

Machine Learning

Week 7

Support Vector Machines

This week is about Support Vector Machines. SVMs are considered by many to be the most powerful 'black box' learning algorithm, and by posing a cleverly-chosen optimization objective, one of the most widely used learning algorithms today.

Large Margin Classification

Optimization Objective

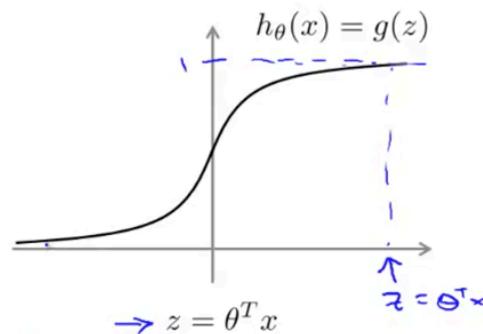
So far, for supervised learning algorithms, we have covered linear regression, logistic regression, and neural networks.

SVMs are also supervised learning algorithms.

SVM sometimes gives us a cleaner more powerful way of learning complex non-linear functions.

Before we start to understand SVM, lets first revisit logistic regression:

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



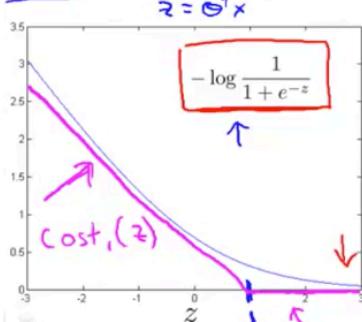
If $y = 1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$
If $y = 0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

If we look at the cost function of logistic regression, each training example (x, y) contributes to our overall cost function.

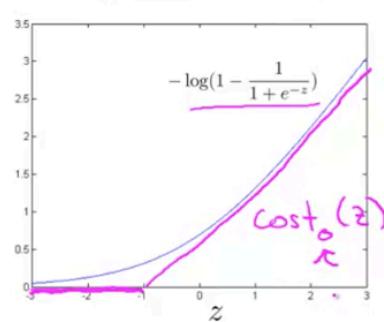
Cost of example: $-(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x))) \leftarrow$

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

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



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



What we do for support vector machines is change the logistic regression cost function curve to be two straight line segments as shown in the magenta color.

- "cost_1_(z)"
- "cost_0_(z)"

This is to give us computational advantages that will allow for easier optimization later on.

Formally, our new cost function for a support vector machine is:

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:

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

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

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

We cross out our $1/m$ term because multiplying the cost function by some constant does not change the optimization problem and as such we don't care about it / don't need it.

Another convention that we change for the SVM is changing from $(A + \lambda * B)$ to $(C * A + B)$. This is just another way to control the way we optimize the first term instead of optimizing the second term.

Think of $(C = 1/\lambda)$ as it plays a similar role.

As such we have our overall cost function for a support vector machine which we want to minimize:

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

Lastly, the support vector machine does not output a probability like logistic regression. What a support vector machine does is make a direct prediction of 1 or 0.

Hypothesis:

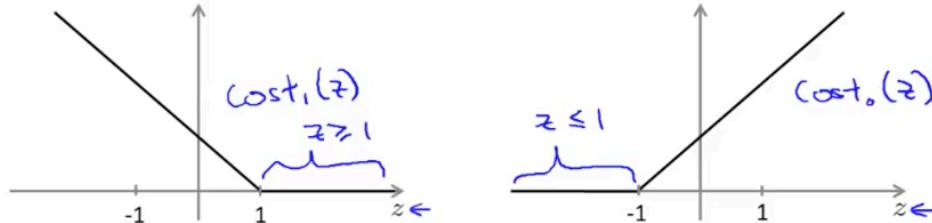
$$h_{\theta}(x) \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Large Margin Intuition

Sometimes the term 'Support Vector Machine' is called 'Large Margin Classifiers'.

