

Recommended approach

- Start with a simple algorithm that you can implement quickly.
Implement it and test it on your cross-validation data.
- Plot learning curves to decide if more data, more features, etc.
are likely to help.
- Error analysis: Manually examine the examples (in cross
validation set) that your algorithm made errors on. See if you
spot any systematic trend in what type of examples it is
making errors on.

Alternative view of logistic regression

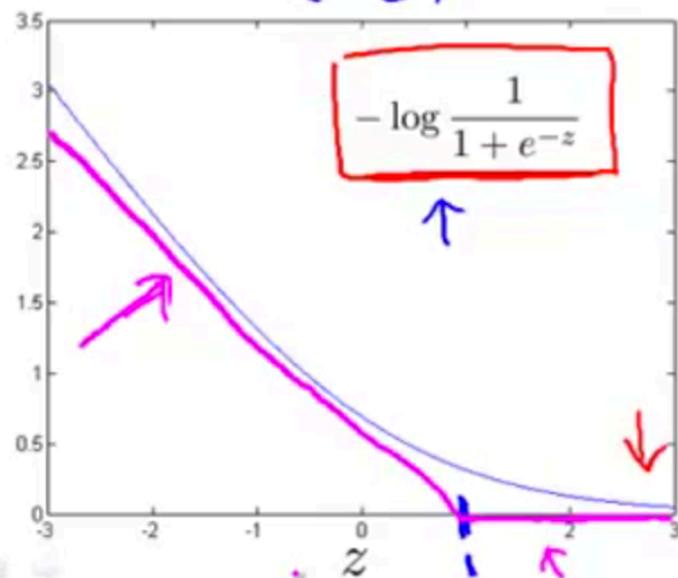
(x, y)

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

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

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

$$z = \Theta^T x$$



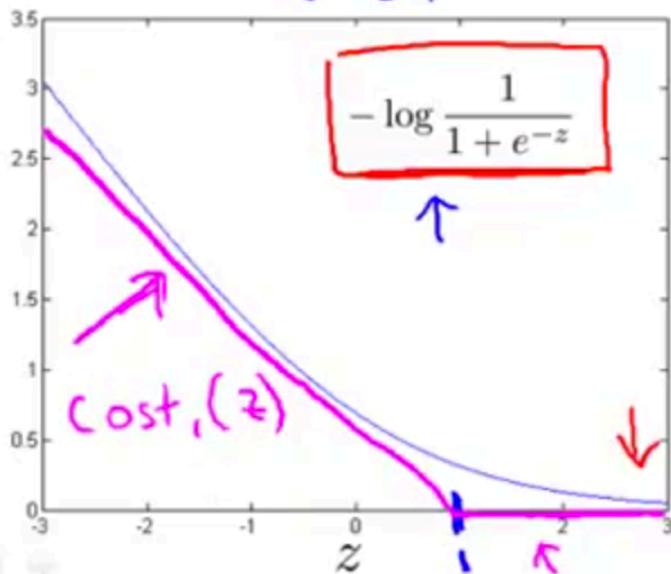
Alternative view of logistic regression

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

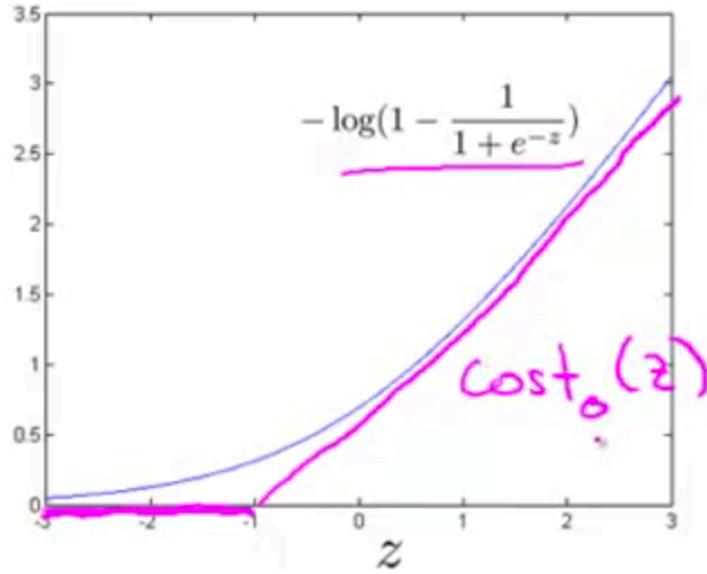
$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$
 ←

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

$$z = \theta^T x$$



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



Support vector machine

Logistic regression:

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


Support vector machine:

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

$$\min_u \frac{(u-s)^2 + 1}{10} \rightarrow u=5 \quad \left| \begin{array}{l} \underline{A} + \underline{\lambda} \underline{B} \leftarrow \\ \rightarrow C \underline{A} + \underline{B} \leftarrow \end{array} \right. \quad C = \frac{1}{\lambda}$$

