

With supervised learning, the performance of many supervised learning algorithms will be pretty similar, and what matters less often will be whether you use learning algorithm a or learning algorithm b, but what matters more will often be things like the amount of data you create these algorithms on, as well as your skill in applying these algorithms. Things like your choice of the features you design to give to the learning algorithms, and how you choose the regularization parameter, and things like that. But, there's one more algorithm that is very powerful and is very widely used both within industry and academia, and that's called the support vector machine.

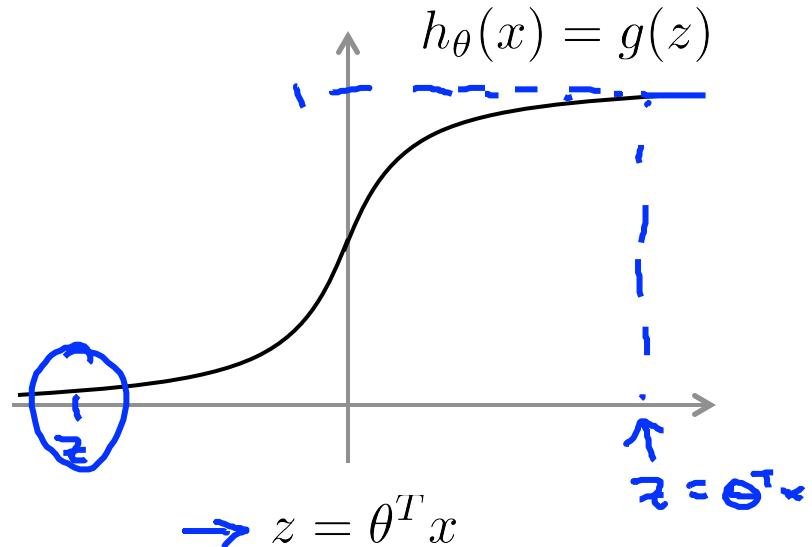
And compared to both logistic regression and neural networks, the Support Vector Machine, or SVM sometimes gives a cleaner, and sometimes more powerful way of learning complex non-linear functions.

Support Vector Machines

Optimization objective

Alternative view of logistic regression

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



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

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

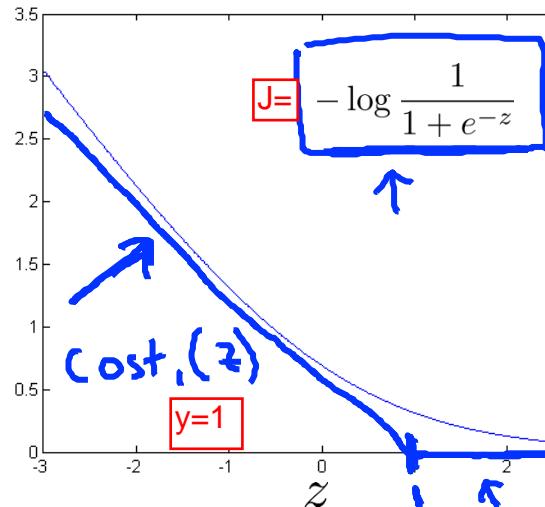
Alternative view of logistic regression

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

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

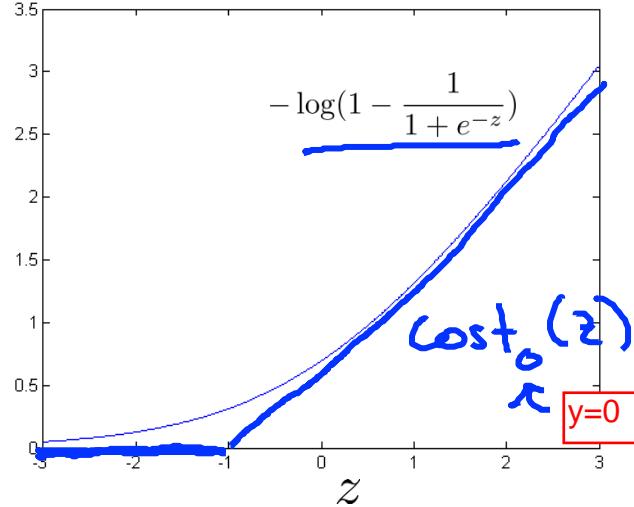
If $y = 1$ (want $\theta^T x \gg 0$):

$$z = \theta^T x$$



when z is large, that corresponds to a small value of cost function. And this explains why, when logistic regression sees a positive example, with $y=1$, it tries to set theta transport x to be very large because that corresponds to this term, in the cost function, being small.

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



Support vector machine

Logistic regression:

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

Support vector machine:

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

Support vector machine

Logistic regression:

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


Support vector machine:

$$\min_{\theta} \underbrace{C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})}_{A} + \frac{1}{2} \sum_{j=0}^n \theta_j^2 \quad B$$

$$\min_u \underbrace{(u-5)^2 + 10}_{10} \rightarrow u=5$$

$$\min_u \underbrace{10(u-5)^2 + 10}_{10} \rightarrow u=5$$

$$A + \frac{\lambda}{2} B \leftarrow C = \frac{1}{\lambda}$$

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

SVM hypothesis

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

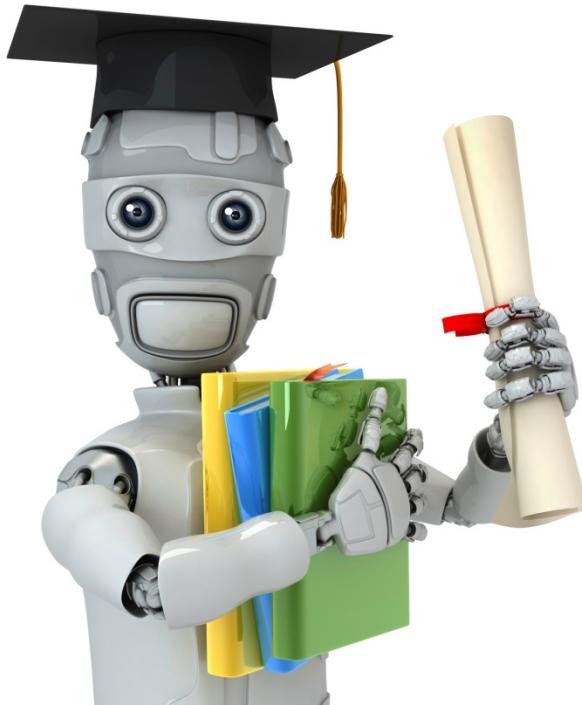
Hypothesis:

SVM hypothesis

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

Hypothesis:

