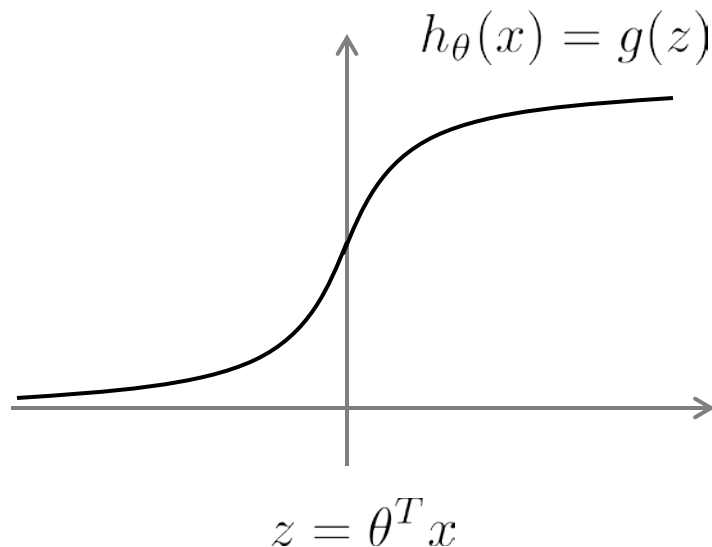# Support Vector Machines

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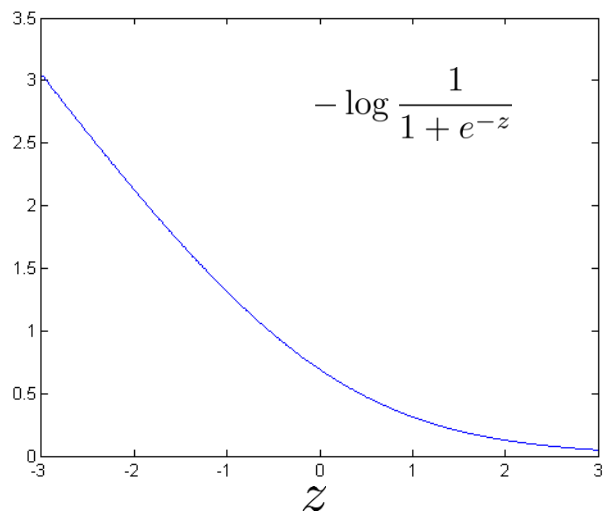If $y = 0$, we want $h_\theta(x) \approx 0$, $\quad \theta^T x \ll 0$

# Alternative view of logistic regression

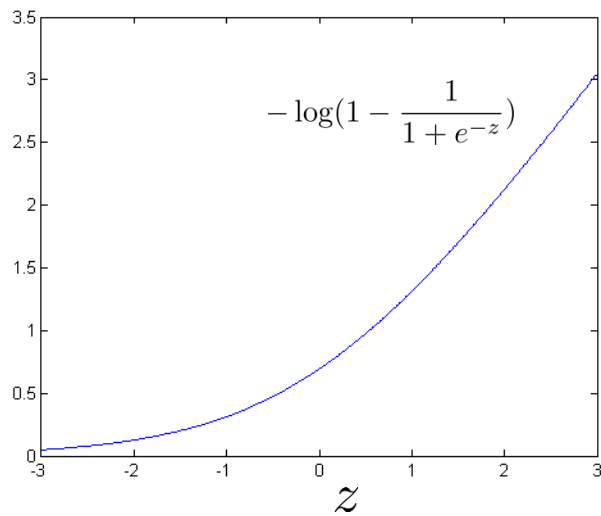Cost of example: $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$

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

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

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

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

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

## Support vector machine

Logistic regression:

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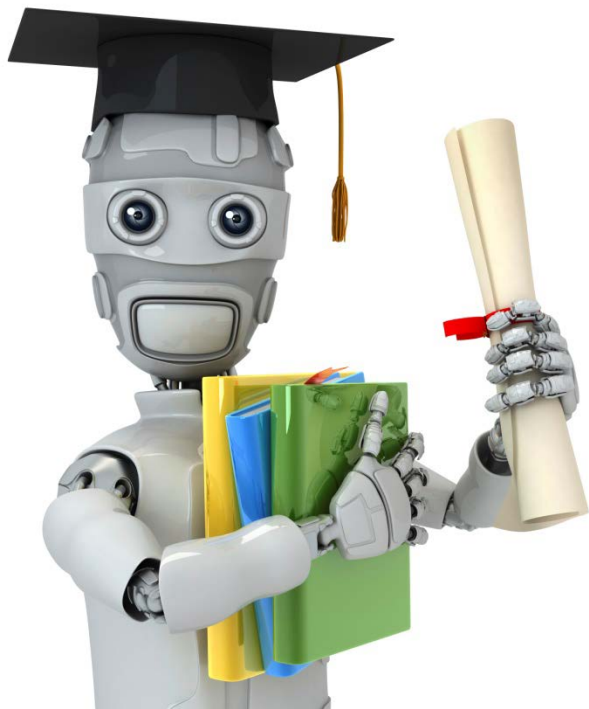
Support vector machine:

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

## SVM hypothesis

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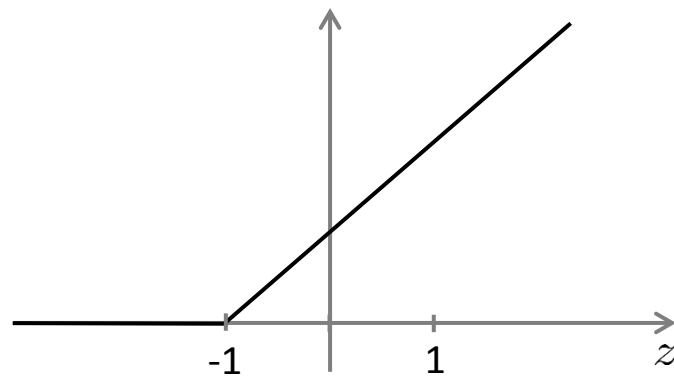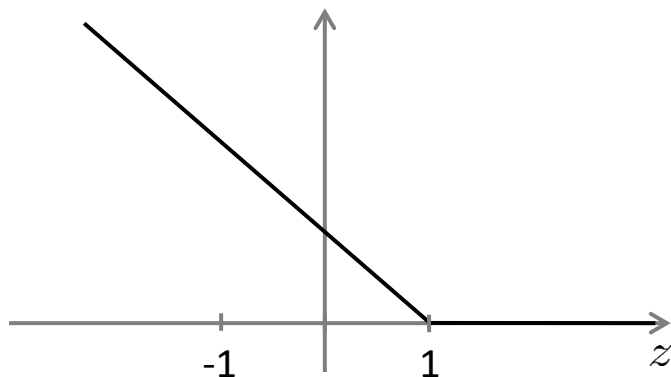
Hypothesis:

# Support Vector Machines

## Large Margin Intuition

Machine Learning

## Support Vector Machine

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



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

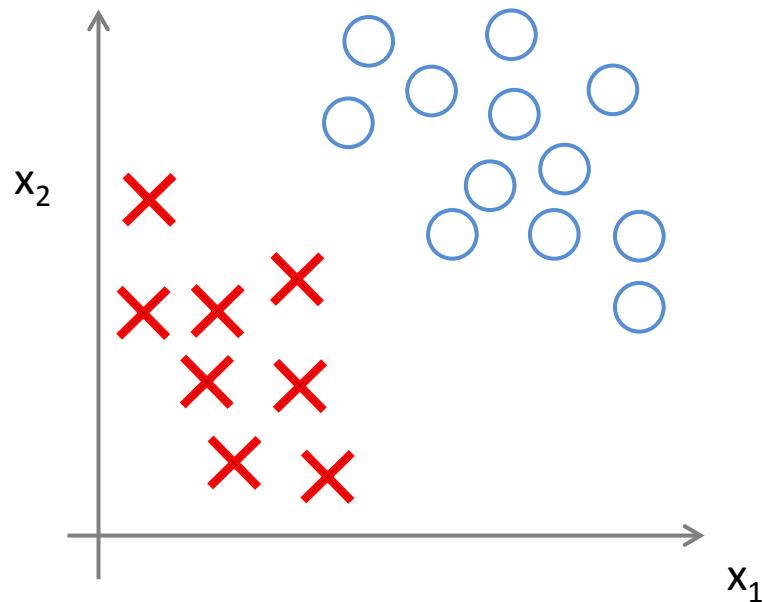If $y = 0$, we want $\theta^T x \leq -1$ (not just $< 0$)

