# Support Vector Machines

## Optimization objective

Machine Learning

# Alternative view of logistic regression

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

$$h_\theta(x) = g(z)$$

$$z = \theta^T x$$

$$z = \theta^T x$$

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

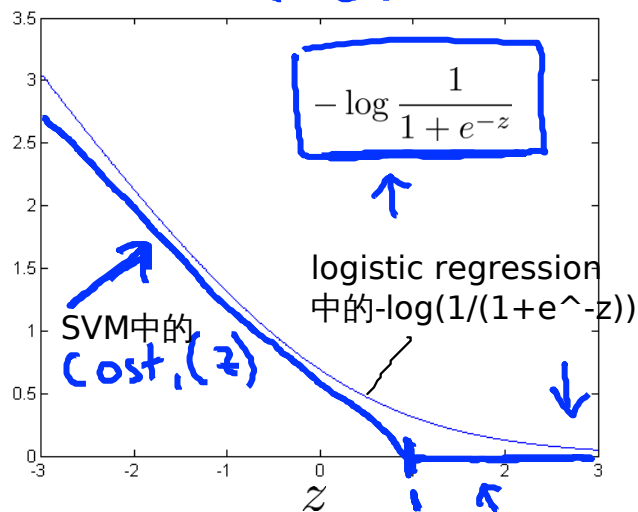If $y = 0$, we want $h_\theta(x) \approx 0$,    $\theta^T x \ll 0$

# Alternative view of logistic regression

$(x, y)$

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

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

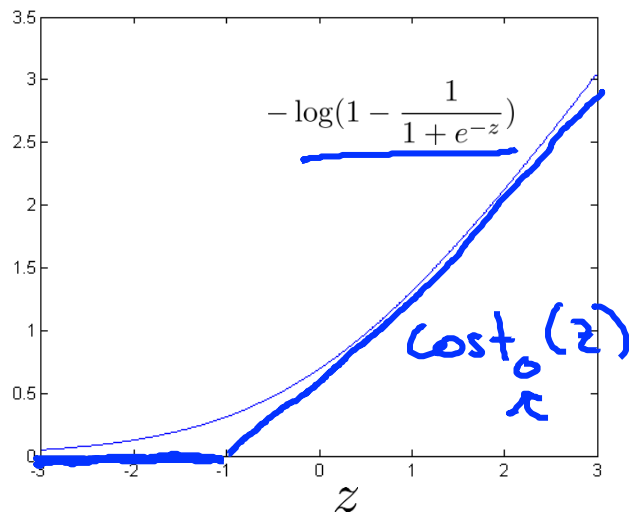If $y = 1$ (want $\theta^T x \gg 0$):  $z = \theta^T x$

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



$-\log \frac{1}{1 + e^{-z}}$

SVM中的 $\text{cost}_1(z)$

logistic regression 中的-log(1/(1+e^-z))

