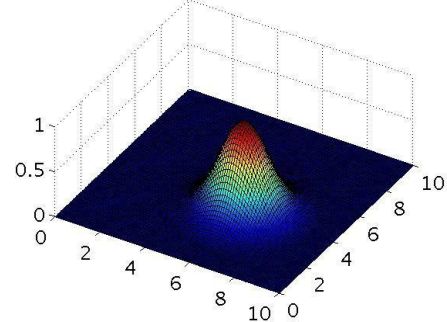
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

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

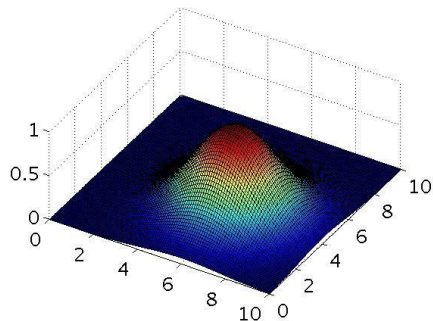


高斯核的公式是由  $\text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$  给出的。

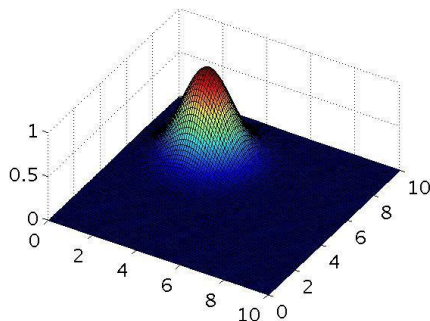
下图显示了当  $\sigma^2 = 1$  时,  $f_1 = \text{similarity}(x, l^{(1)})$  的曲线图。



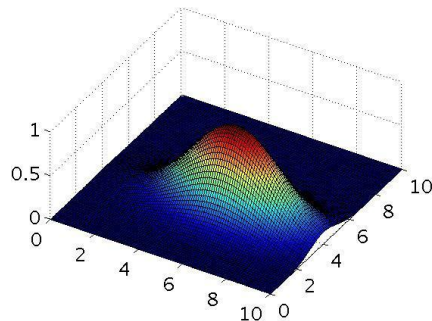
当  $\sigma^2 = 0.25$  时, 下列哪个是  $f_1$  的曲线图?



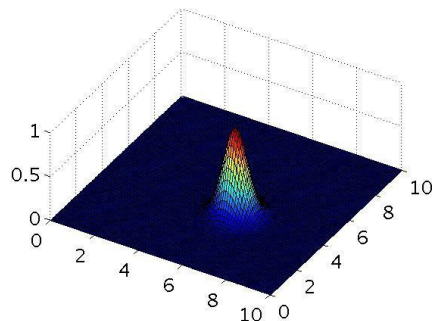
A.



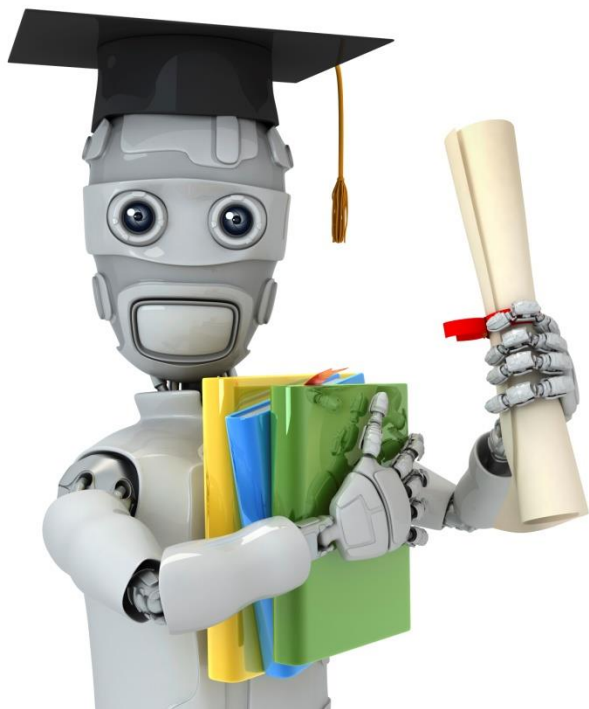
B.



C.



D.