# SVM Decision Boundary

$$\min_\theta C \sum_{i=1}^m \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$
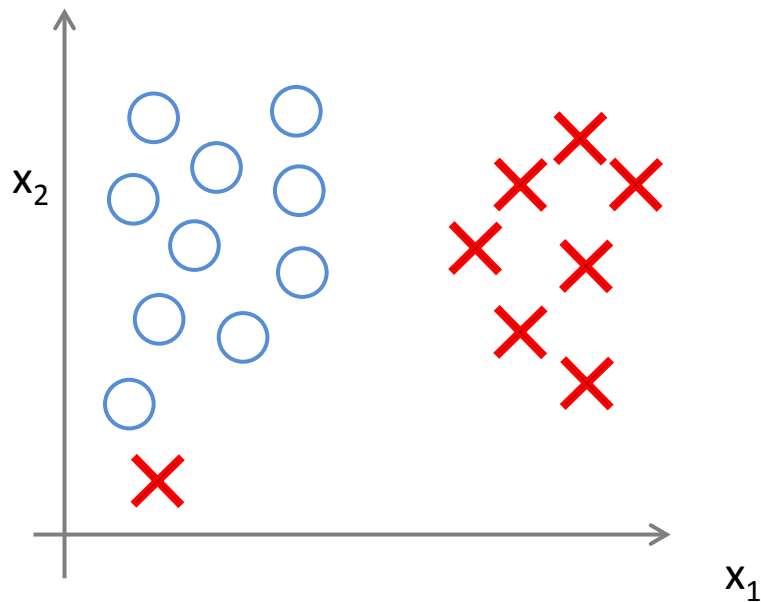
Whenever $y^{(i)} = 1$:
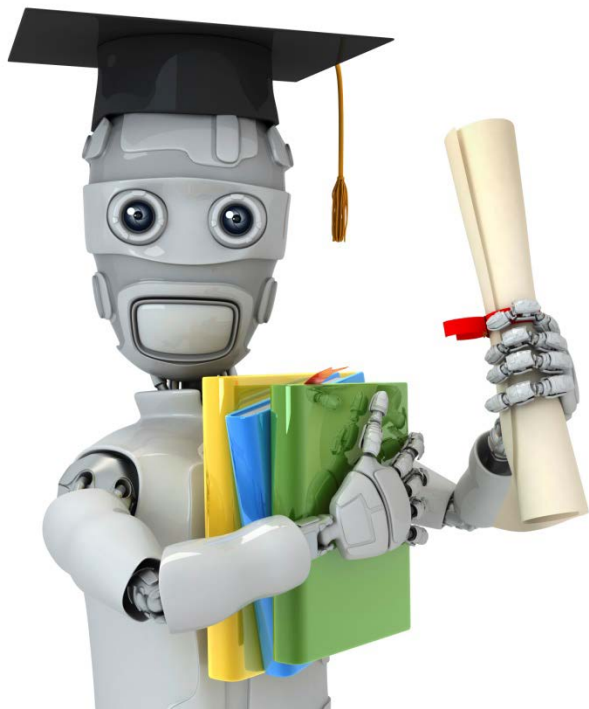
Whenever $y^{(i)} = 0$:

# SVM Decision Boundary: Linearly separable case



Large margin classifier

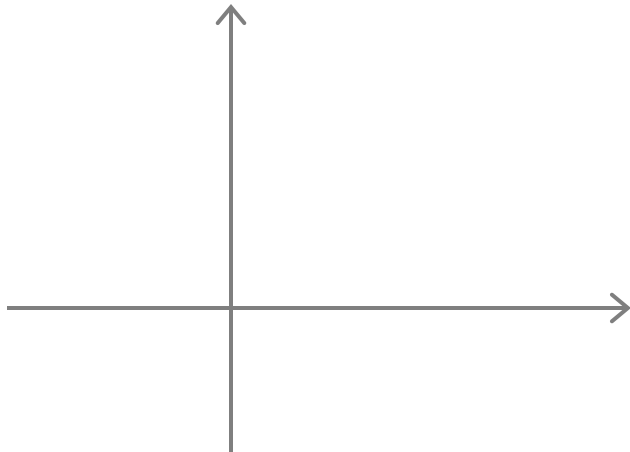# Large margin classifier in presence of outliers

# Support Vector Machines

The mathematics behind large margin classification (optional)

Machine Learning

# Vector Inner Product

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \qquad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$
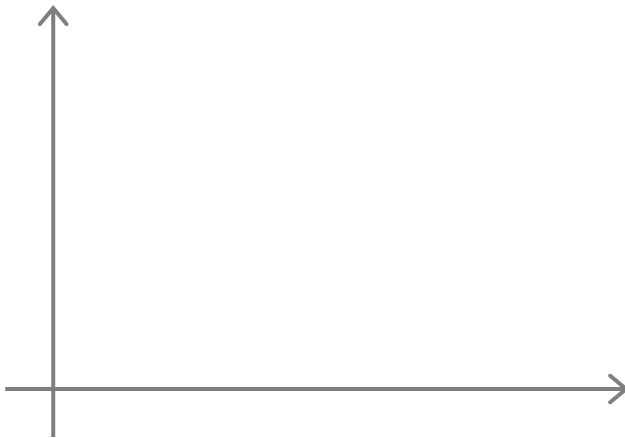
# SVM Decision Boundary

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

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

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$
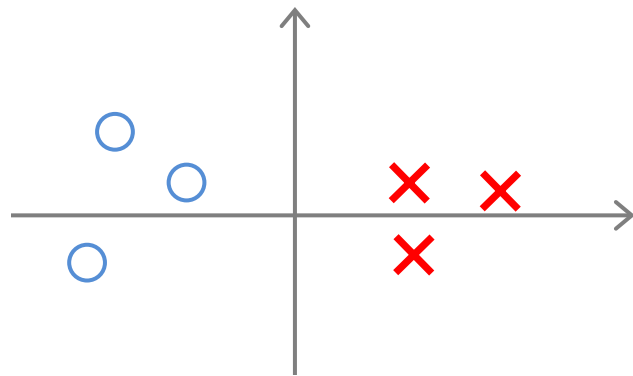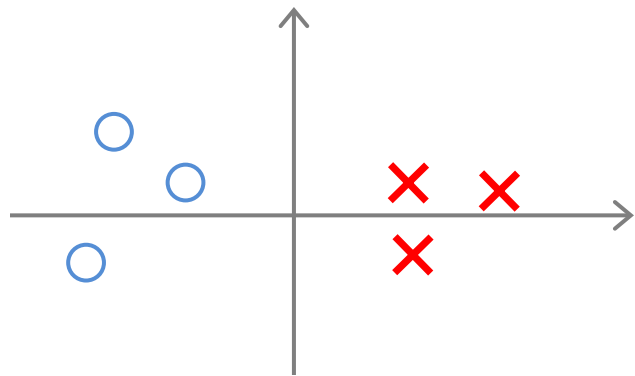
# SVM Decision Boundary
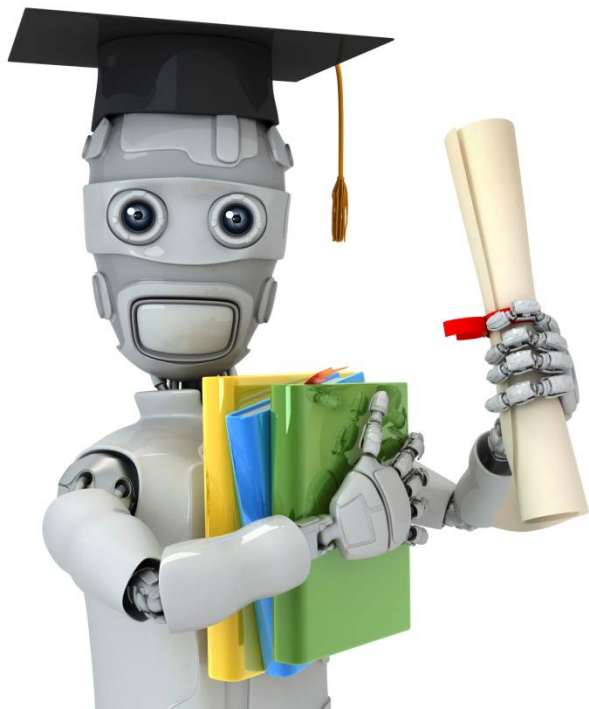
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{n} \theta_j^2$$