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

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

SVM hypothesis

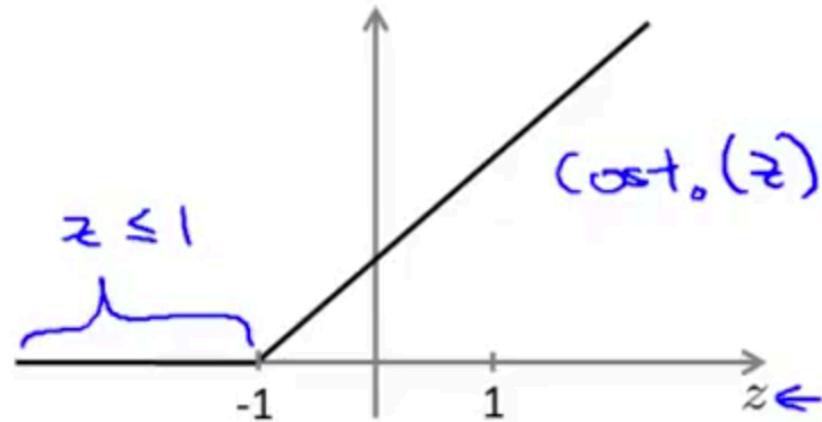
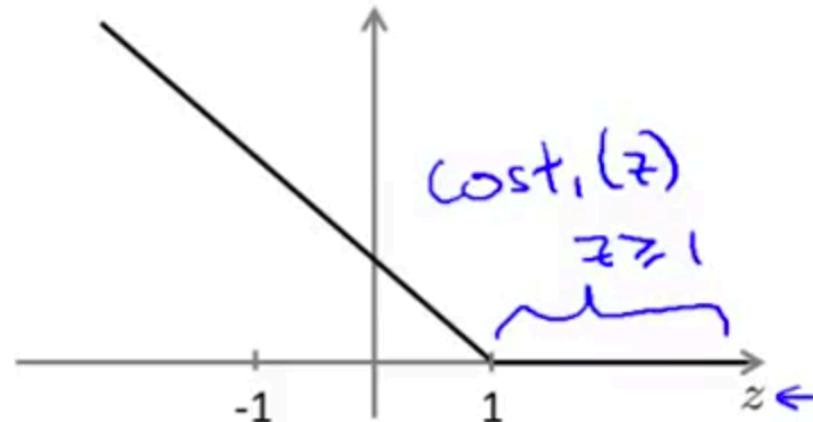
$$\Rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis:

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

Support Vector Machine

$$\rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \underbrace{\text{cost}_1(\theta^T x^{(i)})}_{z \geq 1} + (1 - y^{(i)}) \underbrace{\text{cost}_0(\theta^T x^{(i)})}_{z \leq 1} \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



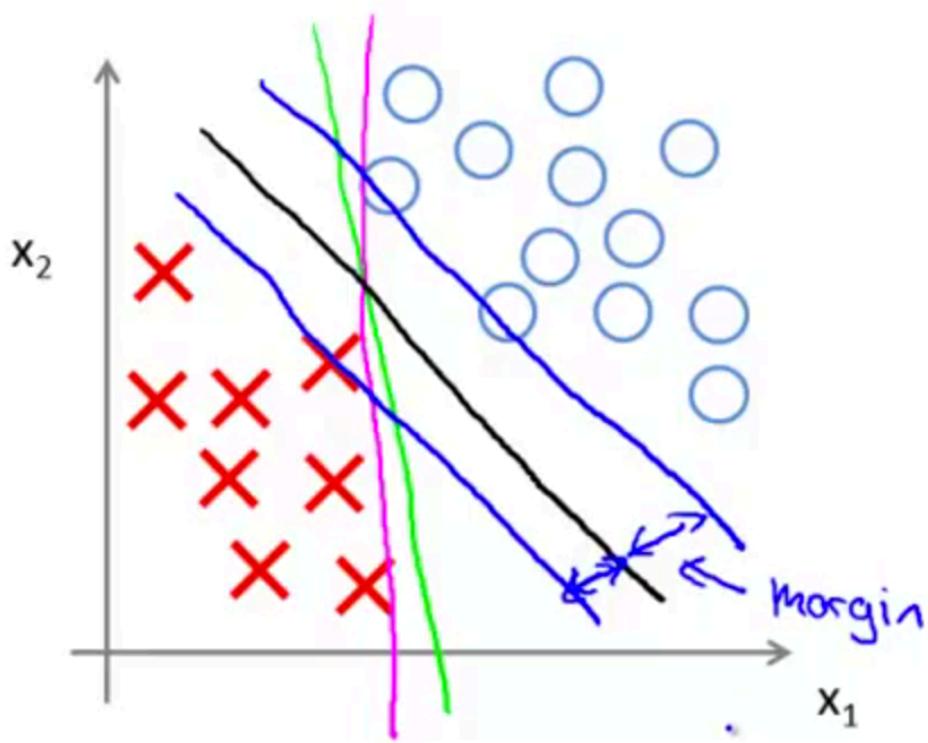
\rightarrow If $y = 1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