$z$

z為θ^T x



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

$\text{cost}_0(z)$

$z$

## 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:

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

Andrew Ng

在这个最小化问题中 无论前面是否有 1/m 这一项 最终我所得到的 最优值θ都是一样的

## SVM hypothesis

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

Hypothesis:     你也可以 把这里的参数C 考虑成 1/λ(regularizaton中的) 同 1/λ 所扮演的 角色相同

有别于逻辑回归 输出的概率 在这里 我们的代价函数 当最小化代价函数 获得参数θ时 支持向量机所做的是
它来直接预测 y的值等于1 还是等于0

# Support Vector Machines
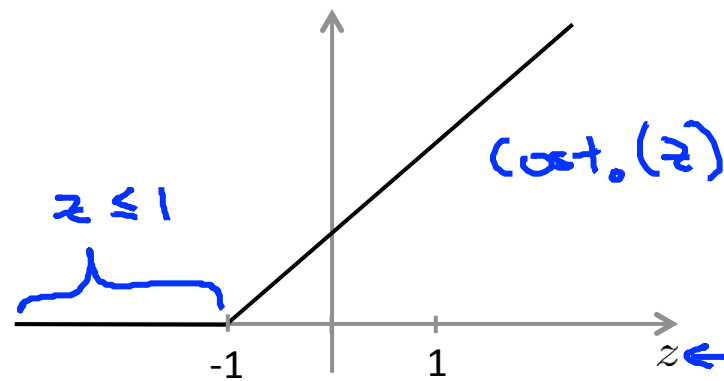
## Large Margin Intuition

Machine Learning

# Support Vector Machine

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$



$cost_1(z)$  $z \geq 1$

$cost_0(z)$  $z \leq 1$

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

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

$C = 100,000$

# SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

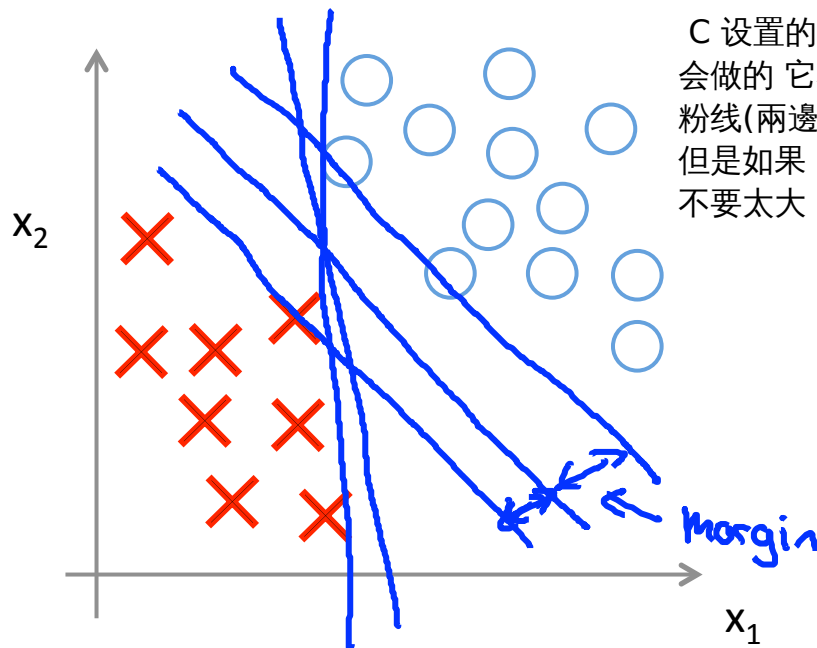$= 0$

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

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

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

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

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

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

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

# SVM Decision Boundary: <u>Linearly separable case</u>

我知道你也许 想知道 求解上一页
幻灯片中的优化问题 为什么会产生
这个结果 它是如何产生这个大间距
分类器的呢
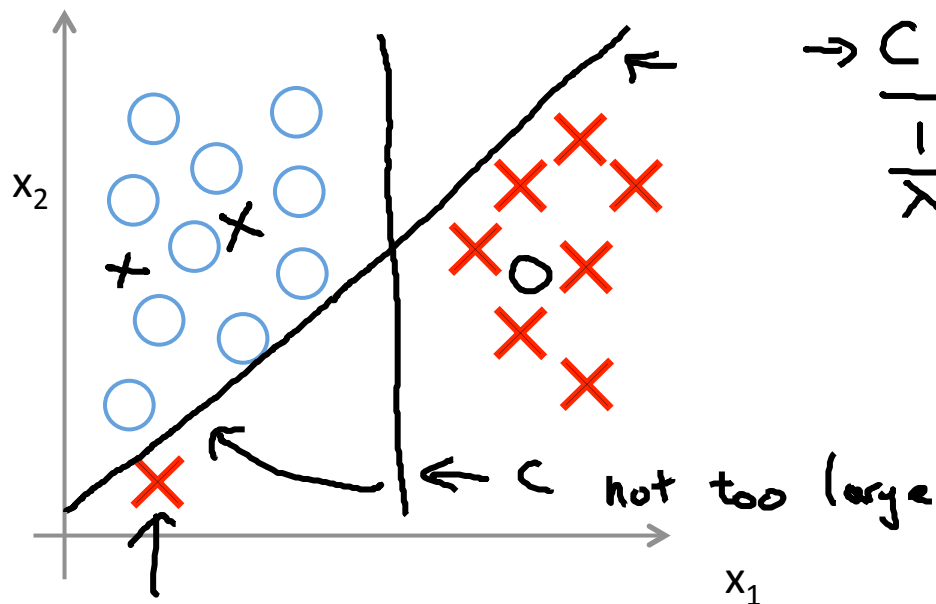我知道我还没有解释这一点
在下一节视频中 我将会从
直观上 略述

