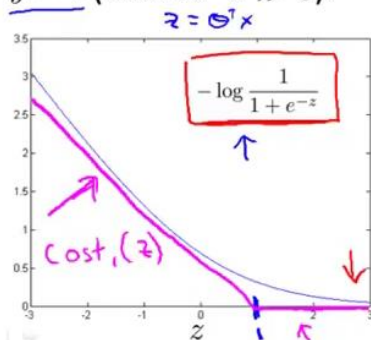Svm 支持向量机

## Alternative view of logistic regression
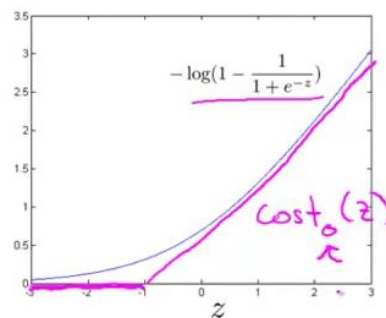
$(x,y)$

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

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

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

$z = \theta^T x$

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

$Cost_1(z)$

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

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

$Cost_0(z)$

Andrew

## Support vector machine

Logistic regression:

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

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

A

Support vector machine:

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

B

$\min_u \left((u-5)^2 + 1\right) \xrightarrow{\times 10} u = 5$

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

$A + \lambda B$

$\Rightarrow C A + B$

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

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

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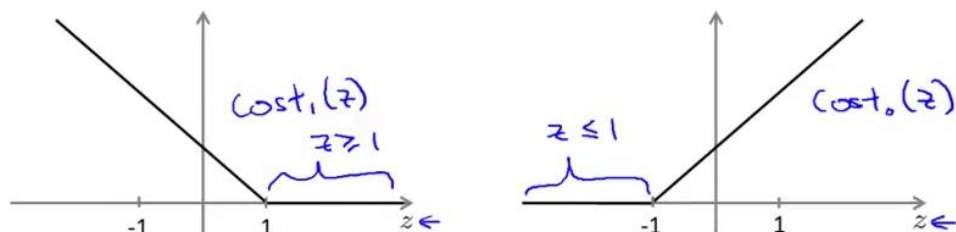
$$\Rightarrow \min_\theta C \sum_{i=1}^{m}\left[y^{(i)} cost_1(\theta^T x^{(i)}) + (1-y^{(i)}) cost_0(\theta^T x^{(i)})\right] + \frac{1}{2}\sum_{i=1}^{n}\theta_j^2$$

Hypothesis:

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

## Support Vector Machine

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

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

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

-1   1   $z \leftarrow$     -1   1   $z \leftarrow$

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

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

## SVM Decision Boundary

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

$\phantom{x} = 0$

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

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

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

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

$\min_{\theta} \cancel{C \cdot 0} + \frac{1}{2} \sum_{i=1}^{n} \Theta_j^2$

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

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

大间距分类器

## Large margin classifier in presence of outliers

$\to$ C very large

$\frac{1}{\lambda}$

$x_2$

$\leftarrow$ C not too large

$x_1$

数学原理

## Vector Inner Product

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

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

$\|u\| = $ length of vector $u$
$\quad = \sqrt{u_1^2 + u_2^2} \quad \in \mathbb{R}$

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

Signed $\quad u^T v = \underline{p} \cdot \underline{\|u\|} \leftarrow \quad = v^T u$
$\quad\quad = u_1 v_1 + u_2 v_2 \leftarrow \quad p \in \mathbb{R}$

$u^T v = p \cdot \|u\|$
$p < 0$

## SVM Decision Boundary

$\omega = (\sqrt{\omega})^2$

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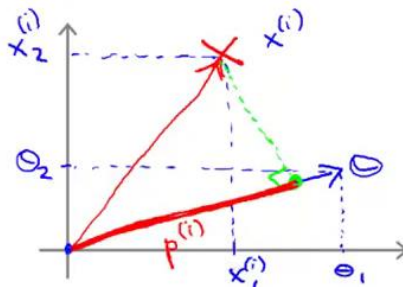
$\quad\quad\quad\quad\quad = \|\theta\|$