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

$\theta^T x \geq 1$

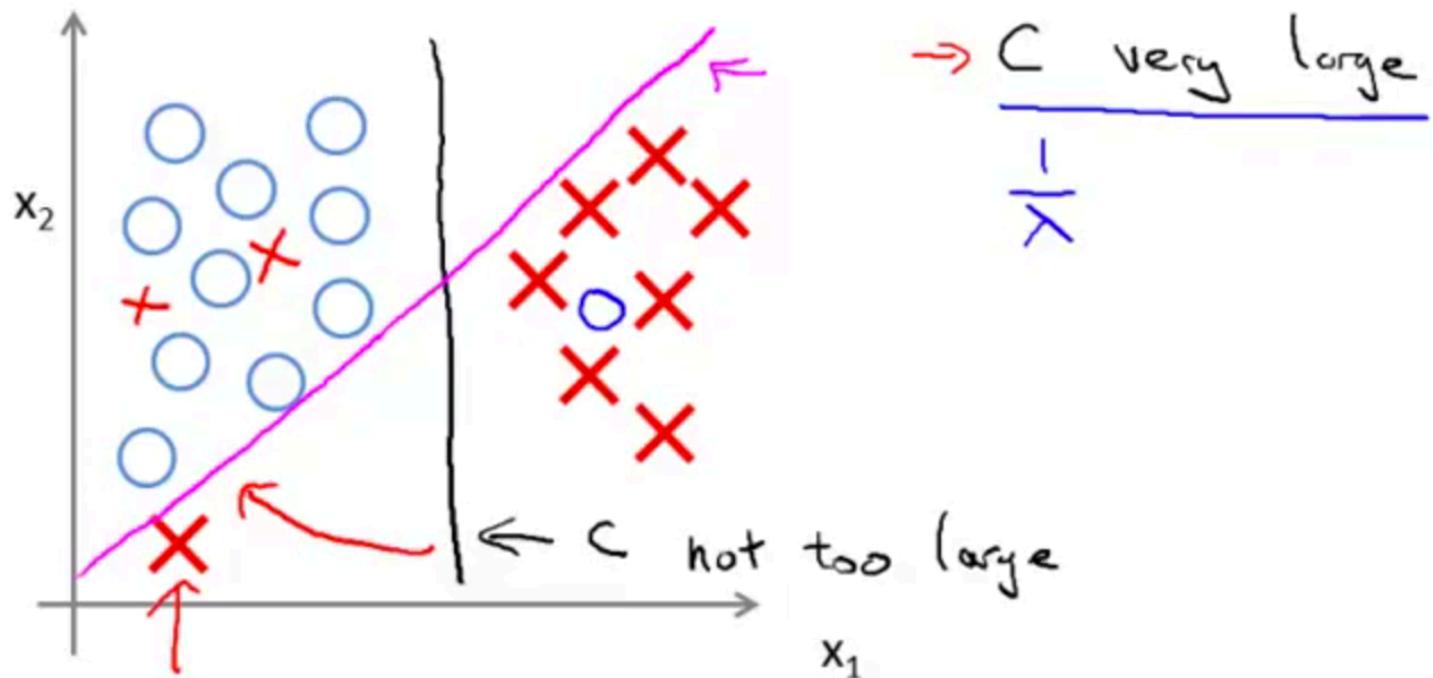
$\theta^T x \leq -1$

SVM Decision Boundary: Linearly separable case

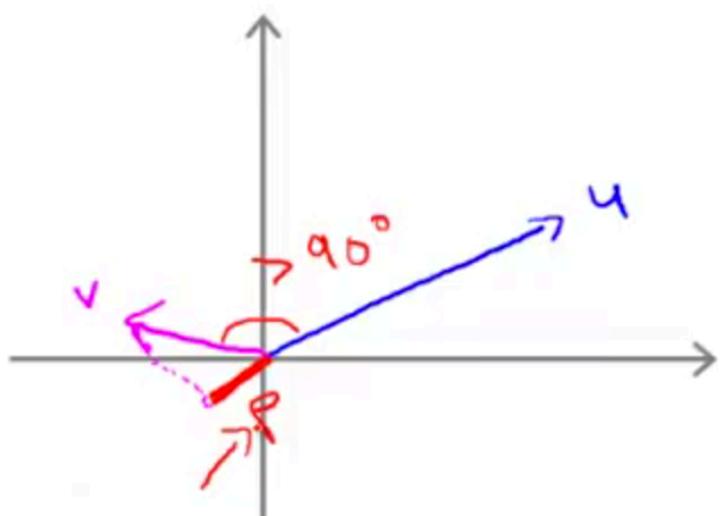
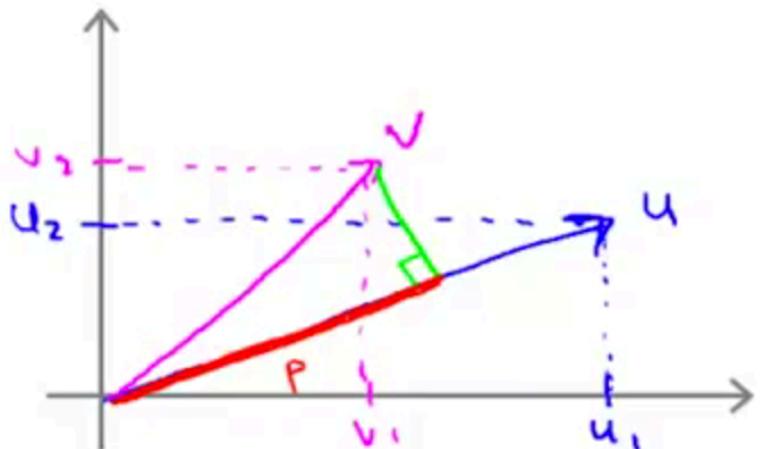


Large margin classifier

Large margin classifier in presence of outliers



Vector Inner Product



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

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

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

p = length of projection of v onto u .