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



Machine Learning

Support Vector Machines

Large Margin Intuition

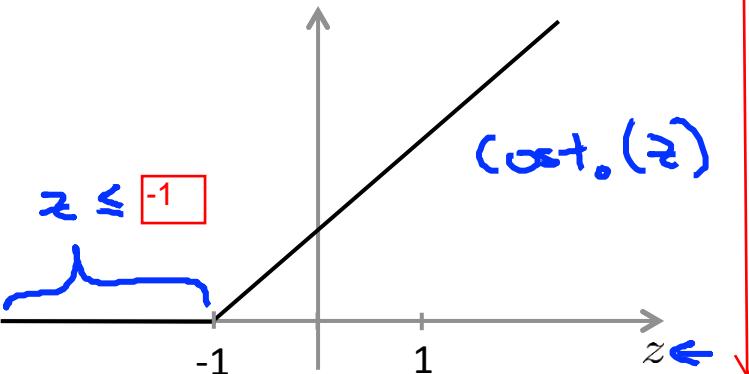
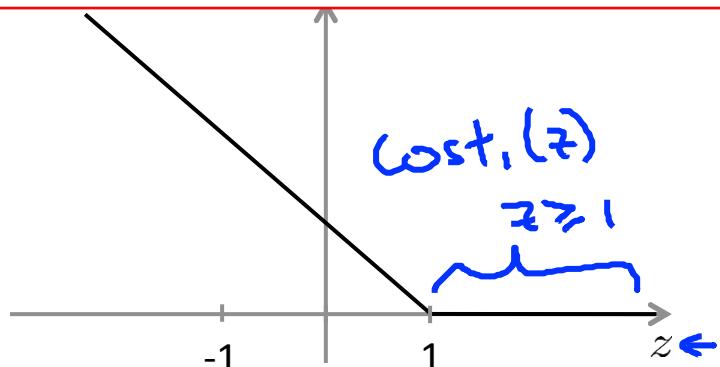
Sometimes people talk about support vector machines, as large margin classifiers

Support Vector Machine

all we need to get the classification right is ≥ 0 , but svm wants more than that (use 1) ! so this builds in an extra safety factor or safety margin factor into SVM, logistic regression can do it too

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

what it takes to make the cost function small ?



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

$$\begin{aligned} \theta^T x &\geq 1 \\ \theta^T x &\leq -1 \end{aligned}$$

$$C = 100,000$$

C is very large,
what SVM will do ?

SVM Decision Boundary

If C is very large, then when minimizing this optimization objective, we're going to be highly motivated to choose a value, so that this first term = 0

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

if C is very large

what does it take to make this first term = 0 (when C huge)

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

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

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

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

now when we choose parameters and sure that this first term = 0, then what left is

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

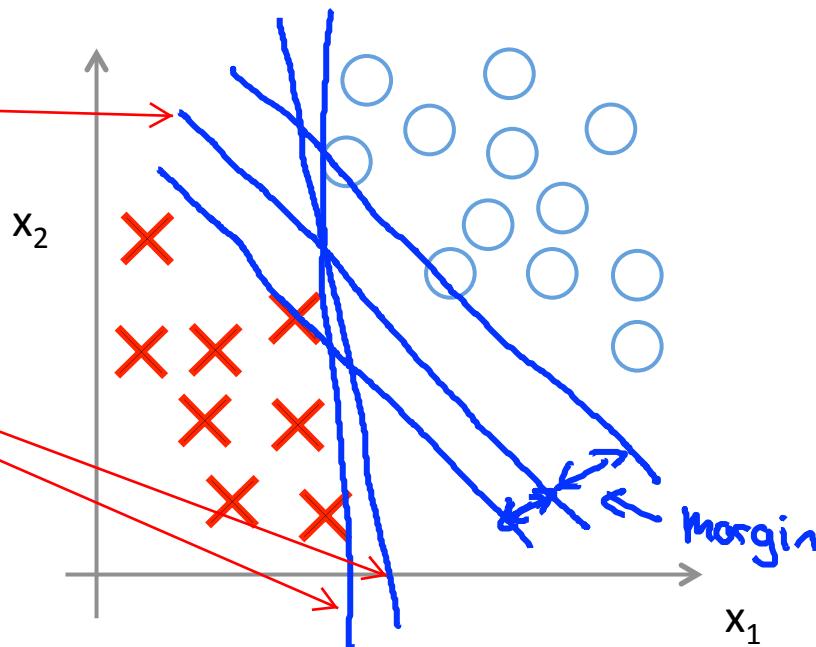
$$\text{s.t. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

subject to the constrain that

linearly separable : there exists a straight line, although there is many a different straight lines, they can separate the positive and negative examples perfectly.

SVM Decision Boundary: Linearly separable case

The Support Vector Machines will instead choose this decision boundary, and that seems like a much better decision boundary than others



