# Support Vector Machines

## Optimization objective

Machine Learning

# Alternative view of logistic regression

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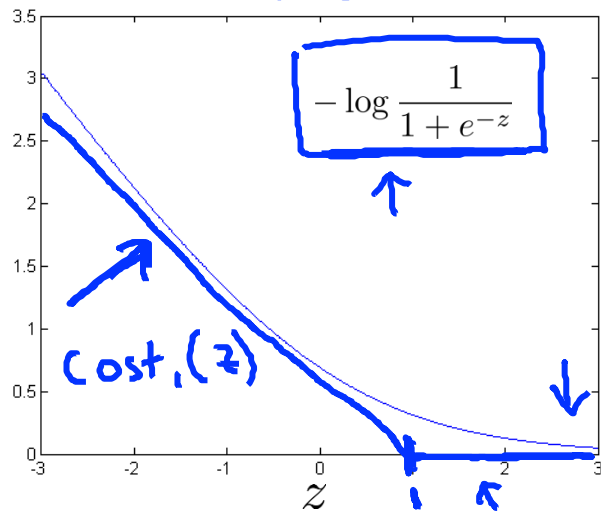
$$h_\theta(x) = g(z)$$

$$z = \theta^T x$$

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

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

Andrew Ng

# Alternative view of logistic regression

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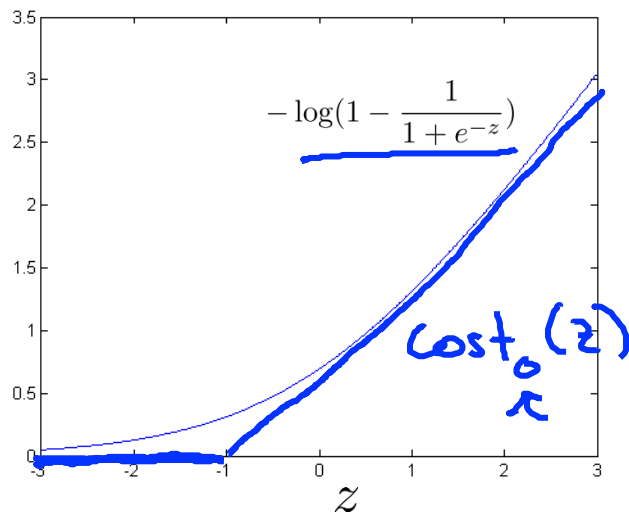
Cost of example: $-(y \log h_\theta(x) + (1-y) \log(1 - h_\theta(x))) \leftarrow$

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

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

$z = \theta^T x$



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

$\text{Cost}_1(z)$

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



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

$\text{Cost}_0(z)$

Andrew Ng

## 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:

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

**SVM hypothesis**

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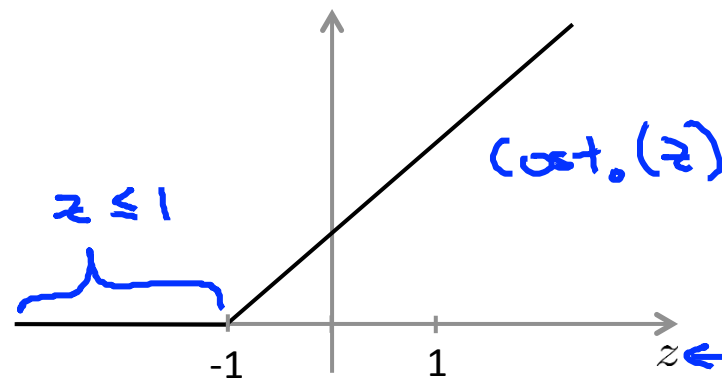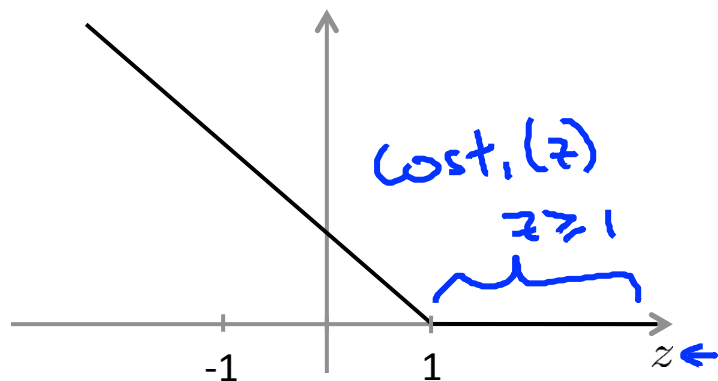
Hypothesis:

# Support Vector Machines

## Large Margin Intuition

Machine Learning

# Support Vector Machine

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



$cost_1(z)$

$z \geq 1$

$cost_0(z)$

$z \leq 1$

-1    1    $z$

-1    1    $z$

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

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

$C = 100,000$

# SVM Decision Boundary

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

$= 0$

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

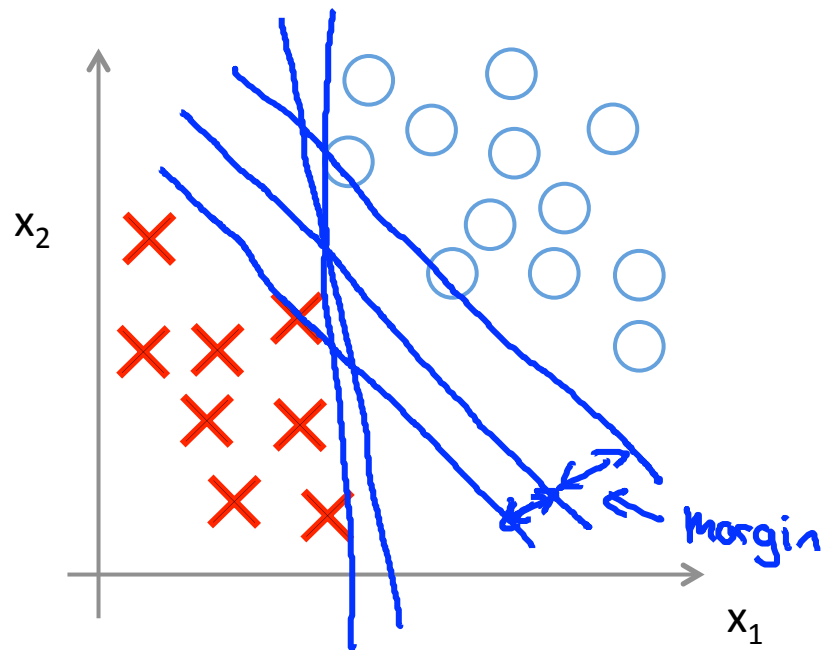$\theta^T x^{(i)} \geq 1$

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

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

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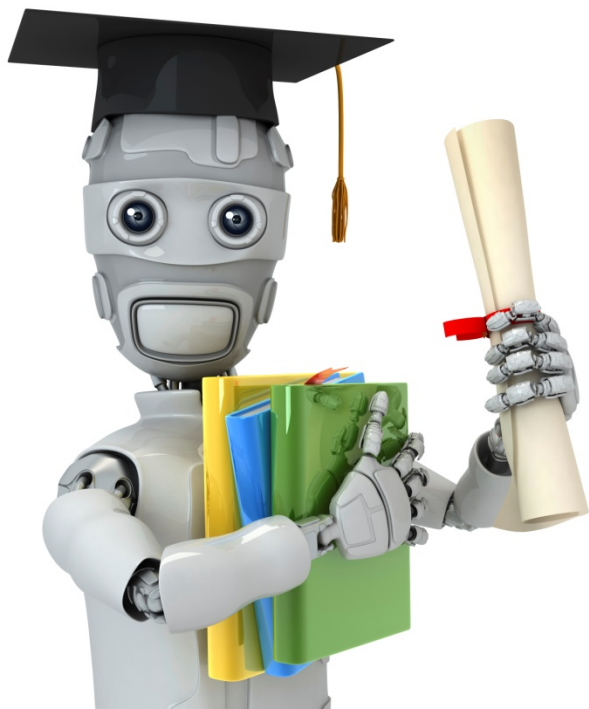
Whenever $y^{(i)} = 0$:

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

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



Large margin classifier

# Large margin classifier in presence of outliers



$$\rightarrow C \text{ very large}$$
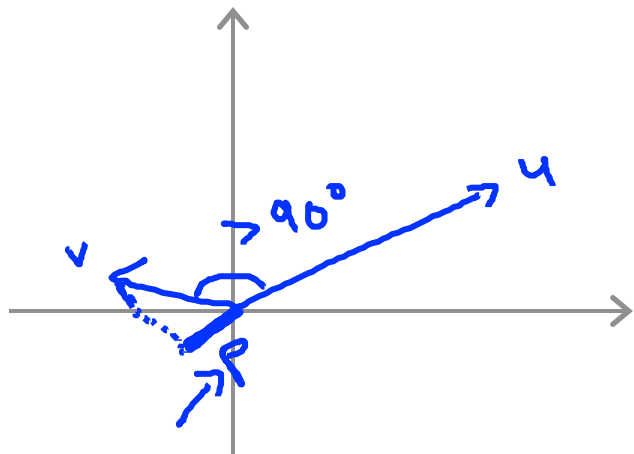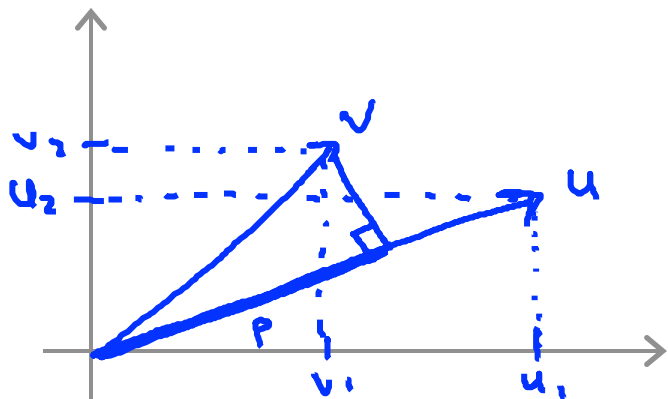
$$\frac{1}{\lambda}$$

$C$ not too large

Machine Learning

# Support Vector Machines

The mathematics behind large margin classification (optional)

# Vector Inner Product



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

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

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

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

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

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

Signed

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

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

$$p < 0$$

Andrew Ng

# SVM Decision Boundary

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

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

$$= \|\theta\|$$

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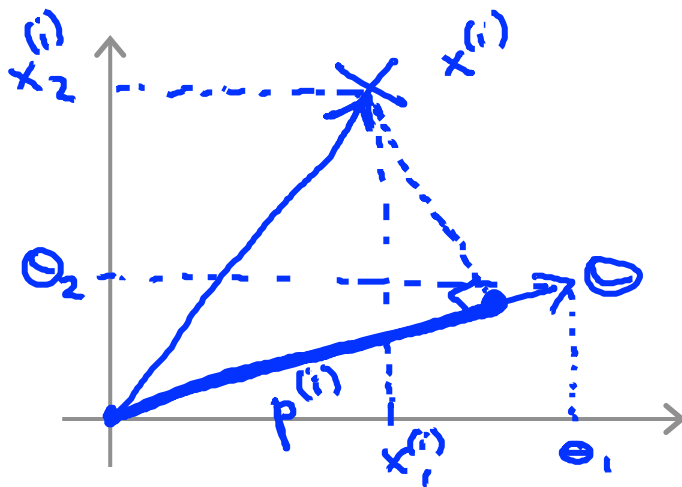
$$\theta^T x^{(i)} \leq -1 \qquad \text{if } y^{(i)} = 0$$