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

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

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

$\theta^T x^{(i)} = ?$
$\quad\uparrow \quad\uparrow$
$\quad u^T v$

$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\| \leftarrow$
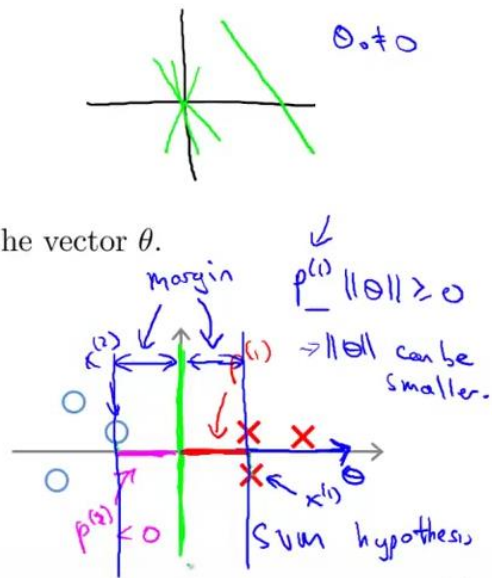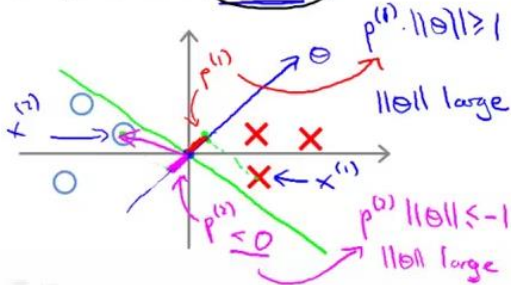$\quad = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \leftarrow$

## SVM Decision Boundary

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

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

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

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

Simplification: $\boxed{\theta_0 = 0}$

$\theta_0 \neq 0$

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

$\|\theta\|$ large

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

$\|\theta\|$ large

margin

$p^{(1)} \|\theta\| \geq 0$

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

$p^{(2)} < 0$

SVM hypothesis

Kernels 核函数

## Kernels and Similarity

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

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

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

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

$\times$

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

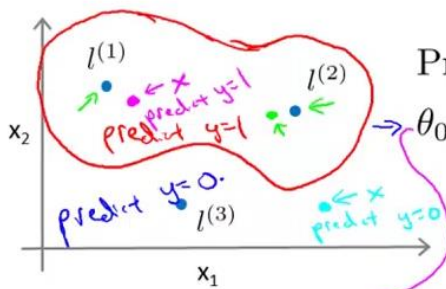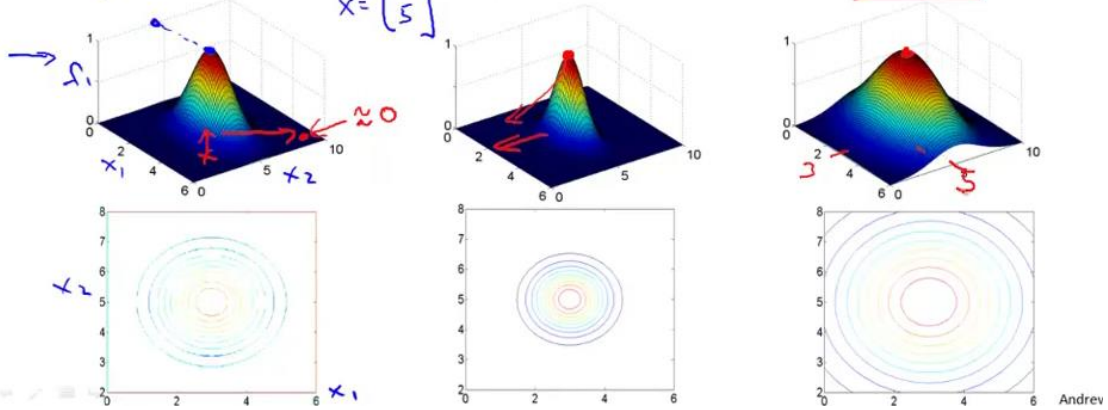$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

**Example:**

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

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

$f_1$    $\approx 0$



$x_1$    $x_2$    $x_1$

Predict "1" when

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

$x_2$    $l^{(1)}$   predict $y=1$   $l^{(2)}$
predict $y=1$
predict $y=0$.
$l^{(3)}$   predict $y=0$

$x_1$

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

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

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

$f_1, f_2, f_3 \approx 0$

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

## 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 $\underline{x}$:

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

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

$\ldots$

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

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

$x^{(i)} \rightarrow$ 
$f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)})$
$f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)})$
$\vdots$
$f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp(-\frac{0}{2\sigma^2}) = 1$   $\leftarrow x^{(i)}$
$f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)})$

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

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

$f_0^{(i)} = 1$

Andrew

---

## SVM with Kernels

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

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

$\underline{\theta_0 f_0 + \theta_1 f_1 + \cdots + \theta_m f_m}$   $n = m$

Training:

→ $\min_{\theta} C \sum_{i=1}^{m} y^{(i)} \text{cost}_1(\underline{\theta^T f^{(i)}}) + (1 - y^{(i)}) \text{cost}_0(\underline{\theta^T f^{(i)}}) + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$

$\not{=} m$

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

$\rightarrow \theta_0$

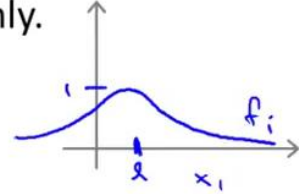$\therefore \quad \sum_j \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$   (ignore $\theta_0$)

$= \|\theta\|^2$

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

## 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.

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

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

线性核函数，高斯核函数

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 $\theta^T x \geq 0$

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

→ $n$ 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
```

$f_i$

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

→ 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

**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: $k(x, l) = (x^T l)^2$

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

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

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

**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.