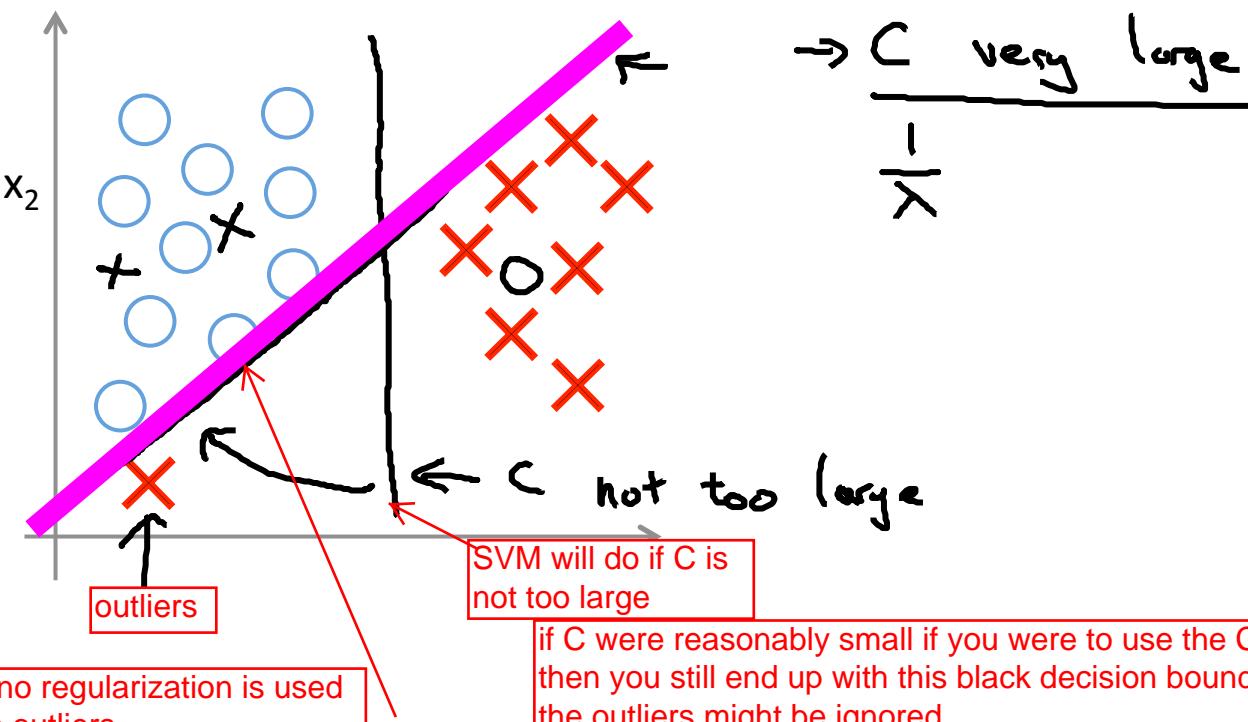
Large margin classifier

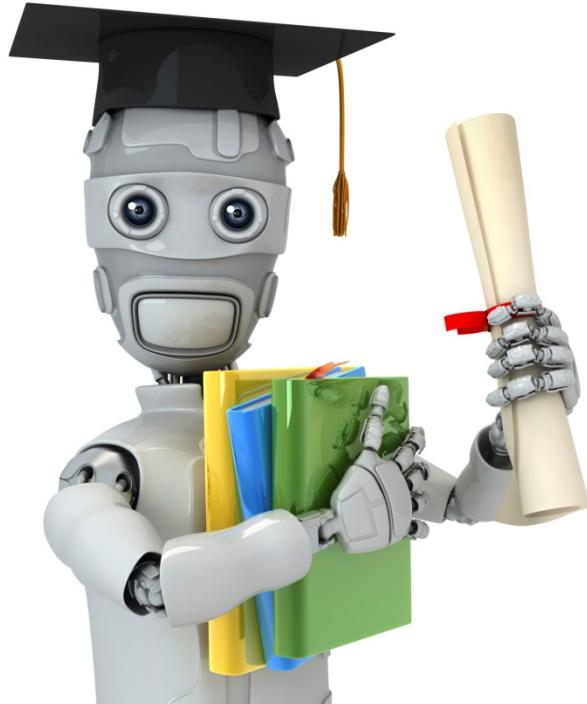
this distance is called "margin" of the SVM and this gives the SVM a certain robustness, because it tries to separate the data with a large a margin, so SVM is also called a "large margin classifier" and this is a consequence of the optimization problem we wrote down on the previous slide.

C parameter is a positive value that controls the penalty for misclassified training examples. A large C tells the SVM to try to classify all the examples correctly.

Large margin classifier in presence of outliers

use large margin classifier it can be sensitive to outliers





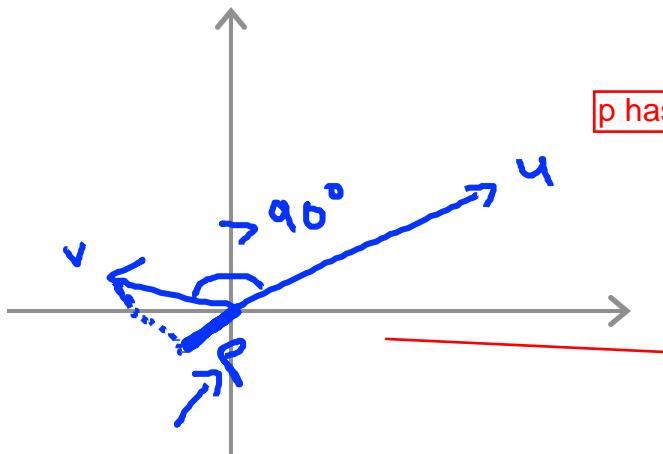
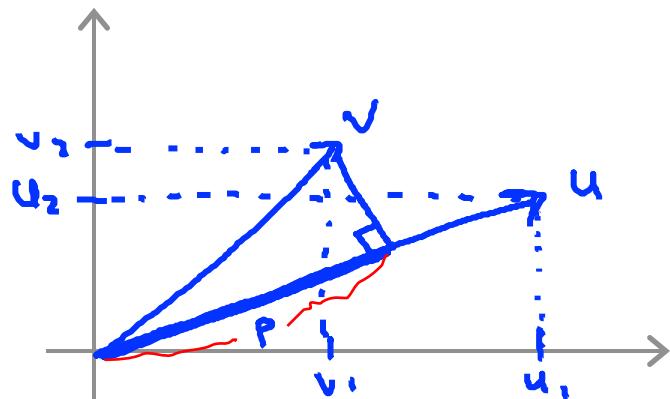
Machine Learning

Support Vector Machines

The mathematics
behind large margin
classification (optional)

why SVM is a large
margin classifier :

Vector Inner Product



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

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

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

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

p = length of projection of v onto u .

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

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

Signed

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

$$p < 0$$

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

svm in the optimization objective is minimizing the squared norm of the square length of the parameter vector theta.

SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_0^2 + \theta_1^2 + \theta_2^2) = \frac{1}{2} (\boxed{\theta_0^2 + \theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

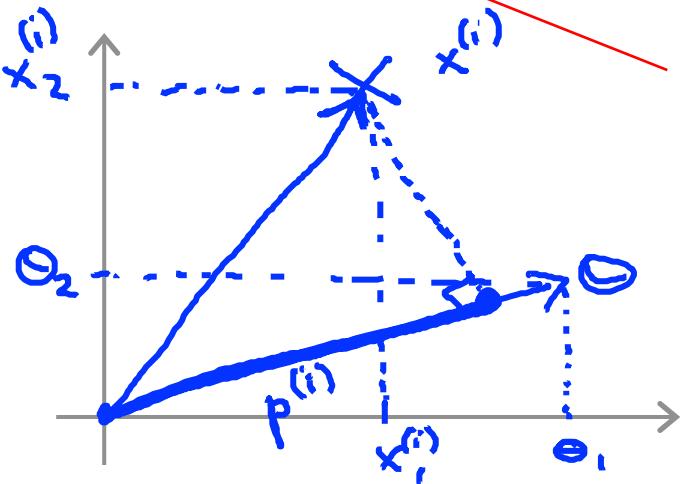
s.t. $\boxed{\theta^T x^{(i)} \geq 1}$ if $y^{(i)} = 1$
 $\rightarrow \theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

Simplification: $\underline{\theta_0 = 0}$. $n=2$

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

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

\uparrow
 $u^T v$



$$\begin{aligned} \underline{\theta^T x^{(i)}} &= \boxed{p \cdot \|\theta\|} \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

x 在 theta 上的投影 * |theta|

by making the margin large (by these P1, P2, P3 large) the SVM can end up with a smaller value for the norm of theta which is what it is trying to do in the objective.