C 设置的非常大 这事实上正是 支持向量机将
会做的 它将决策界 从黑线(中间的) 变到了
粉线(两邊的)
但是如果 C 设置的小一点 如果你将 C 设置的
不要太大 则你最终会得到 这条黑线



Large margin classifier

# Large margin classifier in presence of outliers

当 C 不是 非常非常大的时候
它可以忽略掉一些异常点的影响
得到更好的决策界



$\rightarrow$ C very large

$\frac{1}{K}$

C not too large

# Support Vector Machines

## The mathematics behind large margin classification (optional)

Machine Learning

# Vector Inner Product



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

$$u^T v = ? \qquad [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$\|u\| = \text{length of vector } u$

$$= \sqrt{u_1^2 + u_2^2} \quad \in \mathbb{R}$$

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

Signed

$$u^T v = \underline{p} \cdot \underline{\|u\|} \leftarrow \qquad = v^T u$$

$$= u_1 v_1 + u_2 v_2 \leftarrow \qquad p \in \mathbb{R}$$

$$u^T v = p \cdot \|u\|$$

$$p < 0$$

# SVM Decision Boundary

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

$$\omega = \left(\sqrt{\omega'}\right)^2$$

$$= \|\theta\|$$

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

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

$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$ $\quad \theta_0 = 0$

Simplication: $\theta_0 = 0$ $\quad n = 2$

$\theta^T x^{(i)} = ?$

$\uparrow \quad \uparrow$
$u^T v$

$\theta^T x^{(i)} = \boxed{p^{(i)} \cdot \|\theta\|} \leftarrow$

$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \leftarrow$



Andrew Ng

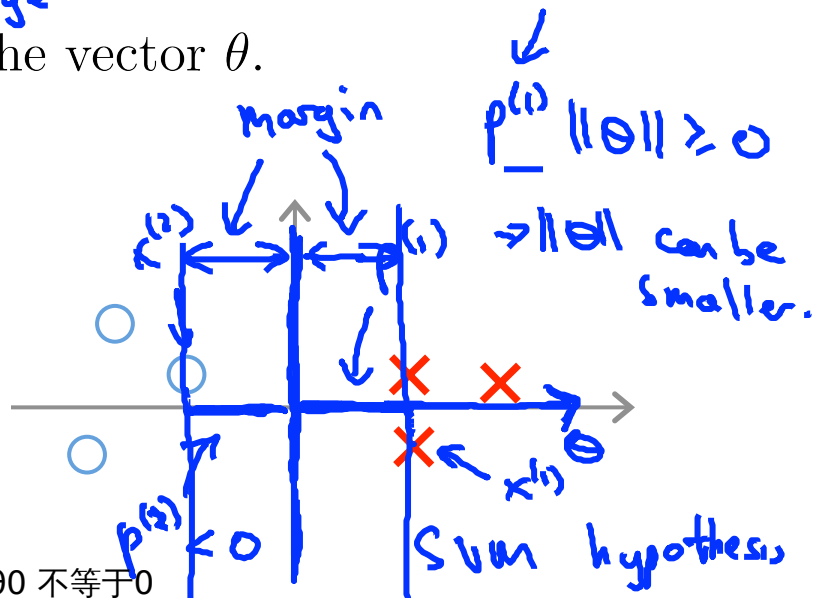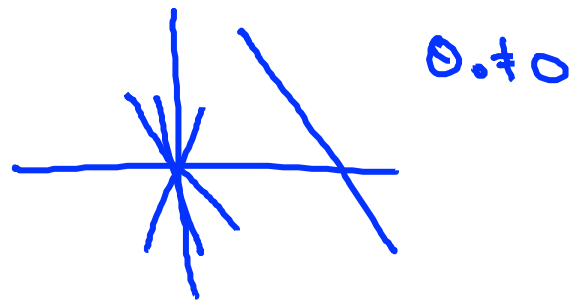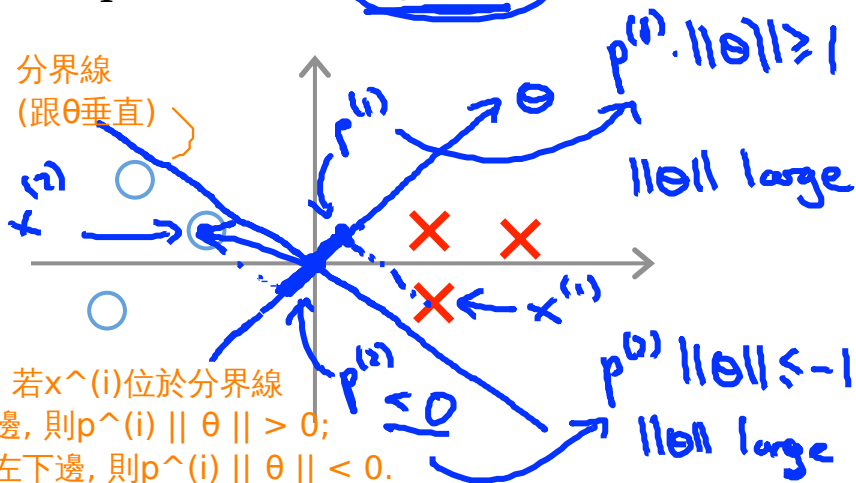# SVM Decision Boundary