Signed

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

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

$$p < 0$$

SVM Decision Boundary

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

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\Theta_0^2 + \Theta_1^2 + \Theta_2^2) = \frac{1}{2} \left(\underbrace{\Theta_0^2 + \Theta_1^2 + \Theta_2^2}_{= \|\theta\|^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

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

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

Simplification: $\Theta_0 = 0$. n=2

$$\begin{bmatrix} \cancel{\Theta_0} \\ \Theta_1 \\ \Theta_2 \end{bmatrix} \quad \Theta_0 = 0$$



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

SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\Theta_0^2 + \Theta_1^2 + \Theta_2^2) = \frac{1}{2} \left(\sqrt{\Theta_0^2 + \Theta_1^2 + \Theta_2^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

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

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

Simplification: $\Theta_0 = 0$. n=2

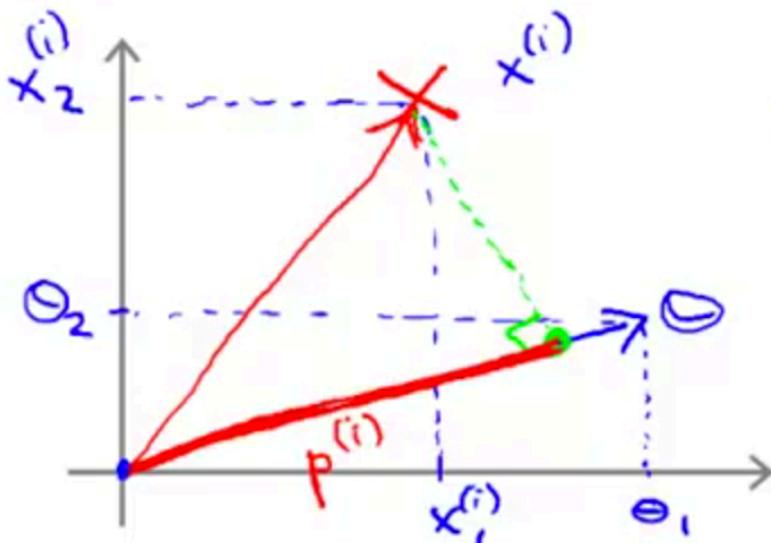
$$= \|\theta\|$$

$$\begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \end{bmatrix}, \Theta_0 = 0$$

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

↑

U^T V



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

$$= \Theta_0 x_0 + \Theta_1 x_1 + \Theta_2 x_2 \leftarrow$$

SVM Decision Boundary

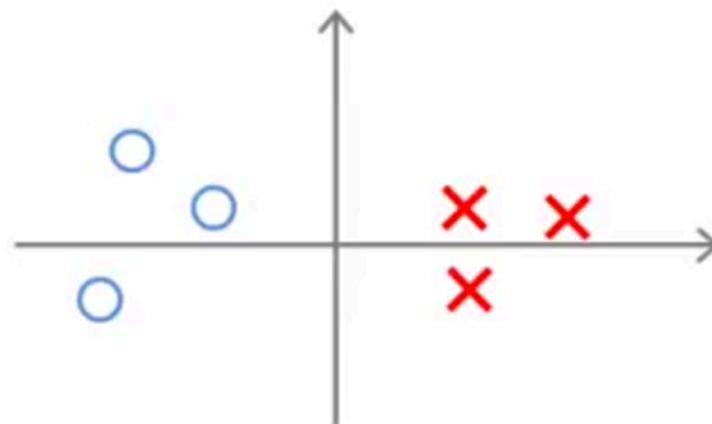
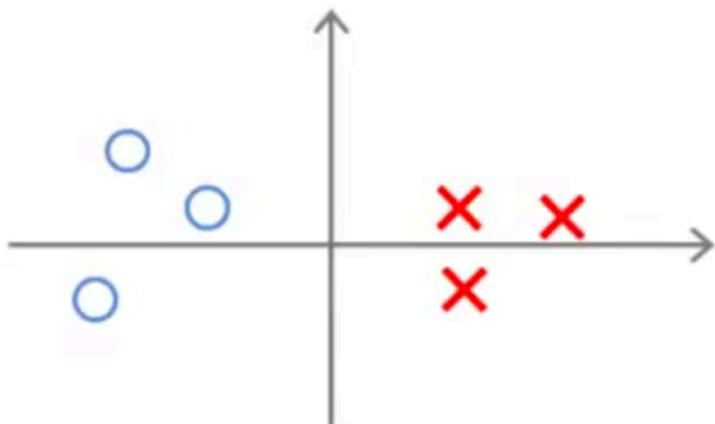
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 =$$