SVM Decision Boundary

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

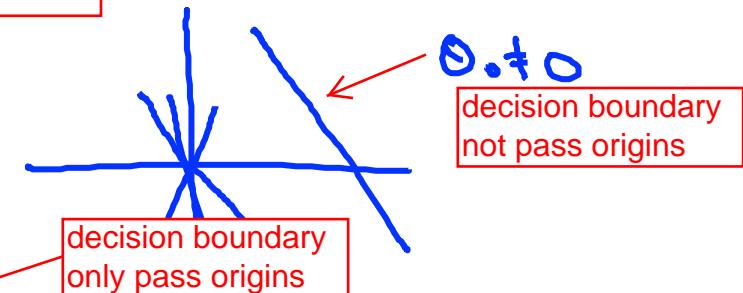
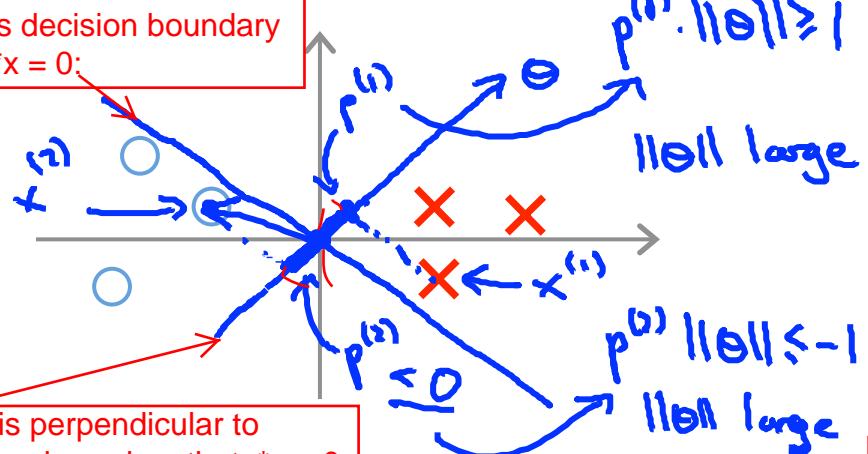
$$\text{s.t. } p^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1$$

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

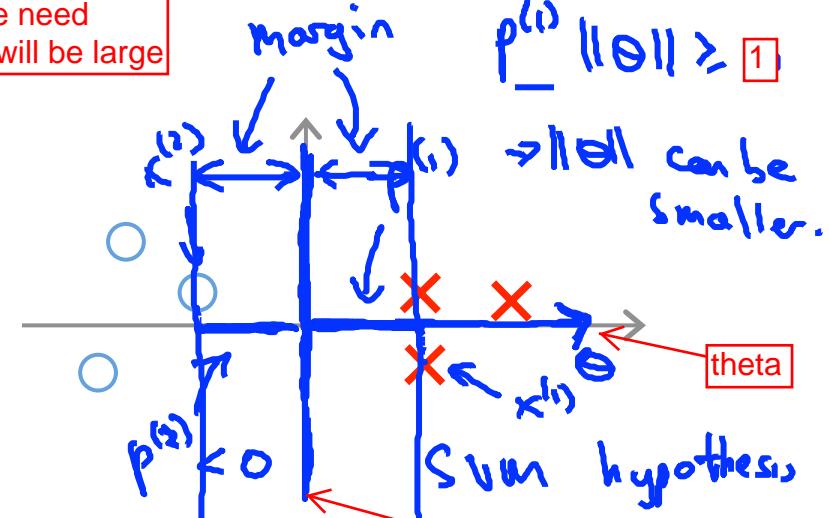
where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$

why svm wont choose this line as decision boundary
 $\theta^*x = 0$:



$p^{(i)} \cdot \|\theta\| \geq 1$, then θ will be smaller



why svm choose this line as decision boundary $\theta^*x = 0$:

SVM Decision Boundary

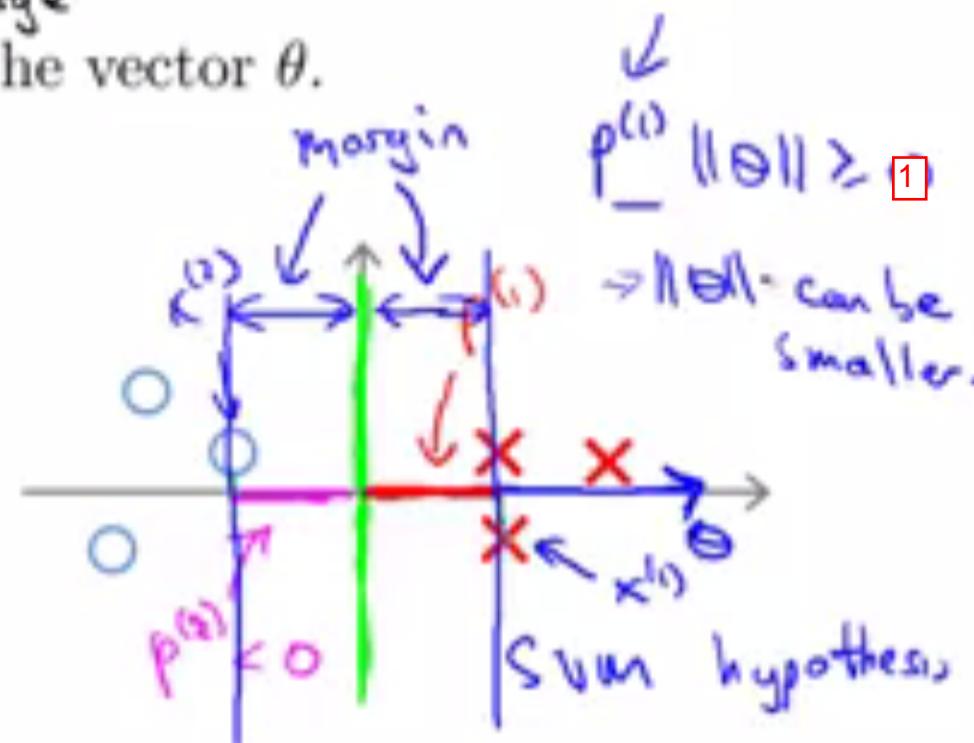
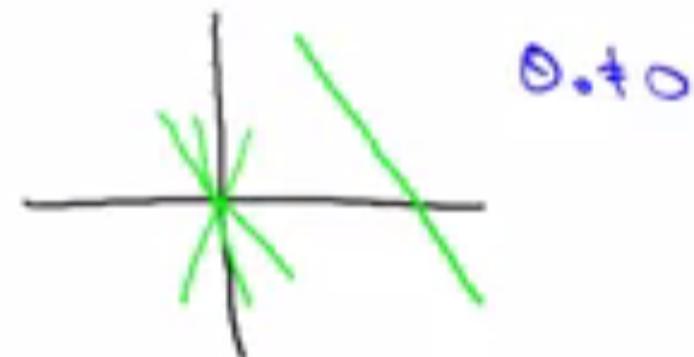
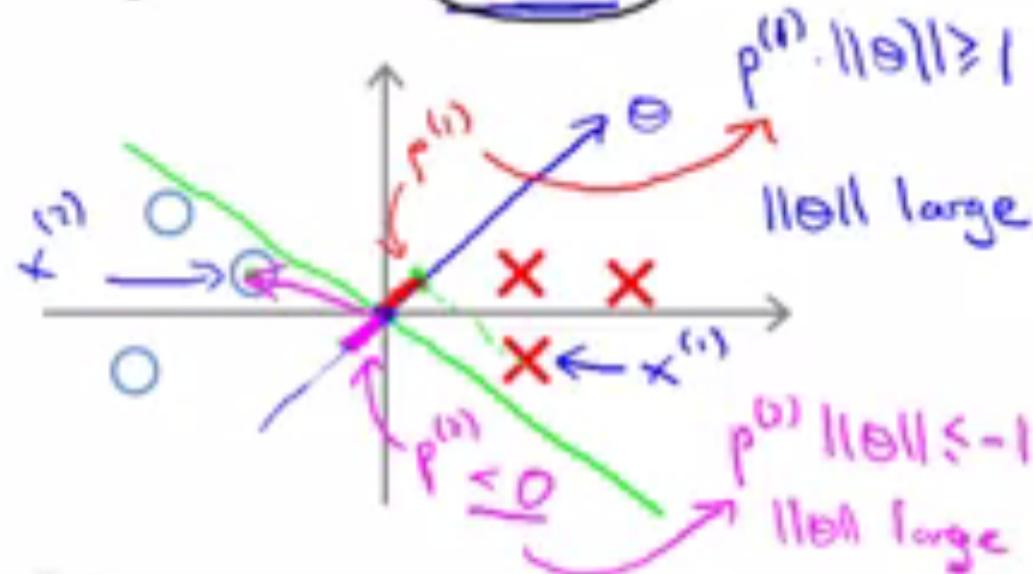
$$\rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