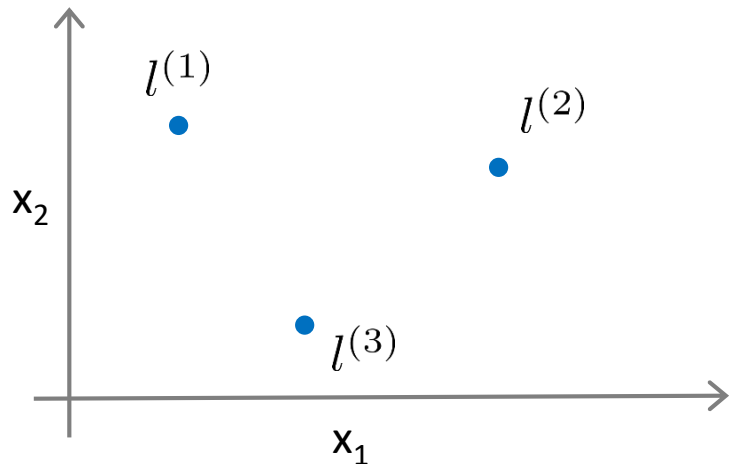
Machine Learning

# Support Vector Machines

---

## Kernels II

## Choosing the landmarks

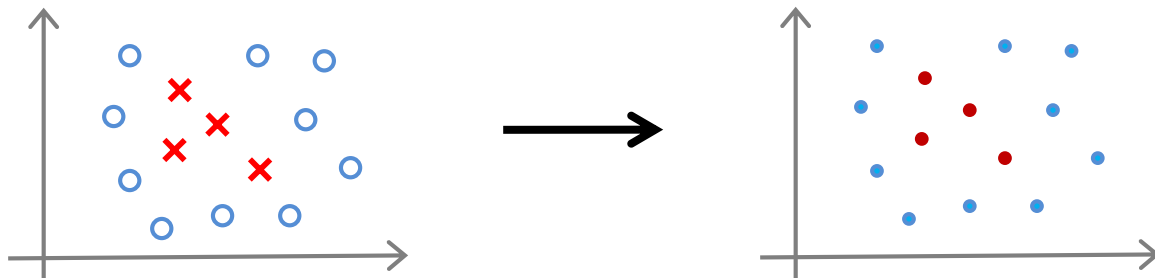


Given  $x$ :

$$f_i = \text{similarity}(x, l^{(i)})$$
$$= \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

Predict  $y = 1$  if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

Where to get  $l^{(1)}, l^{(2)}, l^{(3)}, \dots$ ?



## SVM with Kernels

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

Given example  $x$ :

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

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

...

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

$$f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)})$$

$$x^{(i)} \Rightarrow f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)})$$

$$\vdots \leftarrow f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp\left(-\frac{0^2}{2\sigma^2}\right) = 1$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)})$$



## SVM with Kernels

Hypothesis: Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$

Predict “y=1” if  $\theta^T f \geq 0$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} cost_1(\theta^T f^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

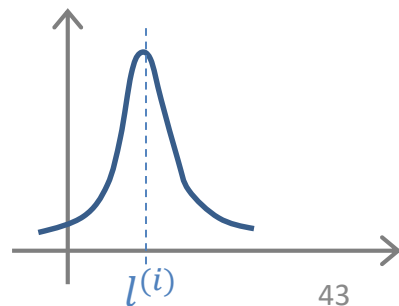
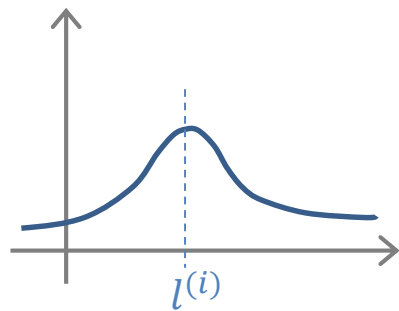
## SVM parameters:

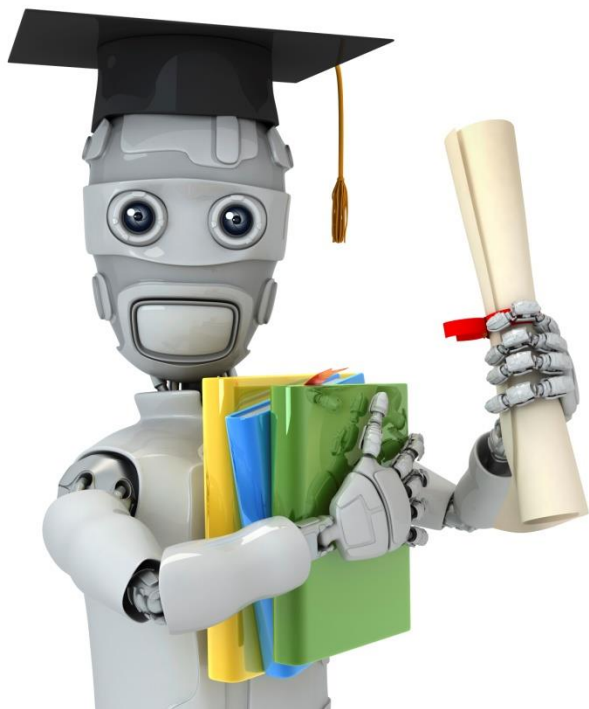
$C \left( = \frac{1}{\lambda} \right)$ . Large  $C$ : Lower bias, high variance.  
Small  $C$ : Higher bias, low variance.