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

$\theta_0 \neq 0$

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

$p^{(i)} \cdot \|\theta\| \leq -1$    if $y^{(i)} = 1$    C very large

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

Simplification: $\theta_0 = 0$

分界線
(跟θ垂直)

$p^{(i)} \cdot \|\theta\| \geq 1$

$\|\theta\|$ large

margin    $p^{(i)} \|\theta\| \geq 0$

→ $\|\theta\|$ can be smaller.

$p^{(i)} \|\theta\| \leq -1$

$\|\theta\|$ large

易知: 若x^(i)位於分界線
右上邊, 則p^(i) || θ || > 0;
若在左下邊, 則p^(i) || θ || < 0.

$p^{(i)} \leq 0$

Svm hypothesis

即便 θ0 不等于0
产生大间距分类器的结论 会被证明同样成立

Andrew Ng

# Support Vector Machines

# Kernels I

Machine Learning

# Non-linear Decision Boundary



Predict $y = 1$ if

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

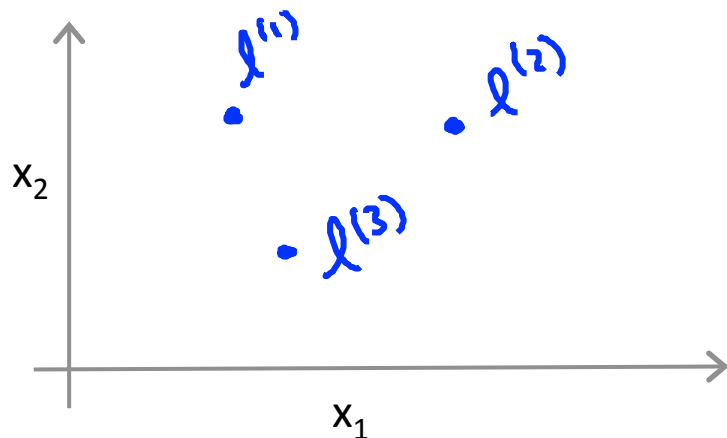$$h_\theta(x) = \begin{cases} 1 & \text{if } \quad \theta_0 + \theta_1 x_1 + \cdots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \cdots$$
$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \cdots$$