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

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

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

Simplification: $\theta_0 = 0$



SVM Decision Boundary

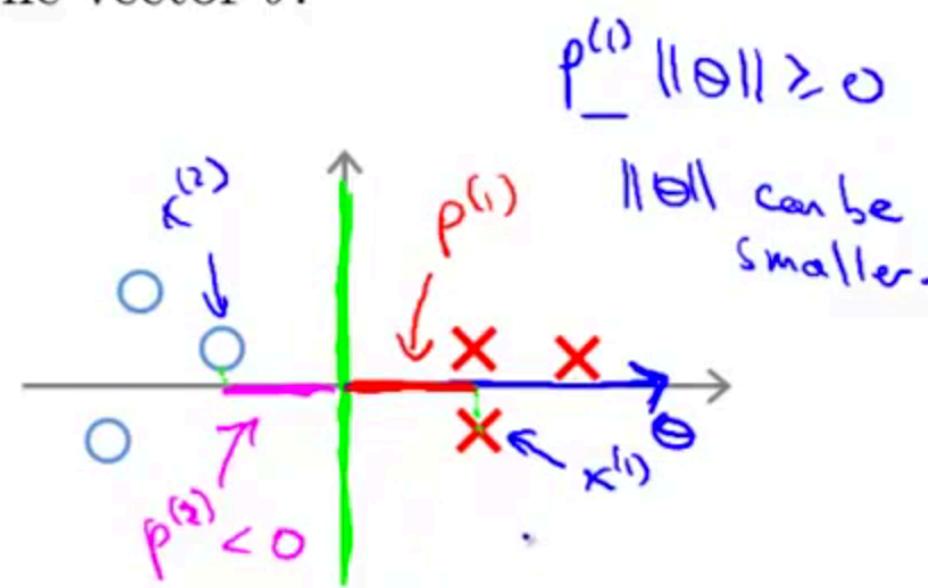
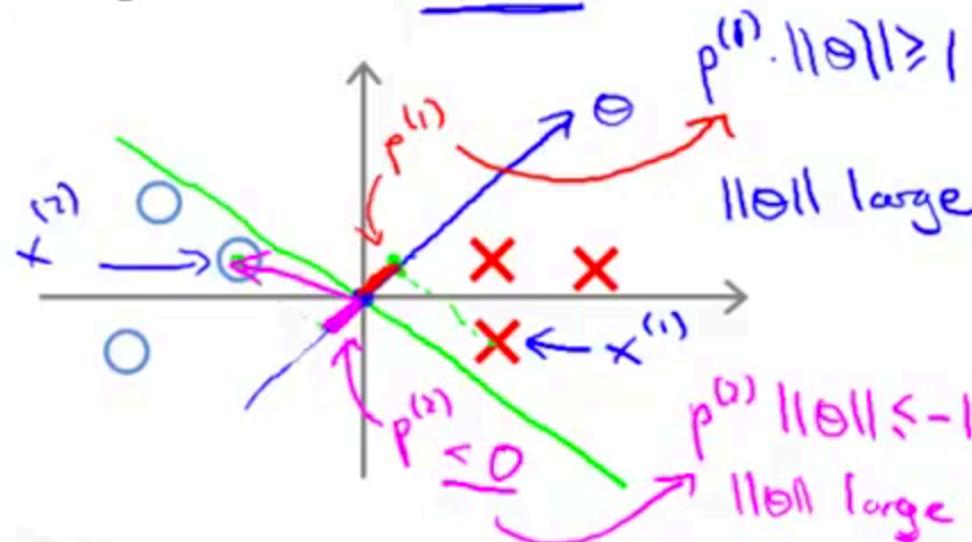
$$\rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

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

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

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

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



SVM Decision Boundary

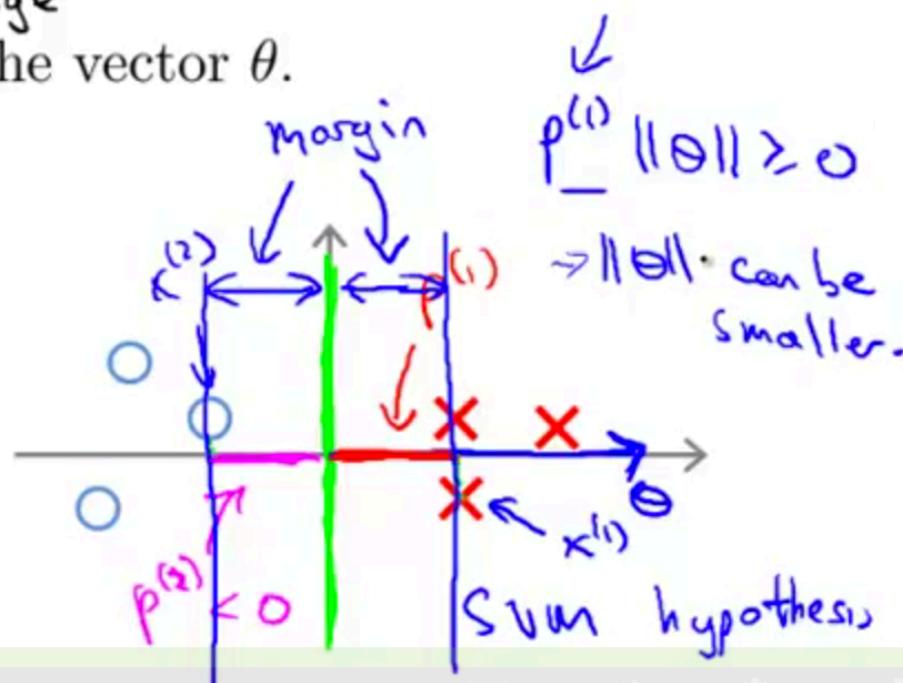
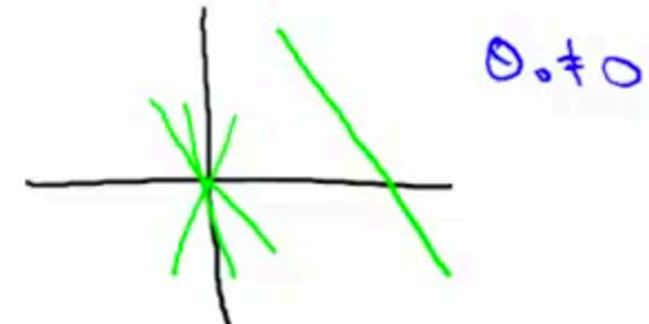
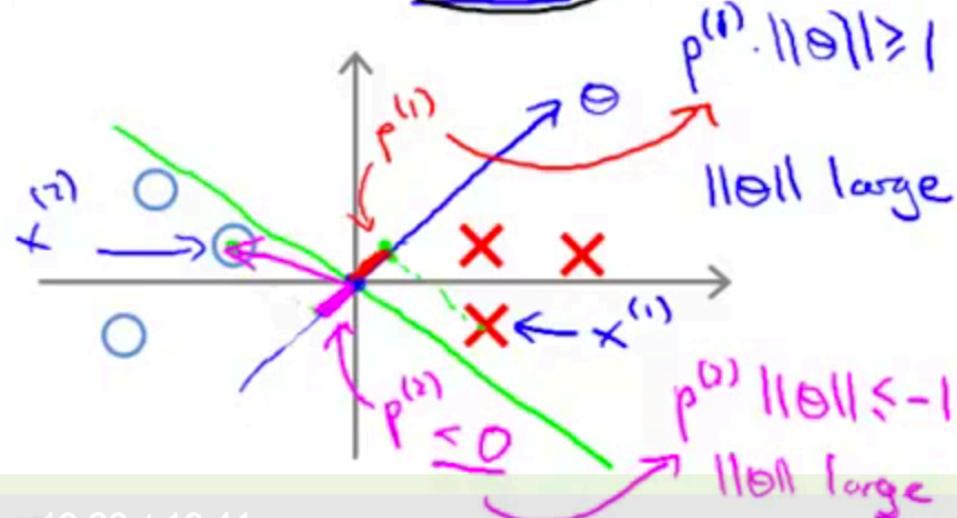
$$\rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