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 .

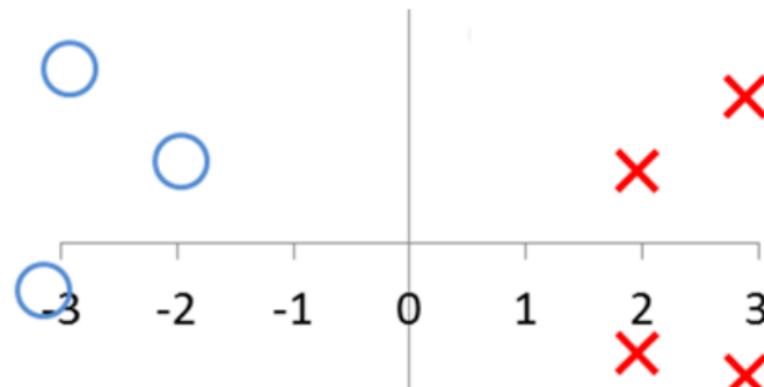
Simplification: $\theta_0 = 0 \leftarrow$



The SVM optimization problem we used is:

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

$$\text{s.t. } \|\theta\| \cdot p^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$
$$\|\theta\| \cdot p^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

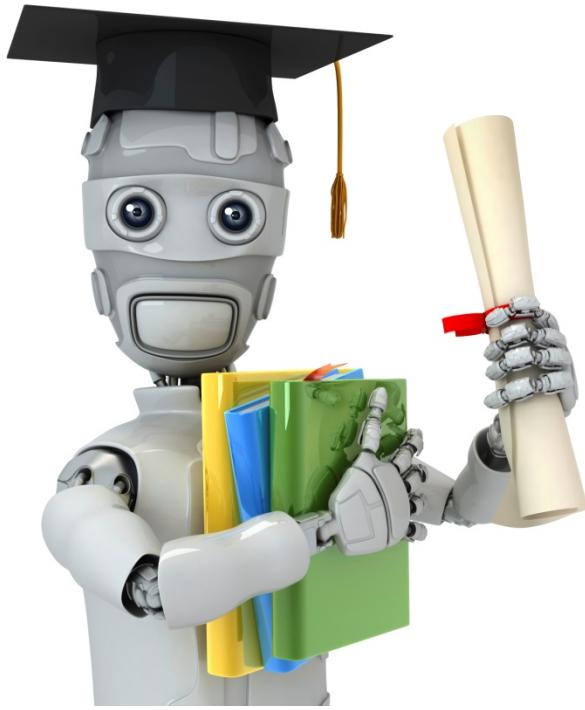


where $p^{(i)}$ is the (signed - positive or negative) projection of $x^{(i)}$ onto θ . Consider the training set above. At the optimal value of θ , what is $\|\theta\|$?

1/4

1/2

Correct

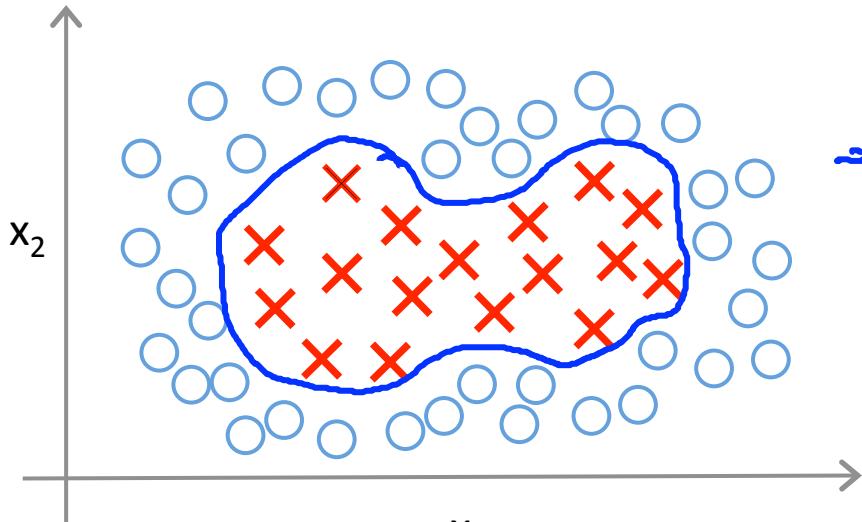


Machine Learning

Support Vector Machines

Kernels I

Non-linear Decision Boundary



you can do complex polynomial features :

Predict $y = 1$ if

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

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

another way to write it :

$$\Rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

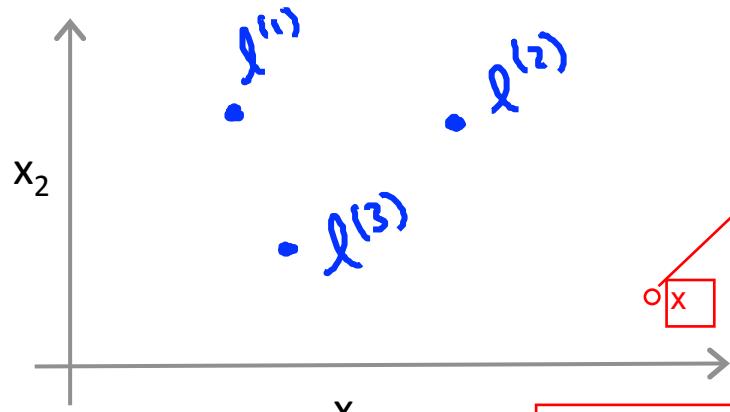
$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \dots$$

is there a different/better choice of features than this high order polynomials because it's not clear that this high order polynomial is what we want, and for computer vision the input is an image with lots of pixels, we saw how using high order polynomials becomes very computationally expensive because there are a lot of these higher order polynomial terms.