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

我打算手动选取一些点 然后将这些点定义为l(1) 再选一个 不同的点 把它定为l(2) 再选第三个点 定为l(3)

# Kernel



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

$\|w\|$

Given $x$:

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

$$f_2 = similarity(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = similarit(x, l^{(3)}) = \exp(\ldots)$$

Kernel (Gaussian kernels)

$k(x, l^{(i)})$

Andrew Ng

# Kernels and Similarity

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

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

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

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

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

$l^{(1)} \rightarrow f_1$

$l^{(2)} \rightarrow f_2$

$l^{(3)} \rightarrow f_3$

**Example:**

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

$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$

$\sigma^2 = 1$     $\sigma^2 = 0.5$     $\sigma^2 = 3$



$\approx 0$

Andrew Ng

对于接近l(1)和l(2)的点 我们的预测值是1. 对于远离 l(1)和l(2)的点 我们最后预测的结果 是等于0的
我们最后会得到 这个预测函数的 判别边界 会像这样
在这个红色的判别边界里面 预测的y值等于1 在这外面预测的y值 等于0

$l^{(1)}$

$l^{(2)}$

$x_2$

predict y=1

predict y=0

$l^{(3)}$

predict y=0

$x_1$

Predict "1" when

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

假设 我已经得到了 这些参数的值

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

$f_1 \approx 1 , \quad f_2 \approx 0, f_3 \approx 0.$ 因为x 接近于l(1) 那么f1 就接近于1

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

$= -0.5 + 1 = 0.5 \geq 0$

$f_1, f_2, f_3 \approx 0$
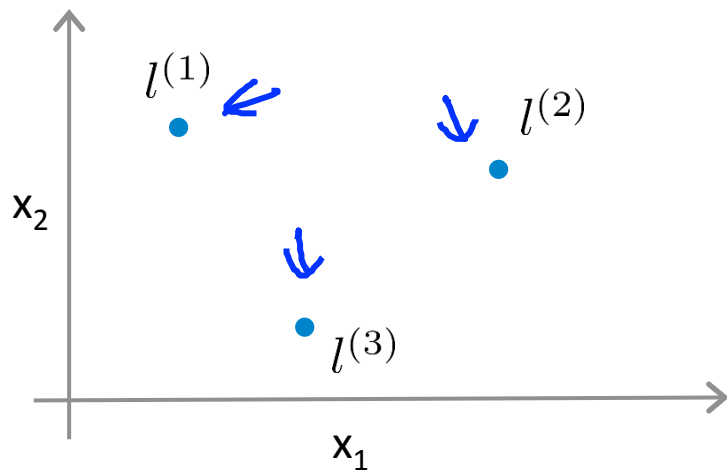
$\theta_0 + \theta_1 f_1 + \dots \approx -0.5 < 0$

Andrew Ng

Support Vector Machines

# Kernels II

Machine Learning

# 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 $l^{(1)}, l^{(2)}, l^{(3)}, \ldots$?



Andrew Ng

**SVM with Kernels**  给定m个训练样本

→ Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})$,

→ choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \ldots, l^{(m)} = x^{(m)}$.

如果我们需要的话 可以添加额外的
特征f0 f0的值始终为1

Given example $x$ :  即x^(i)

$x^{(i)}$

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

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

...

由後面知, 一個f_i就是一個feature.

$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$    $f_0 = 1$

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

$x^{(i)} \rightarrow$

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

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

$x^{(i)}$

$\vdots$

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

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

$f_m^{(i)}$

$x^{(i)} \in \mathbb{R}^{n+1}$   (or $\mathbb{R}^n$)

$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$

$f_0^{(i)} = 1$

Andrew Ng

# SVM with Kernels

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

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

$\theta_0 f_0 + \theta_1 f_1 + \cdots + \theta_m f_m$

Training:

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

$n = m$

$\theta^T f^{(i)}$    $\theta^T f^{(i)}$

$\to \theta_0$

$\sum_j \theta_j^2 = \theta^T \theta \leftarrow$  $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$  (ignore $\theta_0$)

$\|\theta\|^2$

$\theta^T M \theta \leftarrow$    $M = 10,000$

不直接用 θ 的模的平方进行最小化 而是最小化了另一种类似的度量(它主要是为了计算效率, 沒細講)

Andrew Ng

**SVM parameters:**

C ( $= \dfrac{1}{\lambda}$ ). → Large C: Lower bias, high variance.   (small $\lambda$)

→ Small C: Higher bias, low variance.   (large $\lambda$)

$\sigma^2$   Large $\sigma^2$: Features $f_i$ vary more smoothly.