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

$$\exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

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





Machine Learning

# Support Vector Machines

---

## Using an SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .

Need to specify:

- Choice of parameter C.

- Choice of kernel (similarity function):

E.g. No kernel (“linear kernel”)

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

Gaussian kernel:

$$f_i = \exp \left( -\frac{\|x - l^{(i)}\|^2}{2\sigma^2} \right), \text{ where } l^{(i)} = x^{(i)}.$$

Need to choose  $\sigma^2$ .

**Kernel (similarity) functions:**

```
function f = kernel(x1,x2)
```

$$f = \exp \left( -\frac{\| \mathbf{x1} - \mathbf{x2} \|^2}{2\sigma^2} \right)$$

```
return
```

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

## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels. (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

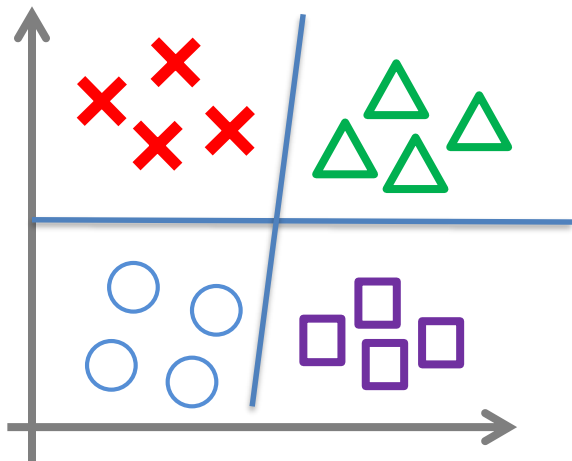
Many off-the-shelf kernels available:

- Polynomial kernel:

$$K(x, l) = (X^T l + m)^d, d = 1, 2, \dots, N$$

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

## Multi-class classification



$$y \in \{1, 2, 3, \dots, K\}$$

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

Otherwise, use one-vs.-all method. (Train  $K$  SVMs, one to distinguish  $y = i$  from the rest, for  $i = 1, 2, \dots, K$ ), get  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$   
Pick class  $i$  with largest  $(\theta^{(i)})^T x$

## Logistic regression vs. SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

If  $n$  is large (relative to  $m$ ):

Use logistic regression, or SVM without a kernel (“linear kernel”)

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

Use SVM with Gaussian kernel

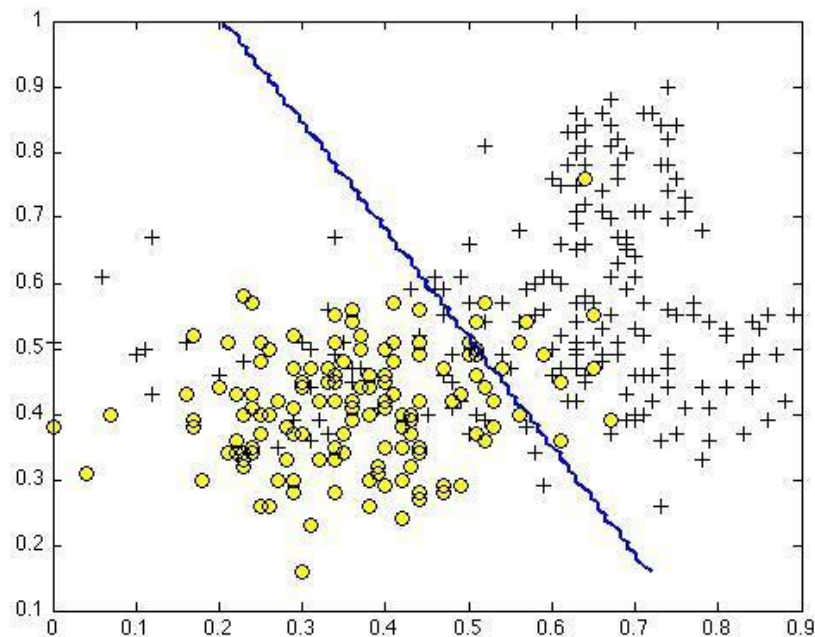
If  $n$  is small,  $m$  is large:

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

Neural network likely to work well for most of these settings, but may be slower to train.



假设您使用训练了一个高斯内核的支持向量机，它在训练集上学习了以下决策边界：



你觉得支持向量机欠拟合了，你应该试着增加或减少 $C$ 吗？或者增加或减少 $\sigma^2$ ？

A. 降低 $C$ ，增加 $\sigma^2$  B. 降低 $C$ ，降低 $\sigma^2$  C. 增加 $C$ ，增加 $\sigma^2$  D. 增加 $C$ ，降低 $\sigma^2$