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

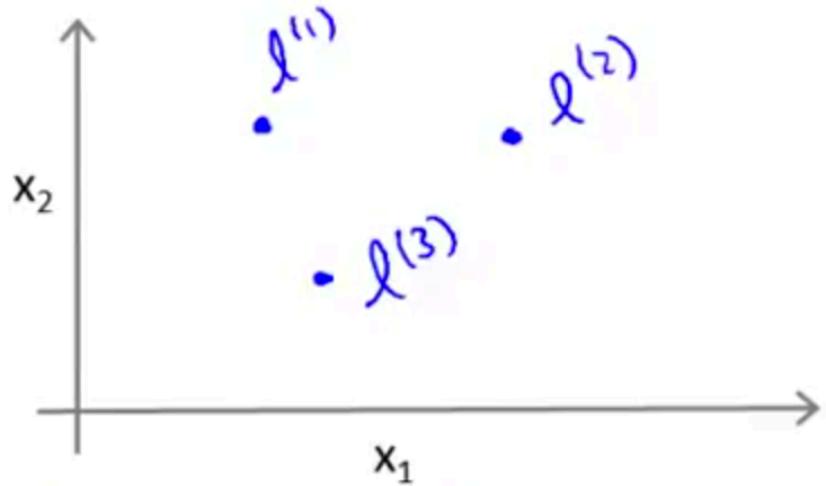
$$p^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = -1$$

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

Simplification: $\underline{\theta_0 = 0} \leftarrow$



Kernel



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

Given x :

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

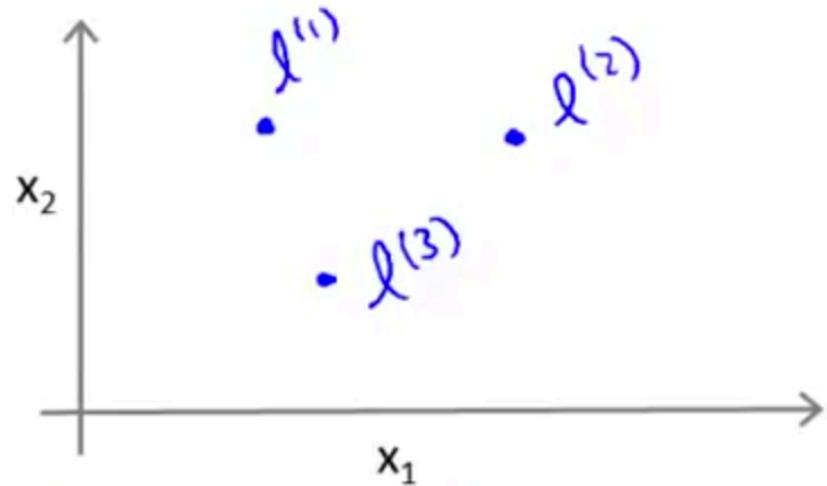
$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

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

⋮

$$\begin{aligned} & \|w\| \\ & \nwarrow \\ & \frac{\|x - l^{(1)}\|^2}{2\sigma^2} \end{aligned}$$

Kernel



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

Given x :

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

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

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

↑ kernel (Gaussian kernels)

$$\|w\|$$

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

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

$$\sqrt{\|x - l^{(3)}\|^2}$$

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

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

$$\begin{array}{lcl} l^{(1)} & \rightarrow & f_1 \\ l^{(2)} & \rightarrow & f_2 \\ l^{(3)} & \rightarrow & f_3 \\ \uparrow & & \uparrow \\ x & & x \end{array}$$

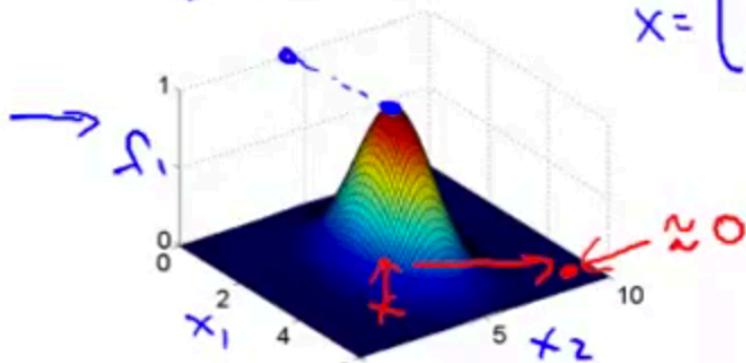
If x if far from $l^{(1)}$:

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

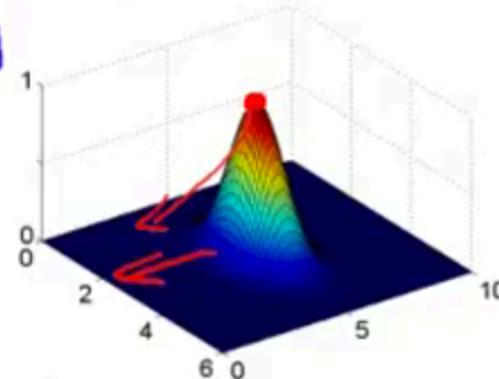
Example:

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

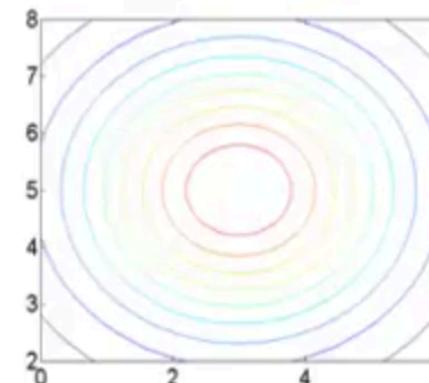
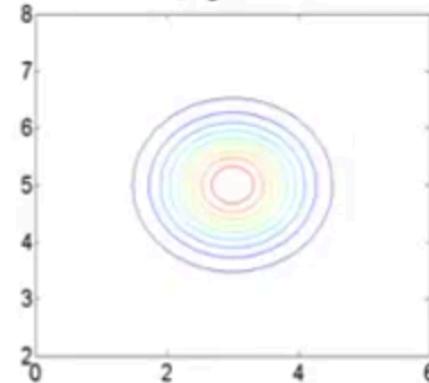
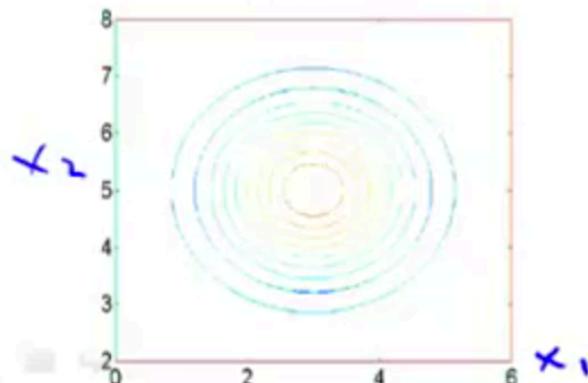
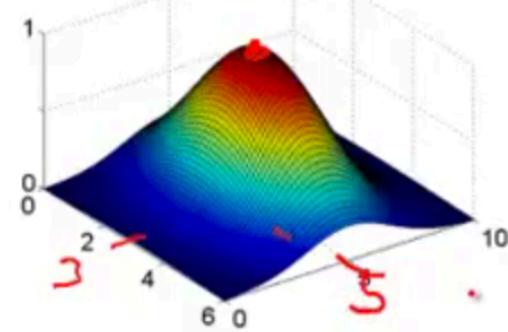
$$\rightarrow \sigma^2 = 1$$



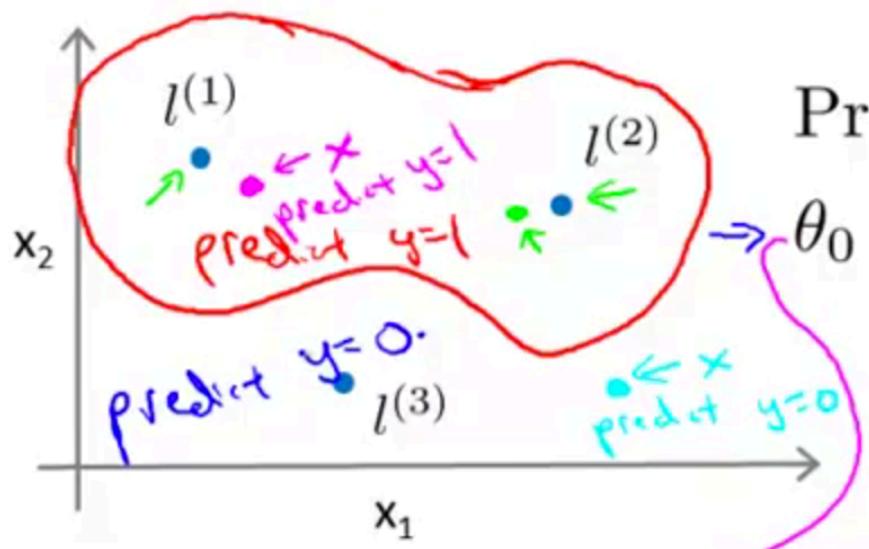
$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \quad \sigma^2 = 0.5$$



$$\sigma^2 = 3$$



Andrew Ng



Predict "1" when

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



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

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

$$\begin{aligned} & \rightarrow \theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0 \\ & = -0.5 + 1 = 0.5 \geq 0 \end{aligned}$$

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

$$\rightarrow \underline{\theta_0 + \theta_1 f_1 + \dots \approx -0.5 < 0}$$

SVM with Kernels

- Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
- choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ \dots & \end{aligned}$$

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

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

$$\begin{aligned} x^{(i)} \rightarrow f_1^{(i)} &= \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots & \\ f_m^{(i)} &= \text{similarity}(x^{(i)}, l^{(m)}) \end{aligned}$$

$$\begin{aligned} x^{(i)} \in \mathbb{R}^{n+1} \quad (\text{or } \mathbb{R}^n) \\ \rightarrow f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \\ f_0^{(i)} = 1 \end{aligned}$$