Once again, here is our cost function for a SVM:

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



- 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)

What does it take to make our cost function small?

For positive examples, $y=1$, the cost is zero when $z \geq 1$.

For negative examples, $y=0$, the cost is zero when $z \leq -1$.

Because of this 'stricter' need SVM have a larger 'safety margin'.

...

What happens if we set our constant C to a very large value if we try to minimize the cost function?

SVM Decision Boundary

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

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

$$\underline{\theta^T x^{(i)}} \geq 1 \quad \min_{\theta} \cancel{C \sum_{i=1}^m} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

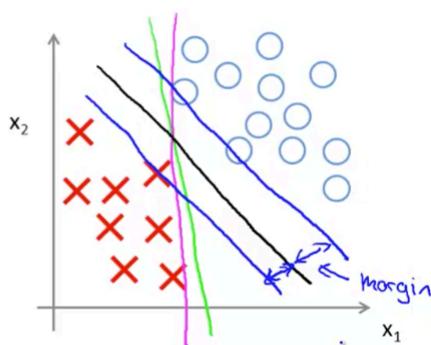
$$\underline{\theta^T x^{(i)}} \leq -1 \quad \text{s.t. } \underline{\theta^T x^{(i)} \geq 1} \text{ if } y^{(i)} = 1$$

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

As you can see the middle term needs to go to zero to be minimize.

We then see a very interesting decision boundary:

SVM Decision Boundary: Linearly separable case



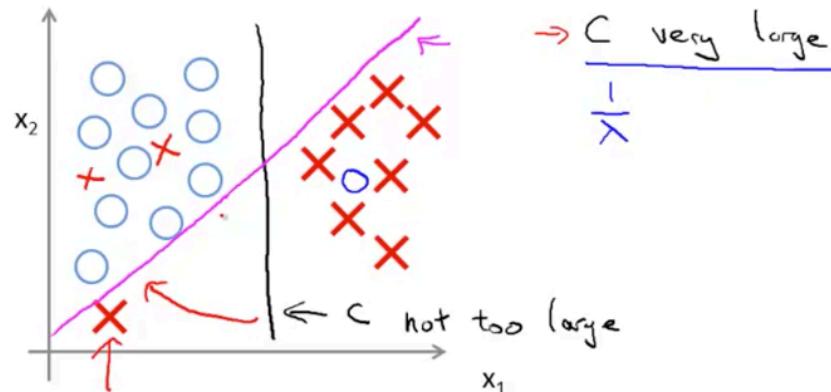
Large margin classifier

The SVM tries to separate the data with as much margin as possible.

Using Large Margin Classifiers can leave your learning algorithm vulnerable to outliers.

If C were very large we would see a decision boundary like the magenta line, but if left C not too large we would see something more like the black line even if your data isn't linearly separable like the graph above:

Large margin classifier in presence of outliers



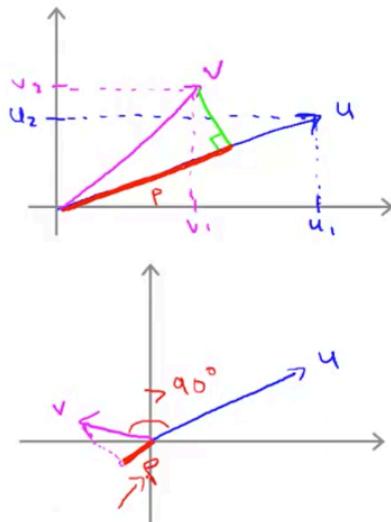
In the next video we will explore the steps to understand how we go from the optimization function to a large margin classifier.

Mathematics Behind Large Margin Classification

This should provide better intuition about how the optimization of SVM leads to large classifiers.