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

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

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector $\theta$.
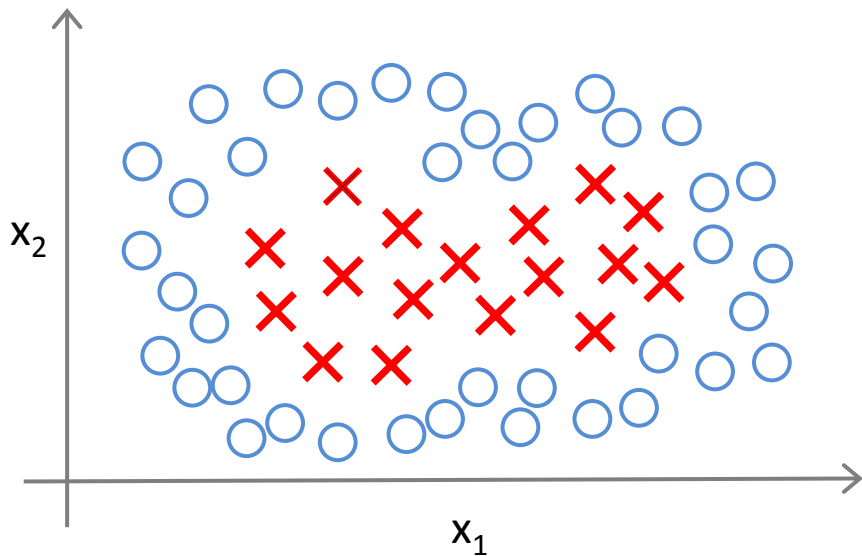
Simplification: $\theta_0 = 0$