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

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

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

$$u^T v$$



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

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

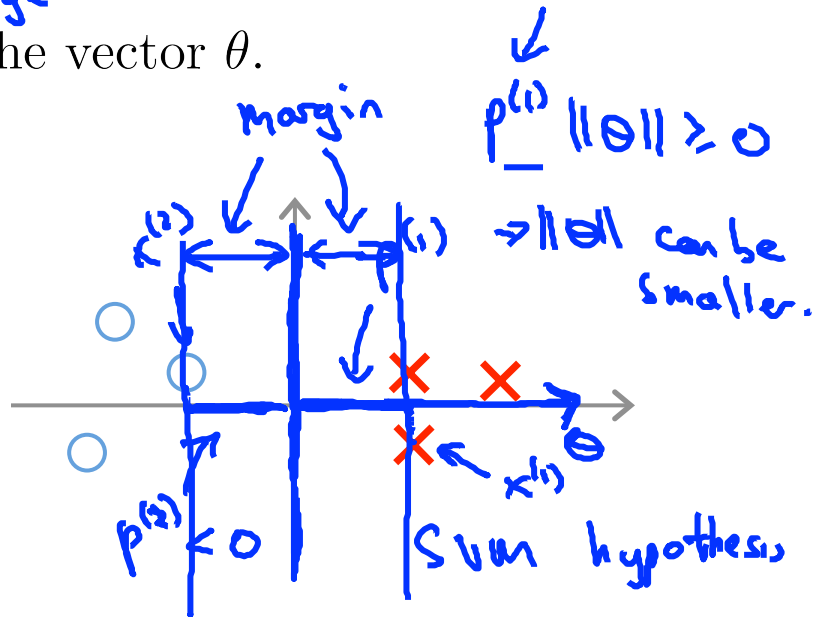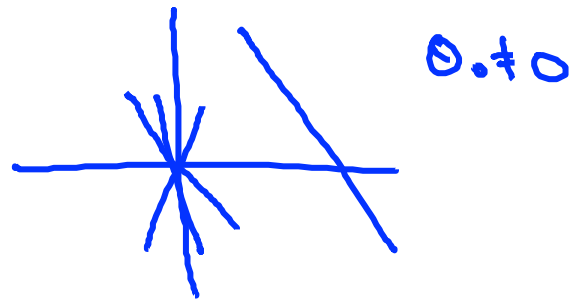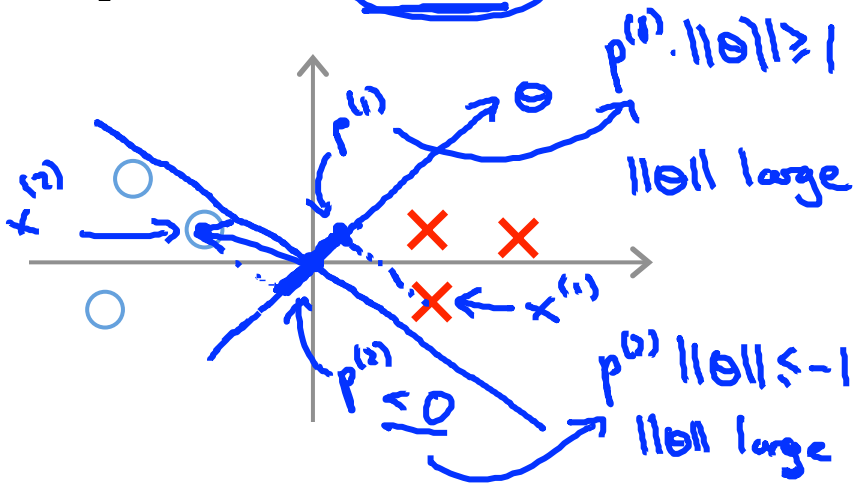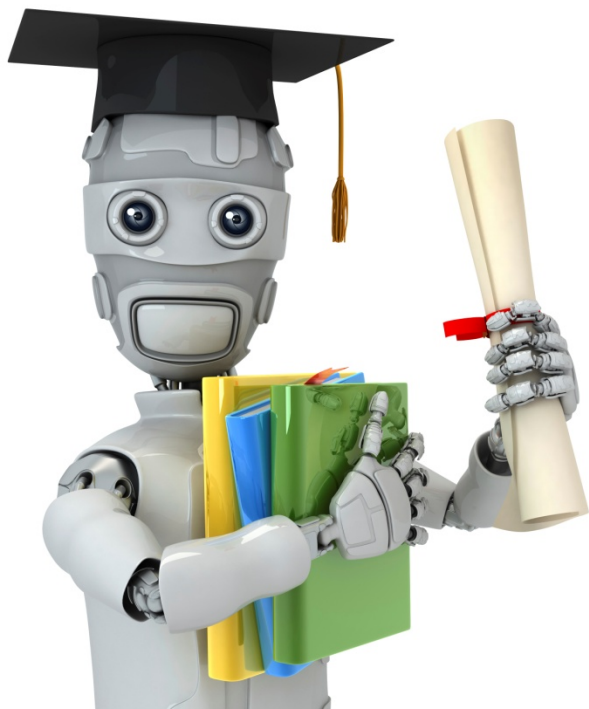**SVM Decision Boundary**