Recall the Vector Inner Product:

Vector Inner Product



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

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

$$\|u\| = \text{length of vector } u$$

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

$$p = \text{length of projection of } v \text{ onto } u.$$

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

Signed

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

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

$$p < 0$$

Andrew'

Now, let us look at the optimization objective of the SVM:

SVM Decision Boundary

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

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

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

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

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

$$\begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix}$$

So essentially what the optimization object is doing is minimizing the squared norm of the squared length of the parameter vector theta.

Now lets look at ($\theta^T x$)

SVM Decision Boundary

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

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

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

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

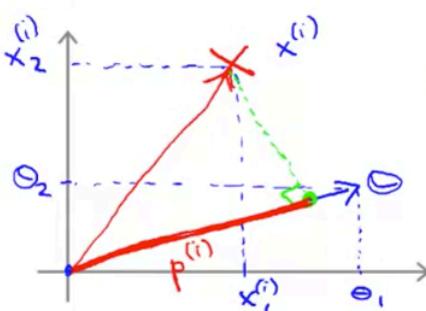
Simplification: $\Theta_0 = \theta$. n=2

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

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

$\uparrow \Theta^T \downarrow$

$\uparrow \Theta^T \downarrow$



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

$$= \Theta_1 x_1^{(i)} + \Theta_2 x_2^{(i)} \leftarrow$$

Andrew N

As such we can write ($\theta_0 + \theta^T x$) into our optimization object we get the following:

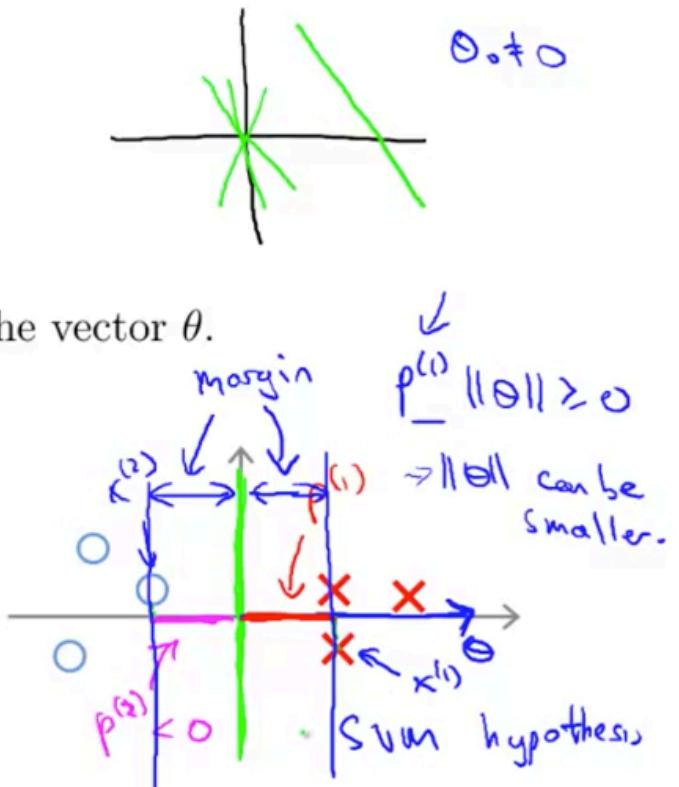
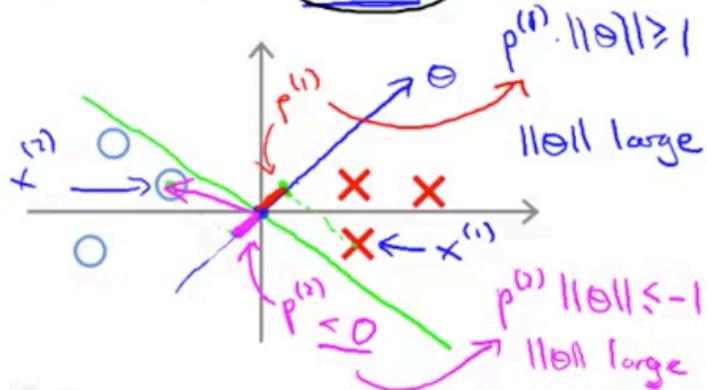
SVM Decision Boundary

