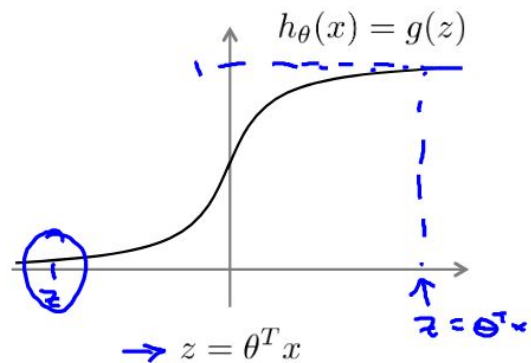# Week Seven
# Support Vector Machines

We have learned logistic regression and neural network on supervised learning. Alongside these we have another algorithm named Support Vector Machines (SVM). Which gives a cleaner, and sometimes more powerful way of learning complex non-linear functions.

## Optimization objective:

### Alternative view of logistic regression

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

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

$$z = \theta^T x$$

$$z = \theta^T x$$
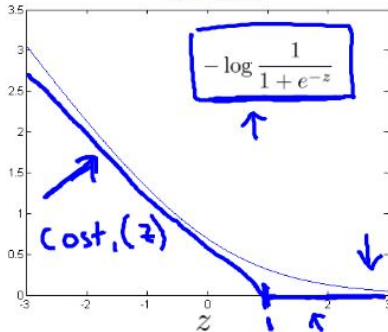
If $y = 1$, we want $h_\theta(x) \approx 1$, $\theta^T x \gg 0$
If $y = 0$, we want $h_\theta(x) \approx 0$, $\theta^T x \ll 0$
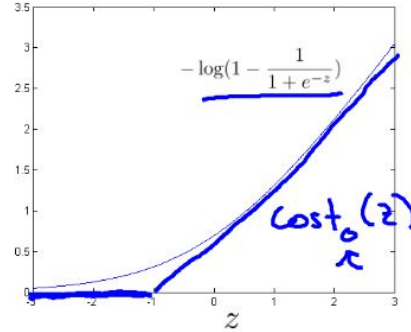
## Alternative view of logistic regression

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 $\theta^T x \ll 0$):

$z = \theta^T x$



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

$Cost_1(z)$



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

$Cost_0(z)$

Andrew Ng

## Support vector machine

Logistic regression:

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

$cost_1(\theta^T x^{(i)})$        $cost_0(\theta^T x^{(i)})$

A

Support vector machine:

$\min_\theta \cancel{\frac{1}{m}} C \sum_{i=1}^m y^{(i)} cost_1(\theta^T x^{(i)}) + (1-y^{(i)}) cost_0(\theta^T x^{(i)}) + \frac{1\cancel{\times}}{2\cancel{\times}}\sum_{j=0}^n \theta_j^2$

B

$\min_u ((u-5)^2 + 1) \times 10 \rightarrow u=5$

$\min_u 10(u-5)^2 + 10 \rightarrow u=5$

$A + \lambda B$ ←

$C A + B$ ←

$C = \frac{1}{\lambda}$

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

Andrew Ng

As m is a constant, so there will be no contribution on min theta. And if C=1/lambda, then the two equations above will give the same output (same min theta).

**SVM hypothesis**

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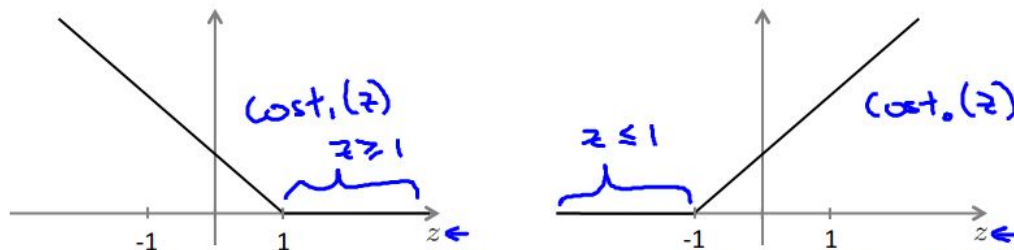
Hypothesis:

# Large Margin Intuition :

In support vector machine we wanna be more confident than logistic regression. So we choose a high threshold value than before.