Andrew Ng

SVM with Kernels

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

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

Training:

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

$n = m$

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

$\rightarrow \theta_0$

$$\begin{bmatrix} - & \sum_j \theta_j^2 \\ - & \end{bmatrix} = \theta^T \theta \quad \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignoring } \theta_0)$$

$$\rightarrow \underline{\theta^T M \theta} \quad \leftarrow \| \theta \|^2$$

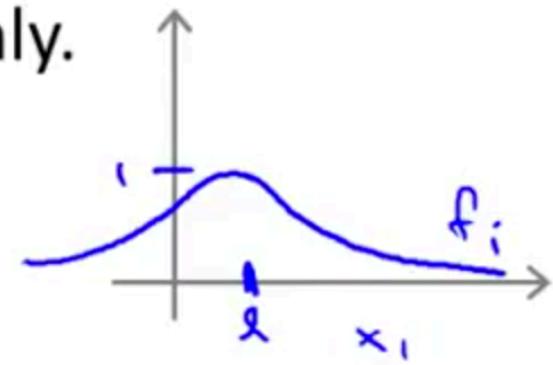
SVM parameters:

$C \left(= \frac{1}{\lambda} \right)$. \rightarrow Large C: Lower bias, high variance. (small λ)
 \rightarrow Small C: Higher bias, low variance. (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly.

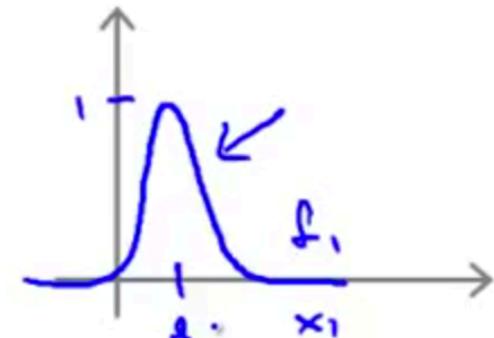
\rightarrow Higher bias, lower variance.

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



Small σ^2 : Features f_i vary less smoothly.

Lower bias, higher variance.



Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

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

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

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

→ n large, m small

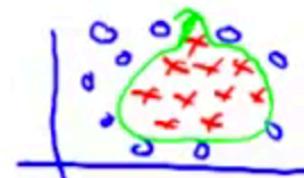
$$x \in \mathbb{R}^{n+1}$$

Gaussian kernel:

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

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

and/or m large



Need to choose $\frac{\sigma^2}{\tau}$.

Kernel (similarity) functions:

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

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

```
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.

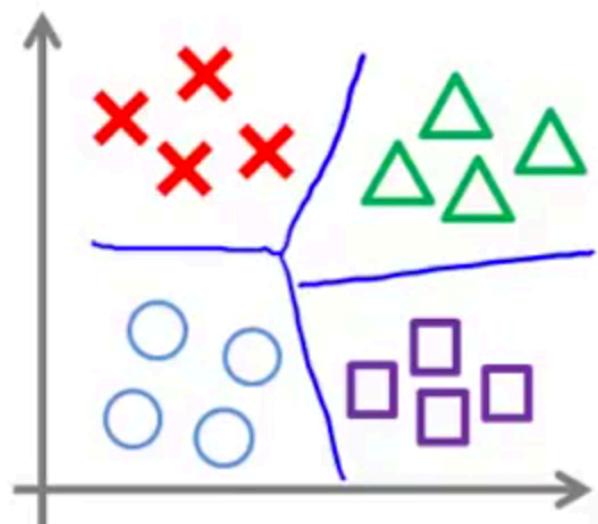
→ (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, l) = \underbrace{(x^T l)^2}_{\text{degree } 2} + \underbrace{\dots}_{\text{degree } 3}, \underbrace{(x^T l + 1)^3}_{\text{degree } 3}, \underbrace{(x^T l + 5)^4}_{\text{degree } 4}$

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

Multi-class classification



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

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 $i = 1, 2, \dots, K$), get $\theta^{(1)}, \theta^{(2)}, \dots, \underline{\theta^{(K)}}$
Pick class i with largest $(\underline{\theta^{(i)}})^T x$

$$\begin{matrix} \uparrow \\ y=1 \end{matrix} \quad \begin{matrix} \nwarrow \\ y=2 \end{matrix} \quad \cdots \quad \begin{matrix} \nwarrow \\ \Theta = K \end{matrix}$$

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. $n \geq m$, $n = \underline{10,000}$, $m = \underline{10} \dots \underline{1000}$)

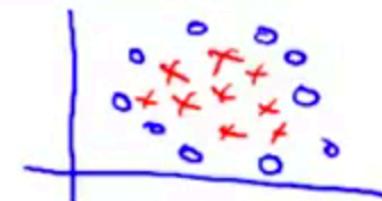
→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If n is small, m is intermediate: ($n = 1 - 1000$, $m = 10 - \underline{10,000}$)

→ Use SVM with Gaussian kernel

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

→ Create/add more features, then use logistic regression or SVM without a 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. $n \geq m$, $n = \underline{10,000}$, $m = \underline{10} \dots \underline{1000}$)
- Use logistic regression, or SVM without a kernel ("linear kernel")
- If n is small, m is intermediate: ($n = \underline{1-1000}$, $m = \underline{10 - 10,000}$)
 - Use SVM with Gaussian kernel
- If n is small, m is large: ($n = \underline{1-1000}$, $m = \underline{50,000+}$)
 - 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.