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

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: $\theta_0 = 0$ ←



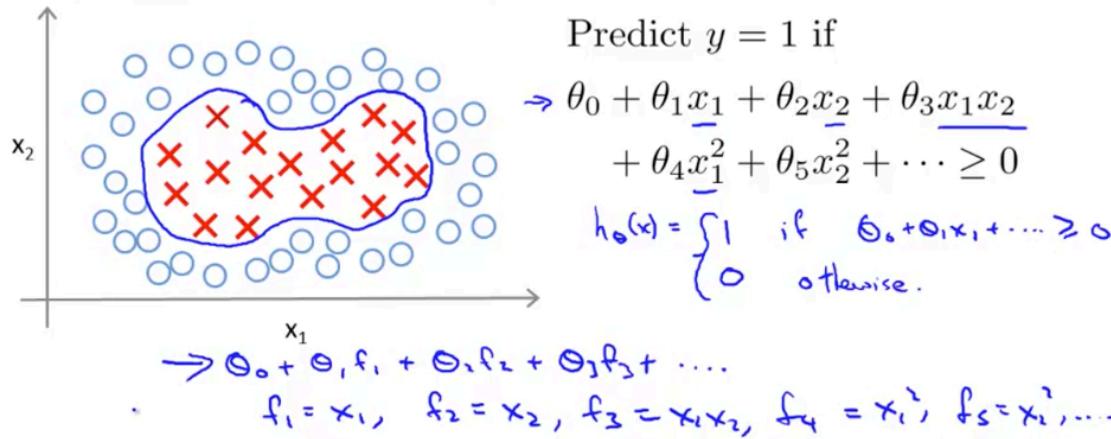
Kernels

Kernels 1

In this video we will start adapting SVMs in order to develop complex non-linear classifiers.

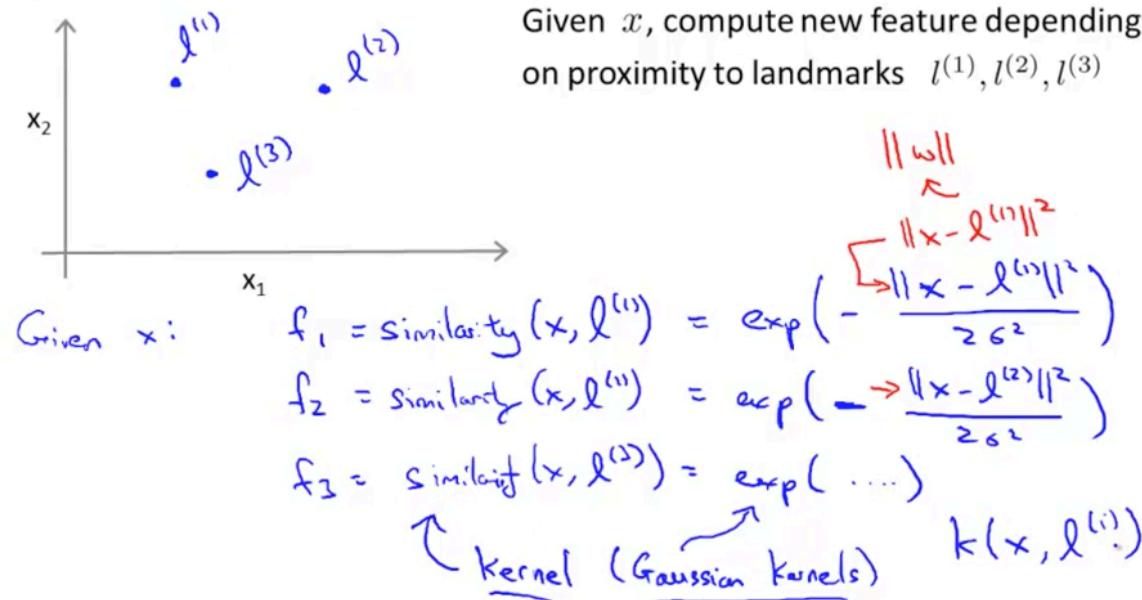
The main technique for doing this is called kernels.