**Support Vector Machine**

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

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

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

-1   1   $z \leftarrow$

-1   1   $z \leftarrow$

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

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

$C = 100,000$

In the picture below, if we wanna set the boxed part of the equation to zero we have two way:

## SVM Decision Boundary

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

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

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

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

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

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

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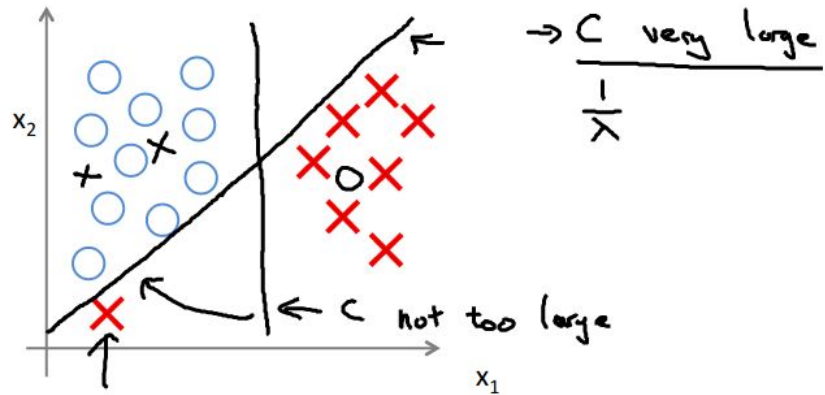
$$\theta^T x^{(i)} \leq -1 \quad if \quad y^{(i)} = 0.$$

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



Large margin classifier

## Large margin classifier in presence of outliers



$\rightarrow$ C very large

$$\frac{1}{\times}$$

← C not too large

# The mathematics behind large margin :

**Vector Inner Product**



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

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

$\|u\| = $ length of vector $u$

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

$p = $ length of projection of $v$ onto $u$.

Signed $\quad u^T v = p \cdot \|u\| \leftarrow \qquad = v^T u$

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

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

$p < 0$

## SVM Decision Boundary

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

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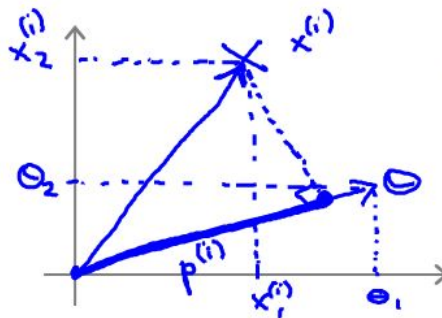
$$= \|\theta\|$$

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

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

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

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

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

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

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

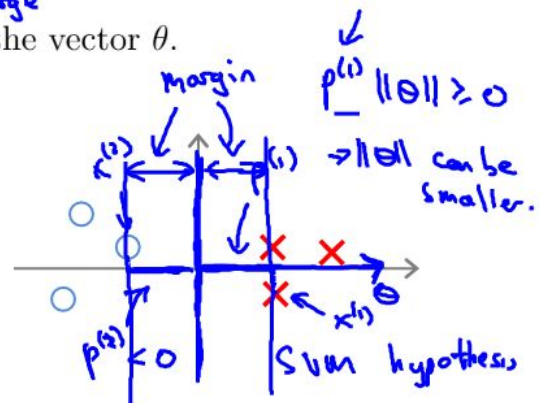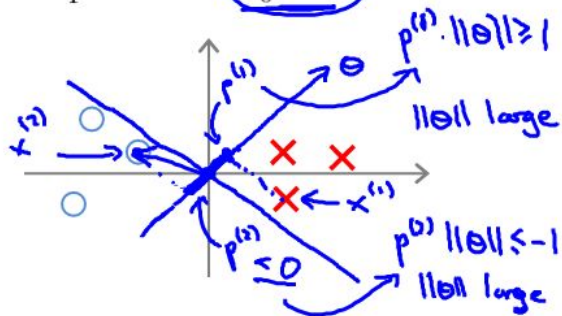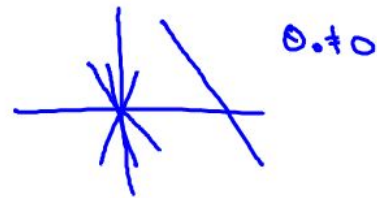$\quad = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$

Andrew Ng

---

## SVM Decision Boundary

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

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

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

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

Simplification: $\theta_0 = 0$

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

$\|\theta\|$ large

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

$\|\theta\|$ large

$\theta_0 = 0$

margin

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

$\Rightarrow \|\theta\|$ can be smaller.