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

its one idea to find out features f_1, f_2, f_3 :

Kernel



given an example x

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

Given x :

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

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

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

$$\begin{aligned} & \frac{\|w\|}{\sim} \\ & \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) \\ & \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right) \\ & \exp(\dots) \end{aligned}$$

Kernel (Gaussian kernels)

$$k(x, l^{(1)})$$

one example of similarity function is Gaussian kernel

Kernels and Similarity

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

if x close to l_1

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

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

if x far from l_1

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

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

given the training example X , we can now compute three new features: f_1, f_2 , and f_3 , given three landmarks

$$\begin{aligned} l^{(1)} &\rightarrow f_1 \\ l^{(2)} &\rightarrow f_2 \\ l^{(3)} &\rightarrow f_3. \end{aligned}$$

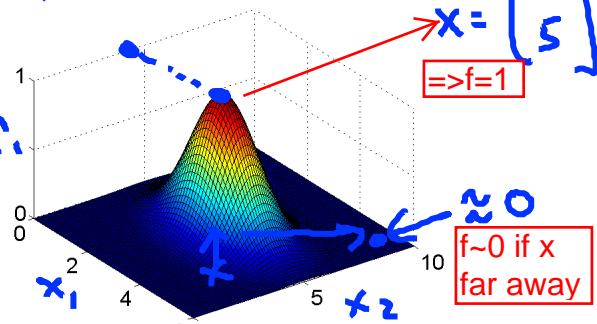
↑
↑
 X

this is a feature, f_1 measures how close X is to the first landmark and it varies between 0 and 1 one depending on how close X is to the first landmark l_1 .

Example:

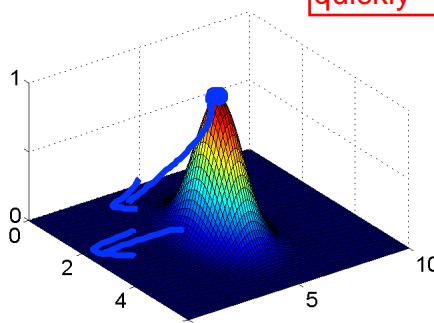
$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

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



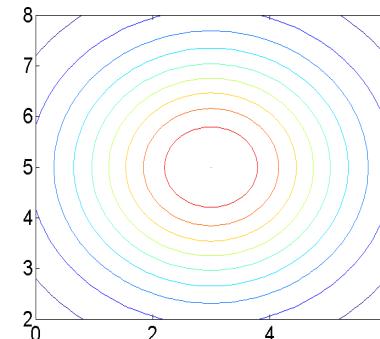
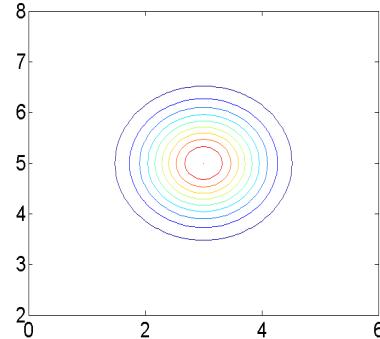
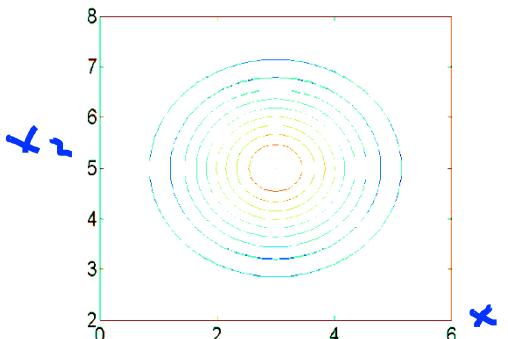
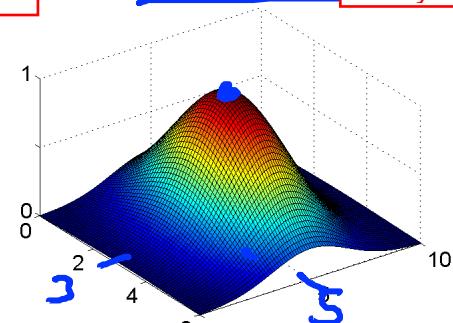
$$\sigma^2 = 0.5$$

smaller => f1 falls quickly

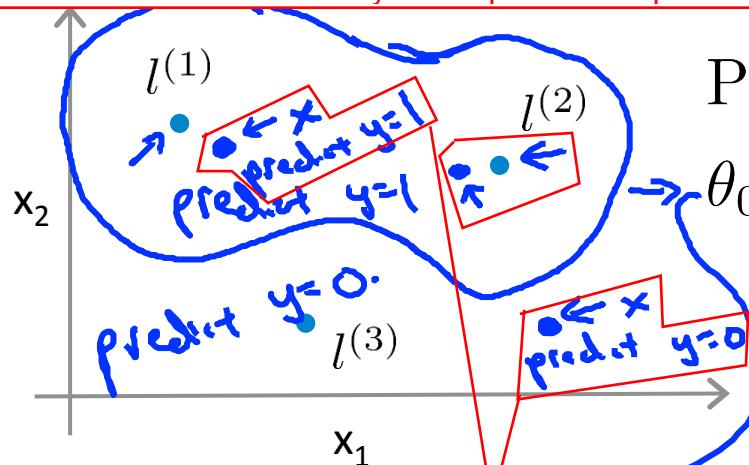


$$\sigma^2 = 3$$

larger => f1 falls slowly

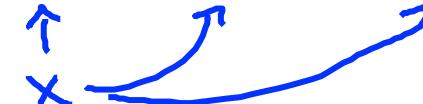


for points near I_1 and I_2 we end up predicting positive. And for points far away from I_1 and I_2 , we end up predicting that the class is equal to 0. so What we end up doing is that the decision boundary of this hypothesis would end up looking something like this where inside this red decision boundary would predict Y equals 1 and outside we predict Y equals 0.