Non-linear Decision Boundary



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

Kernel



The 'similarity' function is what we call the kernel. This particular kernel is called a gaussian kernel.

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

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

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

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

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

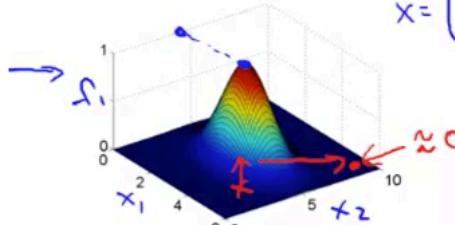
Lets take a better look at the similarity function:

Example:

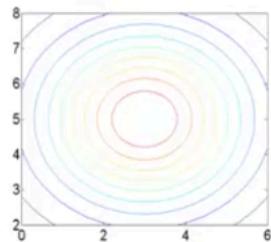
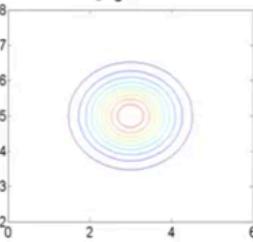
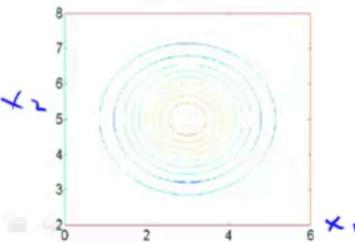
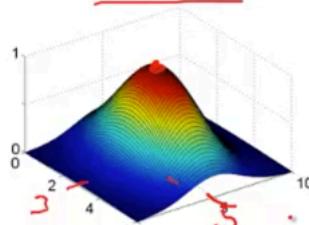
$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x-l^{(1)}\|^2}{2\sigma^2}\right)$$