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

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

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

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

Simplification: $\theta_0 = 0$

$\theta_0 \neq 0$

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

$p^{(2)} \|\theta\| \leq -1$
$\|\theta\|$ large

$p \leq 0$

margin

$p^{(i)} \|\theta\| \geq 0$
$\to \|\theta\|$ can be smaller.

$p^{(2)} < 0$

SVM hypothesis

Andrew Ng

# Support Vector Machines

# Kernels I
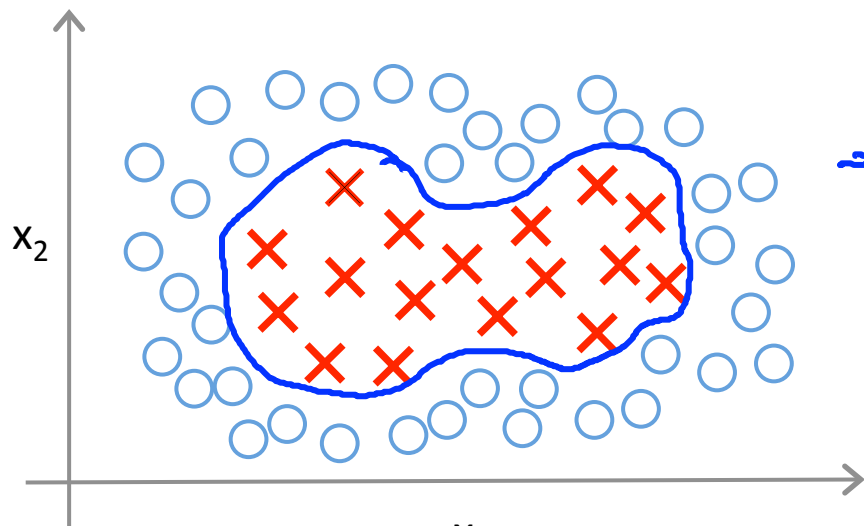
Machine Learning

# Non-linear Decision Boundary



Predict $y = 1$ if

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

$$h_\theta(x) = \begin{cases} 1 & \text{if} \quad \theta_0 + \theta_1 x_1 + \cdots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$
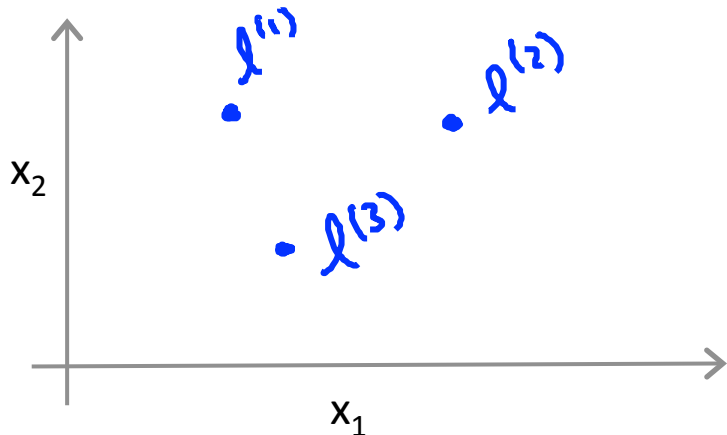
$$\rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \cdots$$
$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \cdots$$

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

# Kernel

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



$$\|w\|$$

Given $x$:

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

$$\|x - l^{(1)}\|^2$$

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

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

kernel (Gaussian kernels)