Predict "1" when

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



assume we already train the model and get theta

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

if i have x here => $f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$.

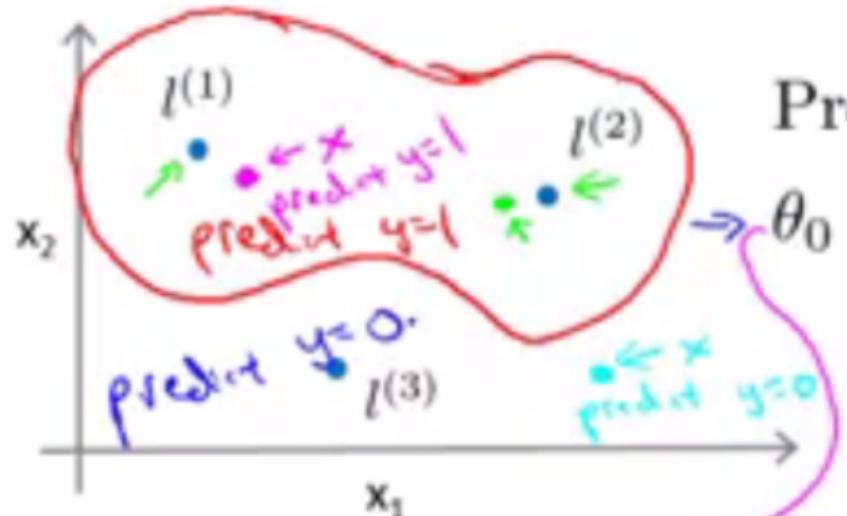
$$\begin{aligned} & \theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0 \\ &= -0.5 + 1 = 0.5 \geq 0 \end{aligned}$$

=> predict $y=1$

if i have x here => $f_1, f_2, f_3 \approx 0$

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

=> predict $y=0$



Predict "1" when

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



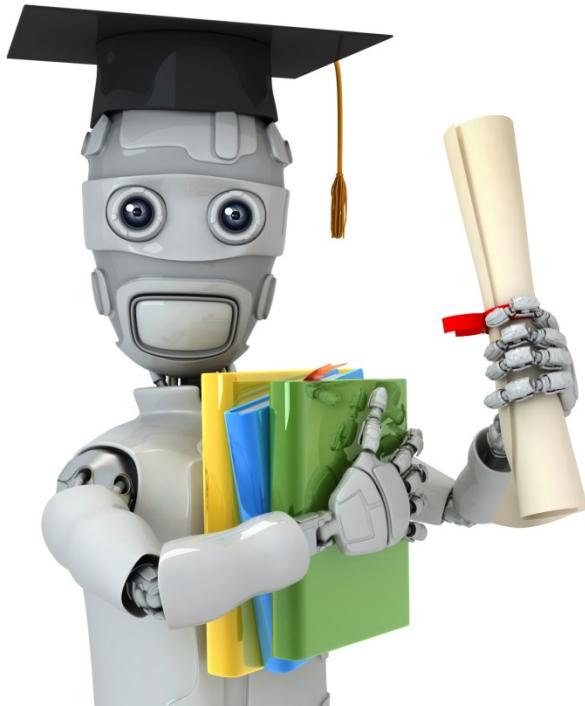
$$\underline{\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0}$$

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

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

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

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

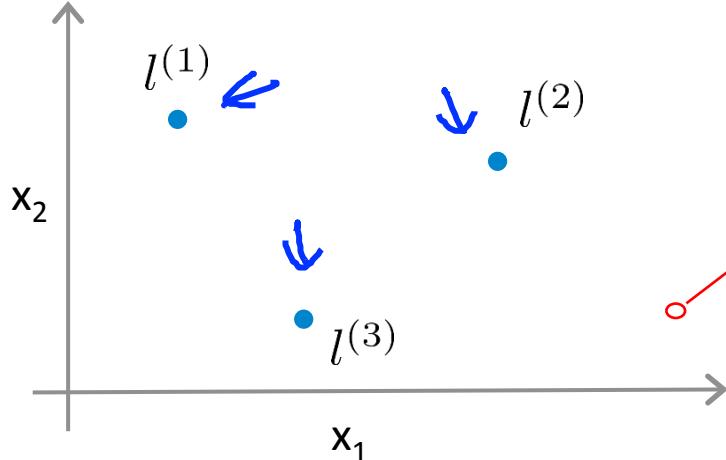


Machine Learning

Support Vector Machines

Kernels II

Choosing the landmarks



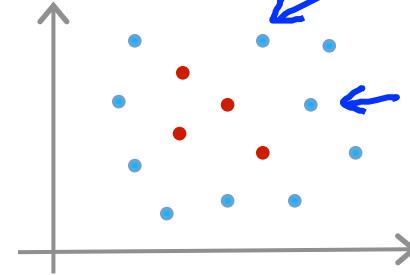
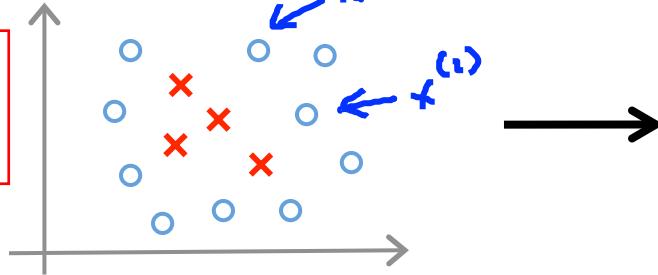
Given x :

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

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

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

We're just going to put landmarks as exactly the same locations as the training examples.



$l^{(1)}$
 $l^{(2)}$
⋮
 $l^{(m)}$

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 \underline{x} :

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

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

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

$$\begin{aligned} x^{(i)} &\rightarrow \boxed{\begin{array}{c} f_1^{(i)} = \overline{\text{sim}}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \overline{\text{sim}}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \overline{\text{sim}}(x^{(i)}, l^{(m)}) \end{array}} \\ \text{第 } i \text{ 個點與其他點有多近? (0~1表示)} \end{aligned}$$

$$\begin{aligned} x^{(i)} &\in \mathbb{R}^{n+1} \quad (\text{or } \mathbb{R}^n) \\ \overline{x}^{(i)} &= \boxed{\begin{array}{c} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{array}} \\ f_0^{(i)} &= 1 \end{aligned}$$

If you already have a learning set of theta, then if you given a value of x and you want to make a prediction:

SVM with Kernels

m training examples
(m landmarks)

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

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

$$\theta \in \mathbb{R}^{m+1}$$

$$\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

Training:

how do
you get
theta: sole
this min
problem

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

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

$n = m$
 $\cancel{\theta_0} = m$

in svm this usually
done differently

don't do regularization of
theta 0

$$\begin{aligned} - \sum_j \theta_j^2 &= \theta^T \theta \quad \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \\ - &\quad \leftarrow \|\theta\|^2 \\ - &\quad \leftarrow \theta^T M \theta \end{aligned}$$

(ignoring θ_0)
 $M = 10,000$

support vector machines and kernels tend go particularly well together. Whereas, logistic regression and kernels, you know, you can do it, but this would run very slowly.