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

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

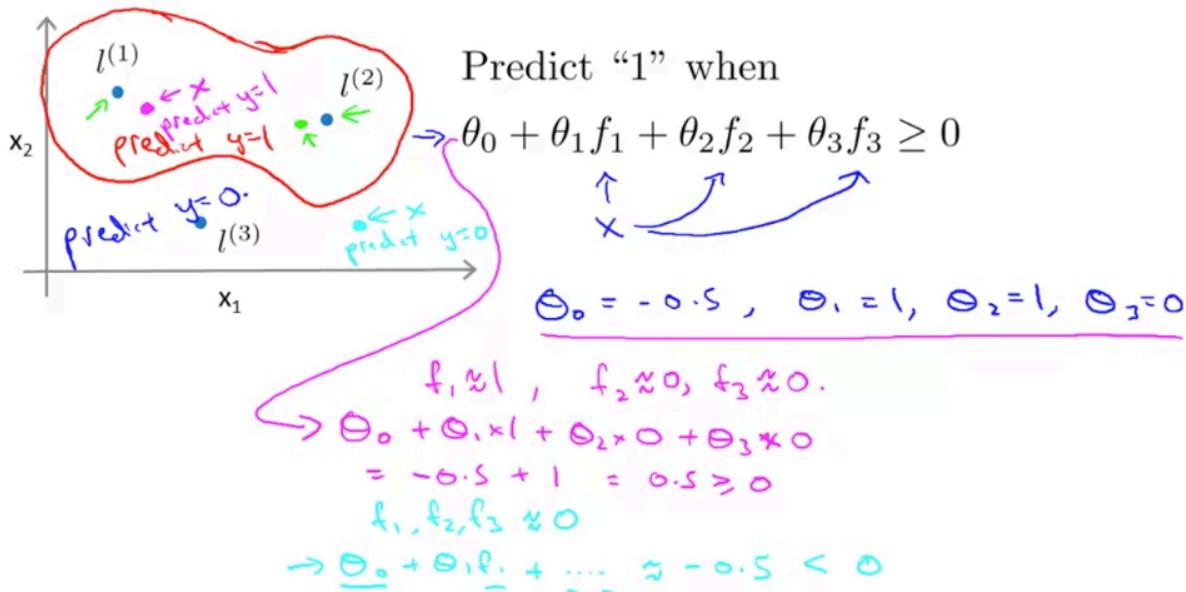


$$\sigma^2 = 3$$



Given the following features, lets say we have f and see what the hypothesis will be. Say we get the following parameters:

What happens if we have a training example at the magenta dot, cyan dot, or green dot?

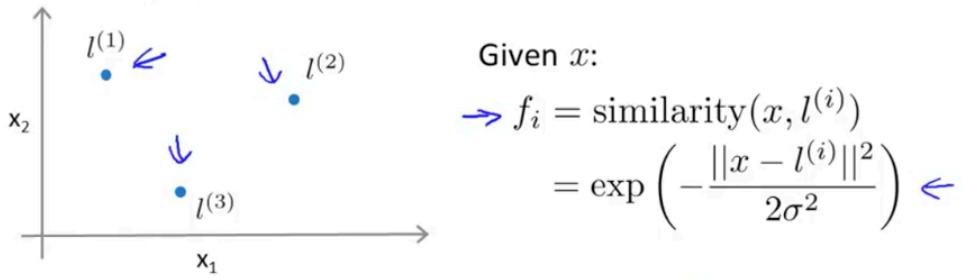


Next we will investigate how do we choose these landmarks and other similarity functions we can use.

Kernels 2

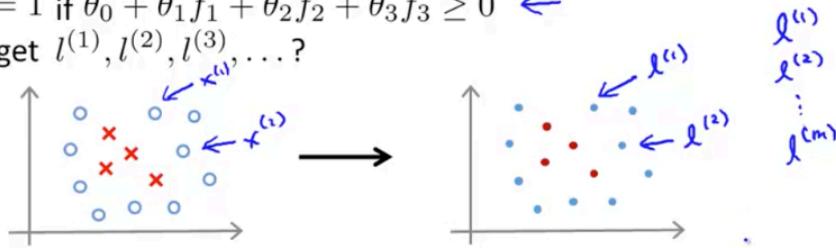
Let us explore how we choose the landmarks and other similarity functions as well as how to use these ideas in practice.

Choosing the landmarks



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

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



What we do is create a landmark in the exact location of every training example we have.

f == feature vector

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 :

$$\rightarrow \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} = \text{similarity}(x, l^{(1)})$$

$$\rightarrow \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} = \text{similarity}(x, l^{(2)})$$

$$\dots$$

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

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

$$\begin{aligned} x^{(i)} \rightarrow & \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \text{similarity}(x^{(i)}, l^{(1)}) \\ & \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \text{similarity}(x^{(i)}, l^{(2)}) \\ & \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \text{similarity}(x^{(i)}, l^{(i)}) = \exp\left(-\frac{\|x^{(i)} - l^{(i)}\|^2}{2\sigma^2}\right) = 1 \end{aligned}$$

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

Andrew Ng

If we want to make a prediction, we first compute a feature vector f for the new x :

SVM with Kernels

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

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

Training:

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

~~$\theta^{(i)}$~~ $\theta^T f^{(i)}$ $\rightarrow \theta_0$

- $\sum_j \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}$ (ignoring θ_0)

- $\theta^T M \theta \leftarrow \| \theta \|^2$ $M = I_{10,000}$

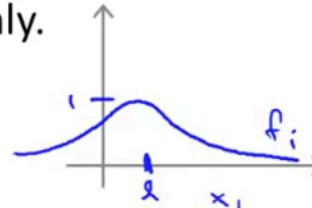
Let's look at the bias and variance trade offs when using a SVM:

SVM parameters:

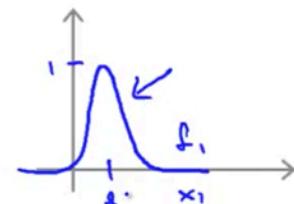
C ($= \frac{1}{\lambda}$). \rightarrow Large C : Lower bias, high variance. (small λ)
 \rightarrow Small C : Higher bias, low variance. (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly.
 \rightarrow Higher bias, lower variance.

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



Small σ^2 : Features f_i vary less smoothly.
Lower bias, higher variance.



Using an SVM

It is not a good idea to write your own software to compute the SVM optimization problem, or any other linear algebra functions. You should use one of the highly optimized libraries for these tasks.

Some SVM software packages are liblinear, libsvm, ... to solve for parameters theta.

The tasks we are responsible for is:

- Choice of parameter C
- Choice of kernel (similarity function)
- Choice of sigma

No kernel (linear kernel) is a valid option.

E.g. No kernel ("linear kernel")
Predict "y = 1" if $\theta^T x \geq 0$ $\Theta_0 + \Theta_1 x_1 + \dots + \Theta_n x_n \geq 0$ $x \in \mathbb{R}^{n+1}$
 $\rightarrow n$ large, m small

Gaussian kernel:

$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$, where $l^{(i)} = x^{(i)}$. $x \in \mathbb{R}^n$, n small
Need to choose σ^2 . and/or m large



If you do decide to use a gaussian kernel, here is what you need to do:

Kernel (similarity) functions:
function $f = \text{kernel}(x_1, x_2)$
 $f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$ $x \rightarrow f_1, f_2, \dots, f_m$
return

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

$$\begin{aligned} & \boxed{\|x - l\|^2} \quad x \in \mathbb{R}^n \\ & v = x - l \\ & \|v\|^2 = v_1^2 + v_2^2 + \dots + v_n^2 \\ & = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2 \\ & \quad \text{1000 feet}^2 \quad 1-5 \text{ bedrooms} \end{aligned}$$

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).

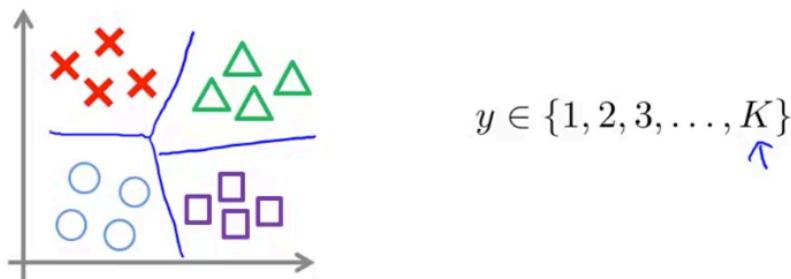
Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, l) = \frac{(x^T l + \text{constant})^{\text{degree}}}{(x^T l)^3}$, $(x^T l + 1)^3$, $(x^T l + 5)^4$
- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...
 $\text{sim}(x, l)$

Usually the polynomial kernel performs worse than the gaussian kernel but it is something you may see.

How do you get an SVM to output decision boundaries for multi-class classification problems?

Multi-class classification



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

When to use Logistic Regression, or SVM?

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. $n \geq m$, $n = 10,000$, $m = 10 \dots 1,000$)
→ Use logistic regression, or SVM without a kernel ("linear kernel")
- If n is small, m is intermediate: ($n = 1-1,000$, $m = 10 - 10,000$)
→ Use SVM with Gaussian kernel
If n is small, m is large: ($n = 1-1,000$, $m = 50,000+$)
→ 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.

** Note SVMs have a convex optimization problem and won't get stuck in local optima