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

**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 $l^{(1)}$ :

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

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

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

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

x

**Example:**

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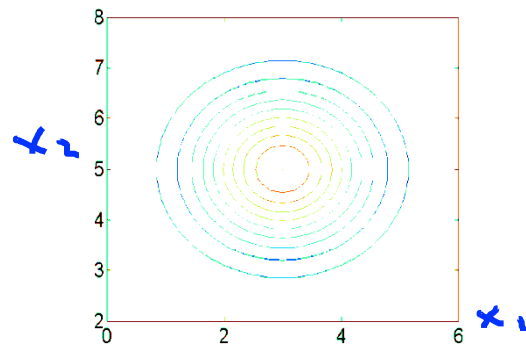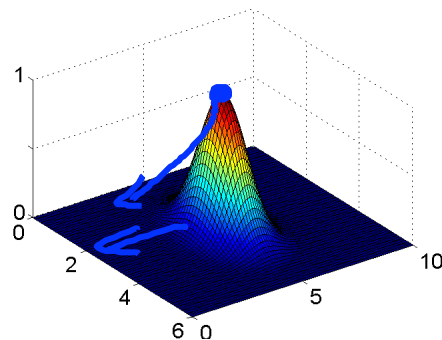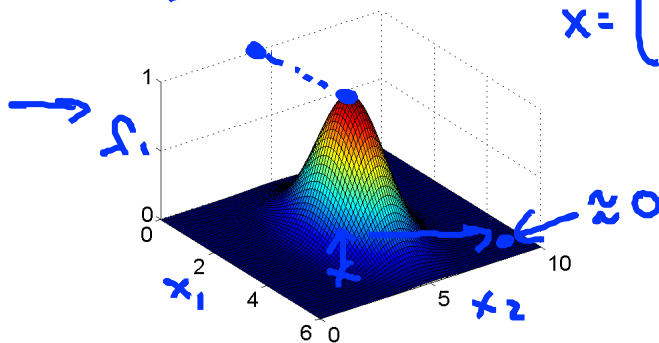
$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$

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

$\approx 0$

$x_1$    $x_2$

$x_2$    $x_1$



Andrew Ng

Predict "1" when

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

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

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

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

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

$f_1, f_2, f_3 \approx 0$

$\theta_0 + \theta_1 f_1 + \ldots \approx -0.5 < 0$

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

$$l^{(1)}$$
$$l^{(2)}$$
$$\vdots$$
$$l^{(m)}$$

**SVM with Kernels**

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

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

Given example $\underline{x}$:

$$\rightarrow f_1 = \text{similarity}(x, l^{(1)}) \quad \leftarrow x^{(i)}$$
$$\rightarrow f_2 = \text{similarity}(x, l^{(2)})$$
$$\ldots$$

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

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

$$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

$$f_1^{(i)} = \sin(x^{(i)}, l^{(1)})$$
$$f_2^{(i)} = \sin(x^{(i)}, l^{(2)}) \quad \leftarrow x^{(i)}$$
$$f_i^{(i)} = \sin(x^{(i)}, l^{(i)}) = \exp\left(-\frac{0}{2\sigma^2}\right) = 1$$
$$\sin(x^{(i)}, l^{(m)})$$

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

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

$$f_0^{(i)} = 1$$

Andrew Ng

# SVM with Kernels

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

$\quad \rightarrow$ Predict "y=1" if $\underline{\theta^T f \geq 0}$ $\qquad \Theta_0 f_0 + \Theta_1 f_1 + \cdots + \Theta_m f_m$

Training:

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

$n = m$

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

$\rightarrow \Theta_0$

$\sum_j \Theta_j^2 = \Theta^T \Theta \quad \Leftarrow \quad \Theta = \begin{bmatrix} \Theta_1 \\ \vdots \\ \Theta_m \end{bmatrix}$ (ignore $\Theta_0$]