SVM hypothesis

$p^{(2)} < 0$

Andrew Ng

---

If you have any confusion how decision boundary is perpendicular to theta then you can go here and see Q3 to find the answer.

The SVM optimization problem we used is:

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

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

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



where $p^{(i)}$ is the (signed - positive or negative) projection of $x^{(i)}$ onto $\theta$. Consider the training set above. At the optimal value of $\theta$, what is $\|\theta\|$?

○ 1/4

◉ 1/2

○ 1

○ 2

Skip    Submit
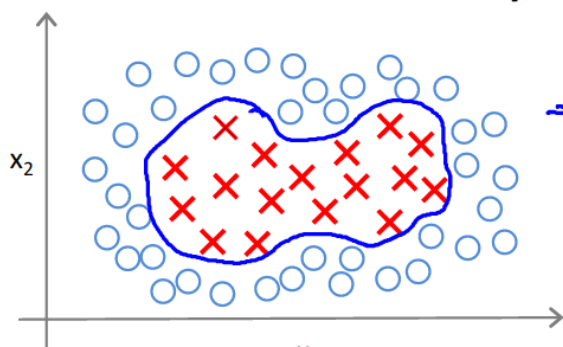
If you can get the clue to solve this quiz then go to <u>here</u> and see Q4b.

# Kernels I :

## Non-linear Decision Boundary



Predict $y = 1$ if

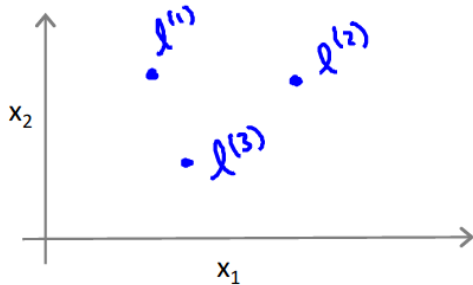$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \cdots \geq 0$$

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

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

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

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 $x$:

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

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

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

Kernel (Gaussian kernels)   $k(x, l^{(i)})$

## Kernels and Similarity

$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\dfrac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\dfrac{\sum_{j=1}^{n}(x_j - l_j^{(1)})^2}{2\sigma^2}\right)$

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

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

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

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

$l^{(1)} \to f_1$
$l^{(2)} \to f_2$
$l^{(3)} \to f_3.$

In the equation sigma squared is the parameter of the Gaussian kernel and as you vary it, you get slightly different effects. As we see with sigma squared the width of the graph is increasing.

**Example:**

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

$\sigma^2 = 1$ $\qquad x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$ $\qquad \sigma^2 = 0.5$ $\qquad\qquad \sigma^2 = 3$



Predict "1" when

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

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

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

$\theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0$
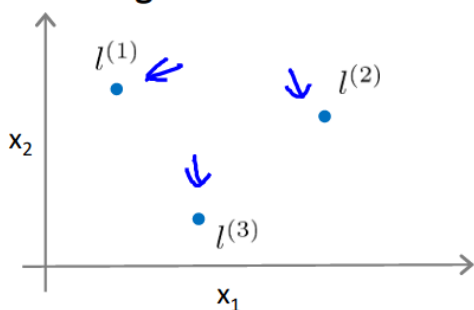
$= -0.5 + 1 = 0.5 \geq 0$