# Support Vector Machines

# Kernels I

Machine Learning

# Non-linear Decision Boundary
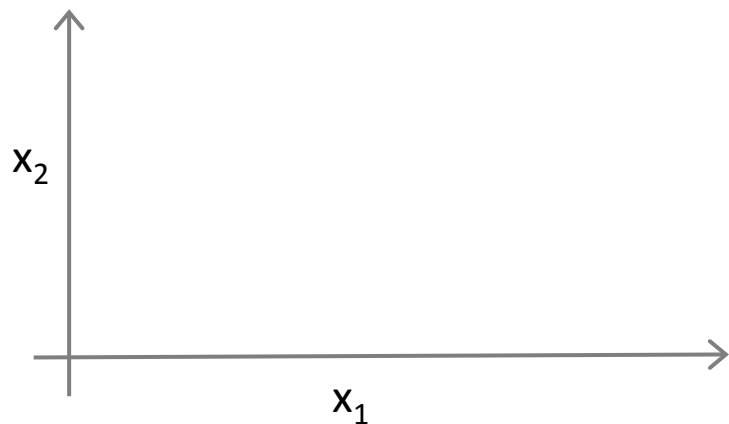


Predict $y = 1$ if

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

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

# Kernel

$x_2$

$x_1$

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

## Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$
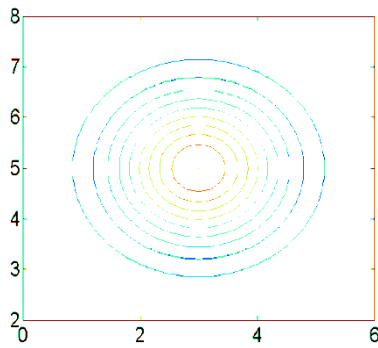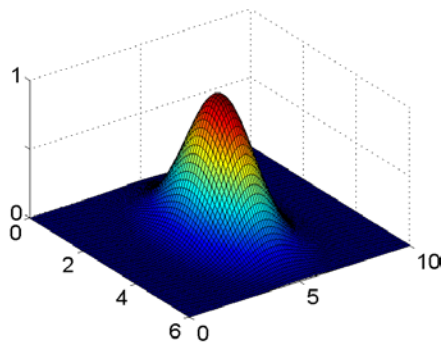
If $x \approx l^{(1)}$ :
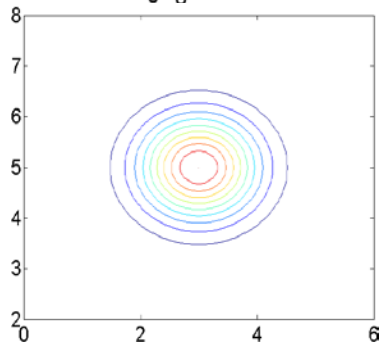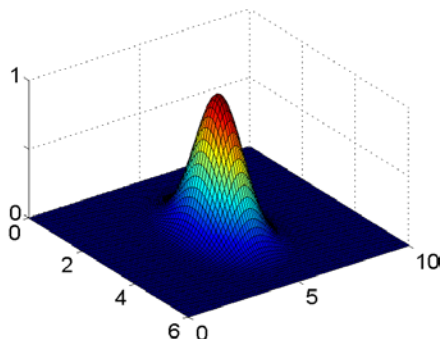
If $x$ if far from $l^{(1)}$ :

**Example:**

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$
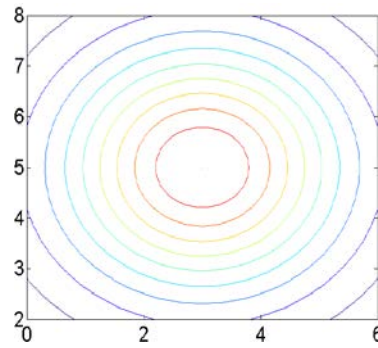
$$\sigma^2 = 1 \qquad\qquad \sigma^2 = 0.5 \qquad\qquad \sigma^2 = 3$$



Andrew Ng

Predict "1" when

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

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

**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 $x$:

$$f_1 = \text{similarity}(x, l^{(1)})$$
$$f_2 = \text{similarity}(x, l^{(2)})$$

$\cdots$

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

## 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 ( $= \frac{1}{\lambda}$ ).    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.

        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 $\theta^T x \geq 0$

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.

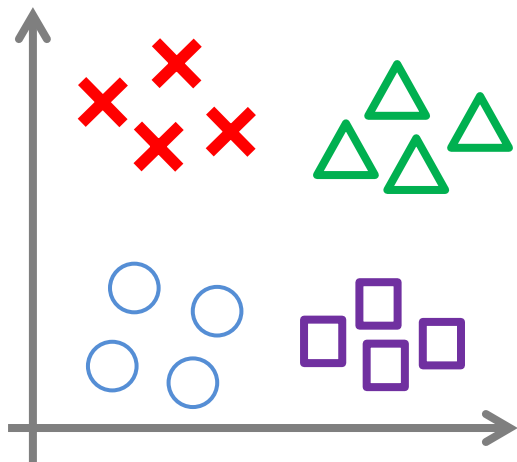**Other choices of kernel**

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:



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

Andrew Ng

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

**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.