→ Higher bias, lower variance.

對不同的x_i值, 對應的f_i值
差別不大, 所以higher bias

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

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

將核函数用于 逻辑回归时 会变得非常的慢

這是{f_1, f_2, …}這一些
數的分佈. f_1是一個數.

# Support Vector Machines

## Using an SVM

Machine Learning

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):

如果你有大量的特征变量 如果 n 很大 而训练集的
样本数 m 很小, 也许你应该拟合 一个线性的判定边界
不要拟合非常复杂的非线性函数 因为没有足够的数据

不用核函数这个作法 也叫线性核函数

E.g. No kernel ("linear kernel")

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

$$\theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \geq 0$$

→ $n$ large, $m$ small       $x \in \mathbb{R}^{n+1}$

→ Gaussian kernel:

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

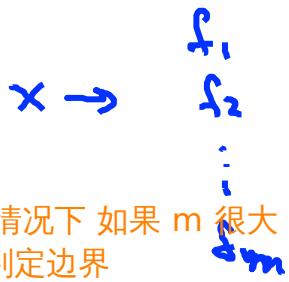Need to choose $\sigma^2$.

$x \in \mathbb{R}^n$, $n$ small

and/or $m$ large

## Kernel (similarity) functions:

目前看来 高斯核函数和线性核函数确实是
最普遍的核函数

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

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

$x^{(i)}$  $l^{(j)} = x^{(j)}$

$f_i$

$x \rightarrow$  $\begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

如果你原来的特征变量 x 是 n 维的 如果 n 很小 并且 理想情况下 如果 m 很大
那么可能你需要用 一个核函数去拟合一个 更复杂的非线性判定边界
那么高斯核函数会是不错的选择

```
return
```

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

归一化

$\|x - l\|^2$

$x \in \mathbb{R}^n$

$V = x - l$

$\|v\|^2 = v_1^2 + v_2^2 + \cdots + v_n^2$

$$= \frac{(x_1 - l_1)^2}{1000 \text{ feet}^2} + \frac{(x_2 - l_2)^2}{1\text{-}5 \text{ bedrooms}} + \cdots + (x_n - l_n)^2$$

Andrew Ng

# Other choices of kernel

Note: Not all similarity functions $\mathrm{similarity}(x, l)$ 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: $k(x, l) = (x^T l)^2$ +0

$$(x^T l + \text{constant})^{\text{degree}}$$

$$(x^T l)^3, \quad (x^T l + 1)^{③}, \quad (x^T l + 5)^{④}$$

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

$$\sin(x, l)$$

## Multi-class classification



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

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 $i = 1, 2, \ldots, K$), get $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(K)}$

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

$y=1 \qquad y=2 \qquad \ldots \qquad \theta \in K$

## 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 $m$): (E.g. $n \geq m$, $n = 10,000$, $m = 10 \cdots 1000$)
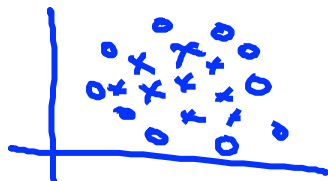
Use logistic regression, or SVM without a kernel ("linear kernel")

If $n$ is small, $m$ is intermediate: ($n = 1-1000$, $m = 10 - 10,000$) ←

Use SVM with Gaussian kernel

If $n$ is small, $m$ is large: ($n = 1-1000$, $m = 50,000+$)

Create/add more features, then use logistic regression or SVM without a kernel    此時高斯核函数 运行起来就会很慢

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

Andrew Ng