$\Theta^T M \Theta \quad \Leftarrow \quad \|\Theta\|^2$ $\qquad M = 10,000$

**SVM parameters:**

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

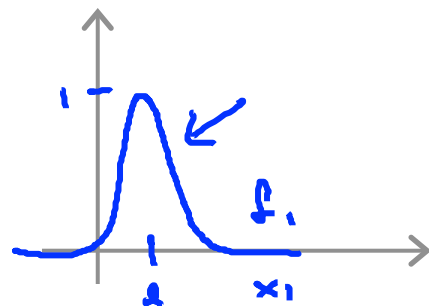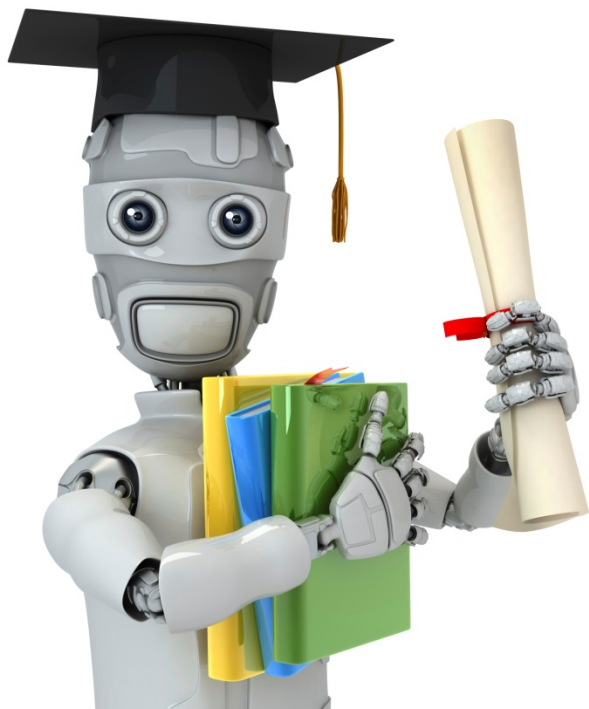→ Small C: Higher bias, low variance.     (large $\lambda$)

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

→ Higher bias, lower variance.

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



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

# 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):

E.g. No kernel ("linear kernel")
Predict "y = 1" if $\theta^T x \geq 0$

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

→ $\underline{n}$ large, $\underline{m}$ 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:**

$x^{(i)}$   $\ell^{(j)} = x^{(j)}$

$f_i$

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

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

$x \to$
$f_1$
$f_2$
$\vdots$
$f_m$

```
return
```

Note: <u>Do perform feature scaling</u> before using the Gaussian kernel.

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

$\boxed{\|x - \ell\|^2}$

$v = x - \ell$

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

$$= \underline{(x_1 - \ell_1)^2} + \underline{(x_2 - \ell_2)^2} + \cdots + (x_n - \ell_n)^2$$

$\underbrace{\qquad}_{1000 \text{ feet}^2}$   $\underbrace{\qquad}_{1-5 \text{ bedrooms}}$

## Other choices of kernel

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

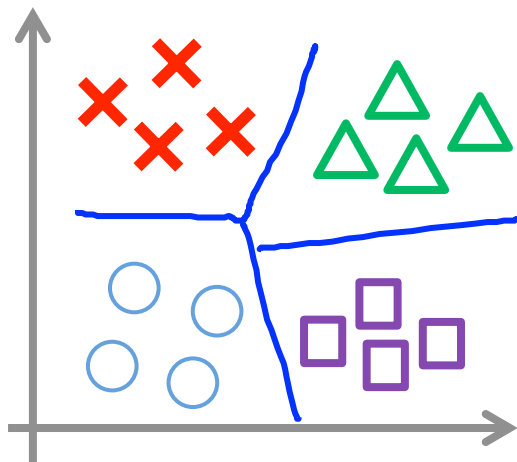Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, l) = (x^T l)^2 + 0$

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

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

- More esoteric: <u>String kernel</u>, <u>chi-square kernel</u>, <u>histogram intersection kernel</u>, …

$$\operatorname{sim}(x, l)$$

# 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 $(\theta^{(i)})^T x$

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

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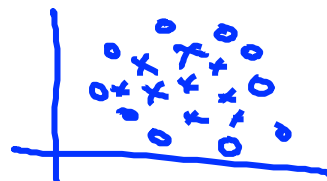
Use logistic regression, or SVM without a kernel ("linear kernel")

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

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