$f_1, f_2, f_3 \approx 0$

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

# Kernel II :

We will choose landmarks as much as our training sets.
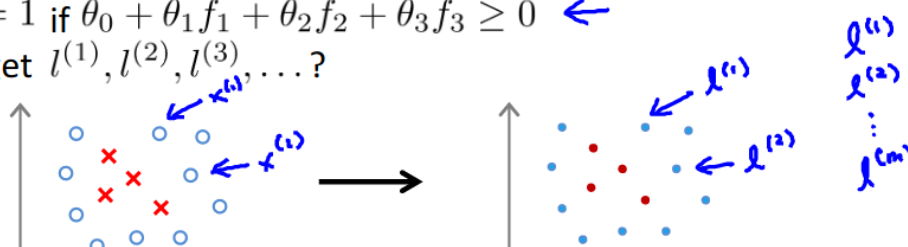
## Choosing the landmarks

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 $x$:

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

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

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

$$x^{(i)} \to \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp(-\frac{0}{2\sigma^2}) = 1 \\ \vdots \\ f_m^{(i)} \quad \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$$

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

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \qquad f_0^{(i)} = 1$$

Andrew Ng

**SVM with Kernels**

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

   → Predict "y=1" if $\underline{\theta^T f \geq 0}$ $\qquad \theta_0 f_0 + \theta_1 f_1 + \cdots + \theta_m f_m$

                                                                           $n = m$

Training:

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

$$\theta^T x^{(i)} \quad \theta^T f^{(i)} \qquad\qquad\qquad \to \theta_0$$

$$\left[ \; - \sum_j \theta_j^2 \; = \; \theta^T \theta \; \leftarrow \; \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignore } \theta_0) \right.$$
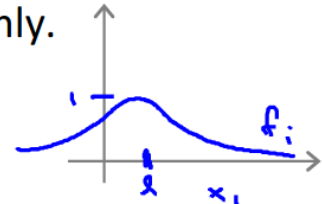
$$\to \theta^T M \theta \leftarrow \|\theta\|^2 \qquad M = 10{,}000$$

 

**SVM parameters:**

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

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

$\sigma^2 \qquad$ 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.

 

# Using an SVM :

We will not write SVM algorithm from scratch. Rather than we will use a built-in library from many of them, like : liblinear, libsvm etc.

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")    $\theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \geq 0$
  Predict "y = 1" if $\theta^T x \geq 0$    → n large, m small    $x \in \mathbb{R}^{n+1}$
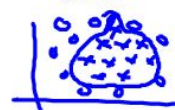
→ Gaussian kernel:
$$f_i = \exp\left(-\frac{||x - l^{(i)}||^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}.$$
Need to choose $\sigma^2$.    $x \in \mathbb{R}^n$, n small  and/or m large

**Kernel (similarity) functions:**    $x^{(i)}$    $l^{(j)} = x^{(j)}$

```
function f = kernel(x1,x2)
```
$$f = \exp\left(-\frac{||\,x1 - x2\,||^2}{2\sigma^2}\right)$$
```
return
```

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

→ Note: Do perform feature scaling before using the Gaussian kernel.
$x \in \mathbb{R}^n$

→ $||x - l||^2$    $v = x - l$
$||v||^2 = v_1^2 + v_2^2 + \cdots + v_n^2$
$= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \cdots + (x_n - l_n)^2$
1000 feet²    1-5 bedrooms

We need to use feature scaling technique if our features may vary in a big range. For example we have a feature which describes flat size another is number of rooms then we need to scale them.

Besides these linear, or Gaussian kernel we have some other kernel like Polynomial kernel, string kernel and so on.

## Other choices of kernel

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.
$\rightarrow$ (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:
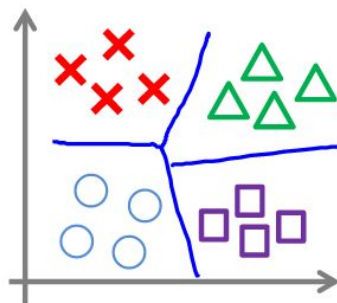- Polynomial kernel: $k(x, l) = (x^T l)^2$, $(x^T l)^3$, $(x^T l + 1)^3$, $(x^T l + 5)^4$, $(x^T l + \text{constant})^{\text{degree}}$

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

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

## Multi-class classification



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

Many SVM packages already have built-in multi-class classification functionality.
$\rightarrow$ 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 $(\theta^{(i)})^T x$

$y=1$  $y=2$  $\cdots$  $\theta = K$

**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.