SVM parameters:

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

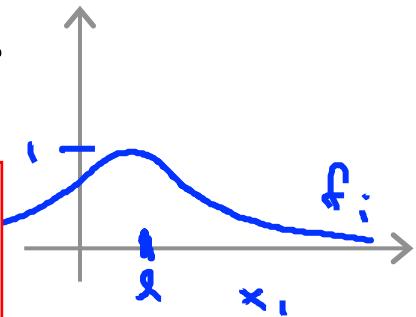
(small λ) not using regularization
(large λ) using regularization

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

cause Higher bias, lower variance.

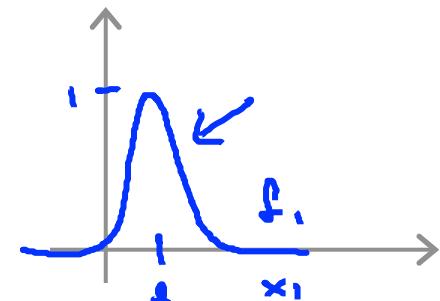
$$\exp \left(- \frac{\|x - \mu\|^2}{2\sigma^2} \right)$$

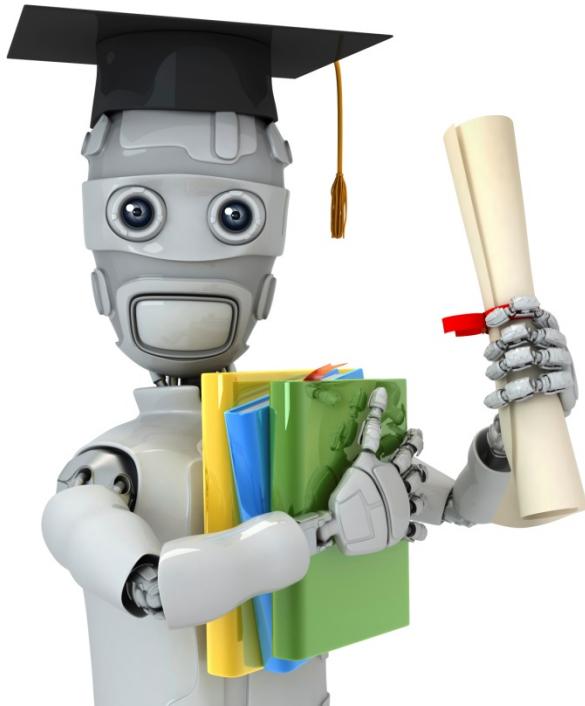
you tend to get a hypothesis that varies slowly, or varies smoothly as you change the input $x \Rightarrow$ not too sensitive to noise



Small σ^2 : Features f_i vary less smoothly.

cause Lower bias, higher variance.





Machine Learning

Support Vector Machines

Using an SVM

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):

standard linear classifier

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

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

you want to fit a linear decision boundary and not try to fit a very complicated nonlinear function, because might not have enough data. And you might risk overfitting, if you're trying to fit a very complicated function

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

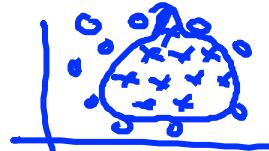
→ n large, m small $x \in \mathbb{R}^{n+1}$

→ Gaussian kernel:

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

Need to choose $\underline{\sigma^2}$

$x \in \mathbb{R}^n$, n small
and/or n large



Kernel (similarity) functions:

function f = kernel(x1,x2)

$$f = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right)$$

return

第 i 個點與其他點有多近？(0~1 表示)

$$x \rightarrow f_1 f_2 \dots e_n$$

→ Note: Do perform feature scaling before using the Gaussian kernel.

$$\Rightarrow \boxed{\|x - l\|^2} \quad V = x - l \quad x \in \mathbb{R}^{n+1}$$

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

$$= \underbrace{(x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2}_{\text{1000 feet}^2} \quad 1-5 \text{ bedrooms}$$

those distances will be almost essentially dominated by the sizes of the houses and the number of bathrooms would be largely ignored. so to avoid this in order to make a machine work well, do perform future scaling.

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:

Usually it is used only for data where X and l are all strictly non negative, and so that ensures that these inner products are never negative.

$$k(x, l) = \underbrace{(x^T l)}_3^1 + \underbrace{(x^T l)^2}_2 + \underbrace{(x^T l + 1)}_3 + \underbrace{(x^T l + 5)}_4$$

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

used if your input data is text strings or other types of strings.

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

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

There are sort of more esoteric kernels that you can use to measure similarity between different objects.

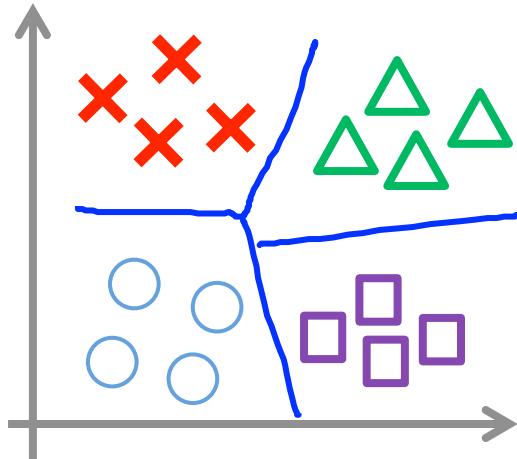
Suppose you are trying to decide among a few different choices of kernel and are also choosing parameters such as C , σ^2 , etc. How should you make the choice?

- Choose whatever performs best on the training data.
- Choose whatever performs best on the cross-validation data.

Correct

- Choose whatever performs best on the test data.
- Choose whatever gives the largest SVM margin.

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

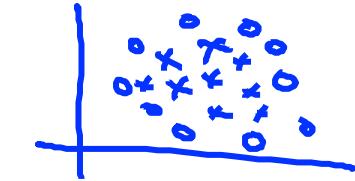
$\underline{y=1}$ $\underline{y=2}$... class $y=K$

if you have many features with smaller training sets, a linear function will probably do fine, and you don't have really enough data to fit a very complicated nonlinear function.

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.



logistic regression and SVM without the kernel, are similar algorithms and they will usually do pretty similar things and give similar performance, but depending on your implementational details, one may be more efficient than the other

optimization problem that the SVM has is a convex optimization problem and so the good SVM optimization software packages will always find the global minimum or something close to it.